# Neural Representations of the Cocktail Party in Human Auditory Cortex

Jonathan Z. Simon

*Department of Biology*
*Department of Electrical & Computer Engineering*
*Institute for Systems Research*

University of Maryland

# Acknowledgements

## Grad Students
Francisco Cervantes
Alex Presacco
Krishna Puvvada

## Past Grad Students
Nayef Ahmar
Claudia Bonin
Maria Chait
Marisel Villafane Delgado
Kim Drnec
Nai Ding
Victor Grau-Serrat
Ling Ma
Raul Rodriguez
Juanjuan Xiang
Kai Sum Li
Jiachen Zhuo

## Undergraduate Students
Abdulaziz Al-Turki
Nicholas Asendorf
Sonja Bohr
Elizabeth Camenga
Corinne Cameron
Julien Dagenais
Katya Dombrowski
Kevin Hogan
Kevin Kahn
Andrea Shome
Madeleine Varmer
Ben Walsh

## Collaborators' Students
Murat Aytekin
Julian Jenkins
David Klein
Huan Luo

## Past Postdocs
Dan Hertz
Yadong Wang

## Collaborators
Catherine Carr
Monita Chatterjee
Alain de Cheveigné
Didier Depireux
Mounya Elhilali
Jonathan Fritz
Cindy Moss
David Poeppel
Shihab Shamma

# Introduction
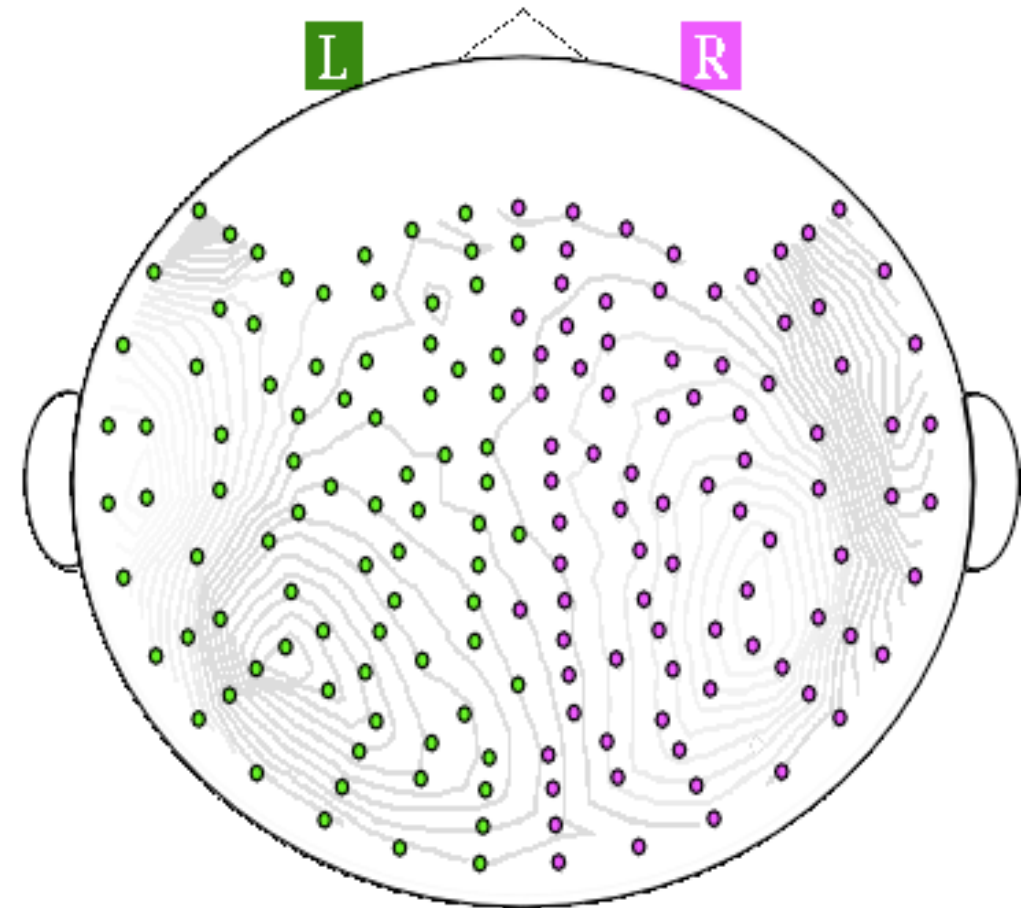
- Magnetoencephalography (MEG)

- Auditory Objects

- Neural Representations of Auditory Objects in Cortex: Decoding

- Neural Representations of Auditory Objects in Cortex: Encoding

- Neural Representations of Speech in Noise
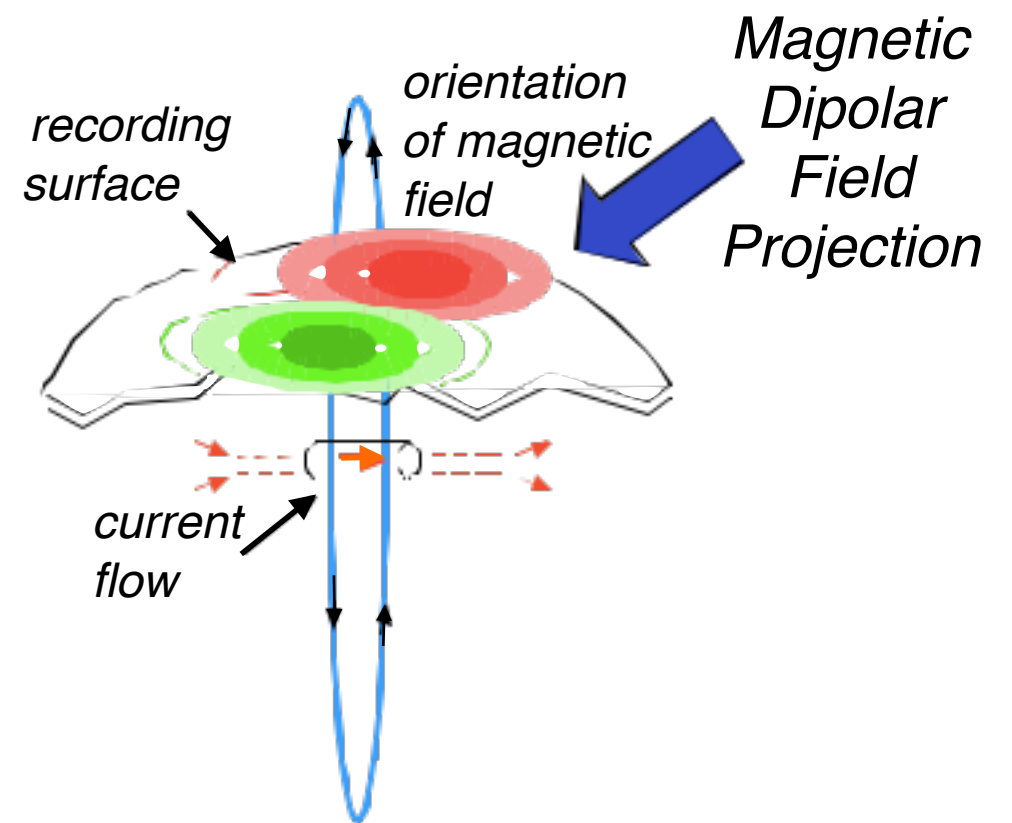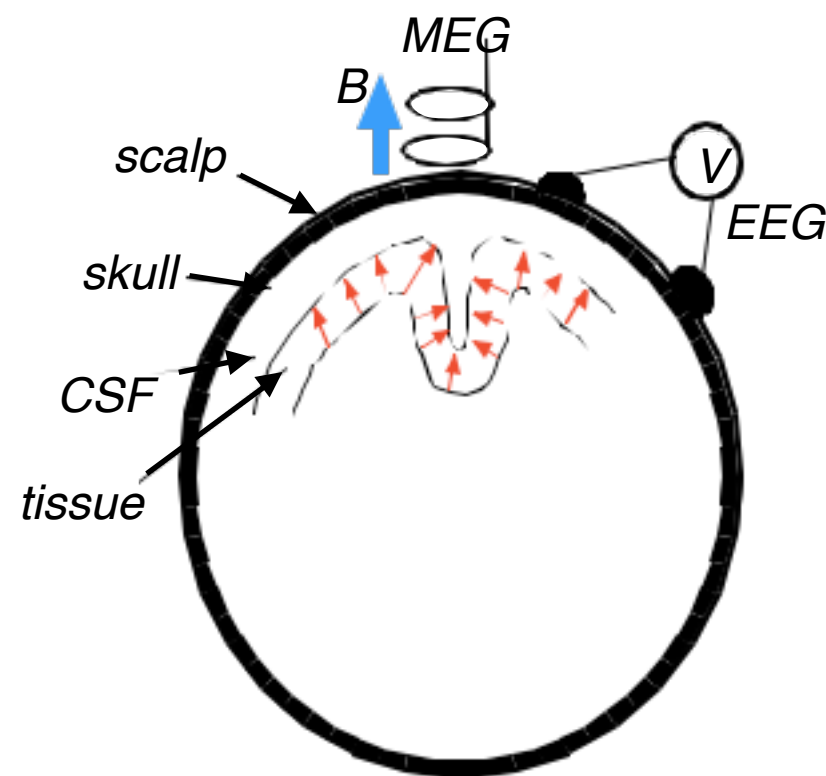
# Magnetoencephalography

- Non-invasive, Passive, Silent Neural Recordings

- Simultaneous Whole-Head Recording (~200 sensors)

- Sensitivity
  - high: ~100 fT ($10^{-13}$ Tesla)
  - low: ~$10^4$ – ~$10^6$ neurons

- Temporal Resolution: ~1 ms

- Spatial Resolution
  - coarse: ~1 cm
  - ambiguous

# Neural Signals & MEG



*Photo by Fritz Goro*

MEG

B

scalp

skull

CSF

tissue

V

EEG

recording surface

orientation of magnetic field
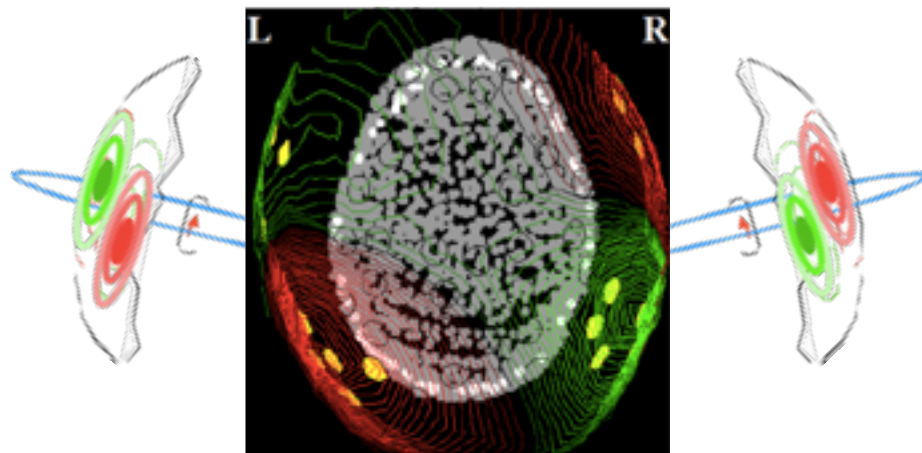
*Magnetic Dipolar Field Projection*

current flow

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

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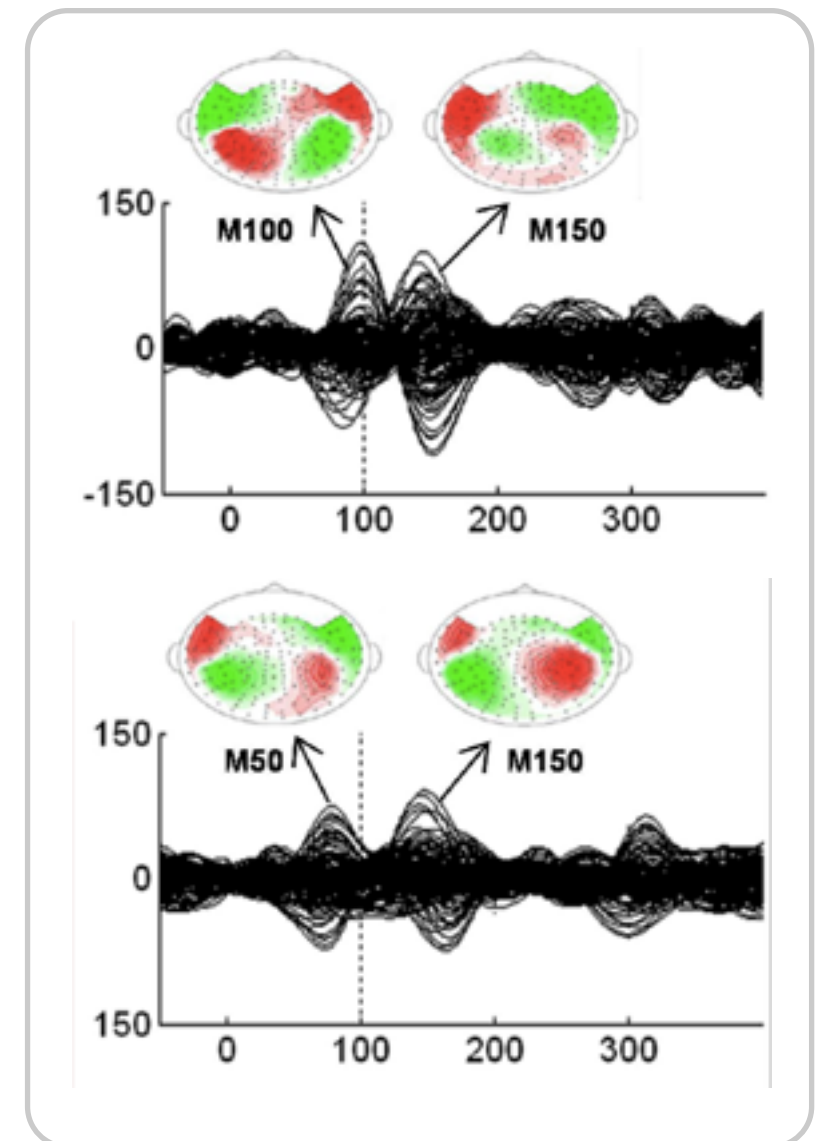- Fine temporal resolution (~ 1 ms)
- Moderate spatial resolution (~ 1 cm)

# Time Course of MEG Responses

**Auditory Evoked Responses**

- MEG Response Patterns Time-Locked to Stimulus Events

- Robust

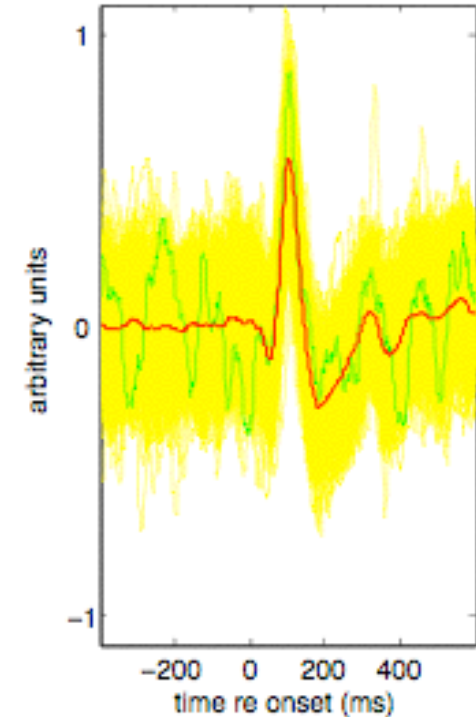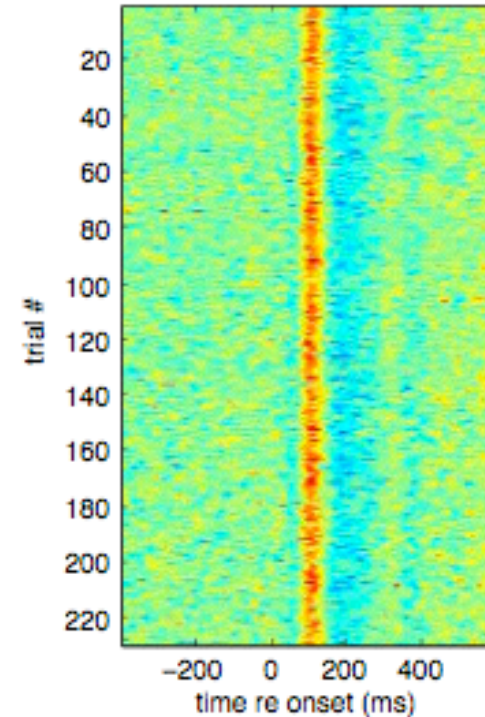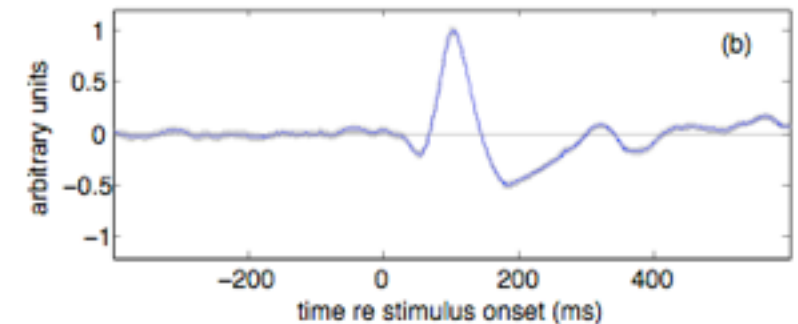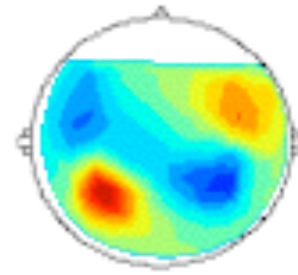- Strongly Lateralized



Pure Tone



Broadband Noise

# MEG Component Analysis

- Data driven spatial filtering:
  many available methods—ICA, PCA, DSS

- Generate spatial filters & their outputs ("components")

- DSS: Denoising Source Separation:
  Särelä & Valpola (2005)

- DSS components ordered by reproducibility
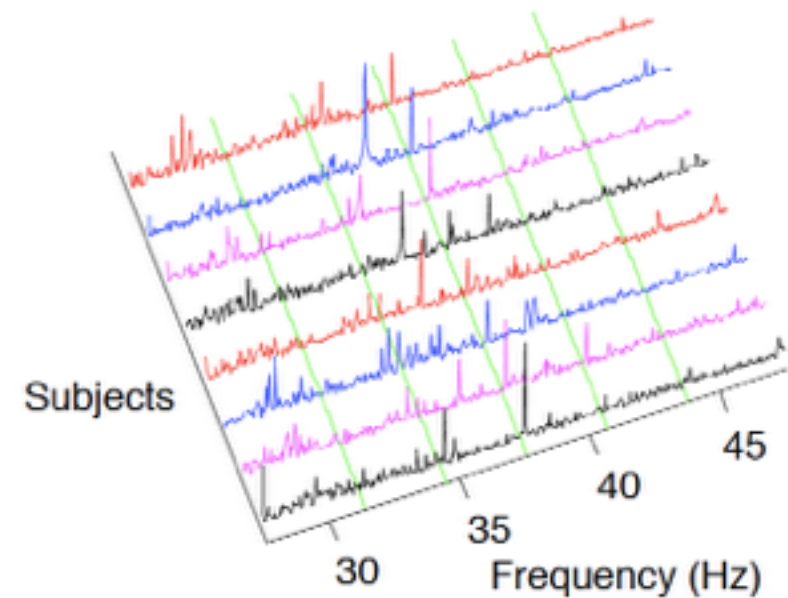  - 1st component "maximally reproducible" = most stimulus driven

# DSS Example

- Most reproducible filter & component

- Optimally filters out trial-to-trial-variable signal = neural noise

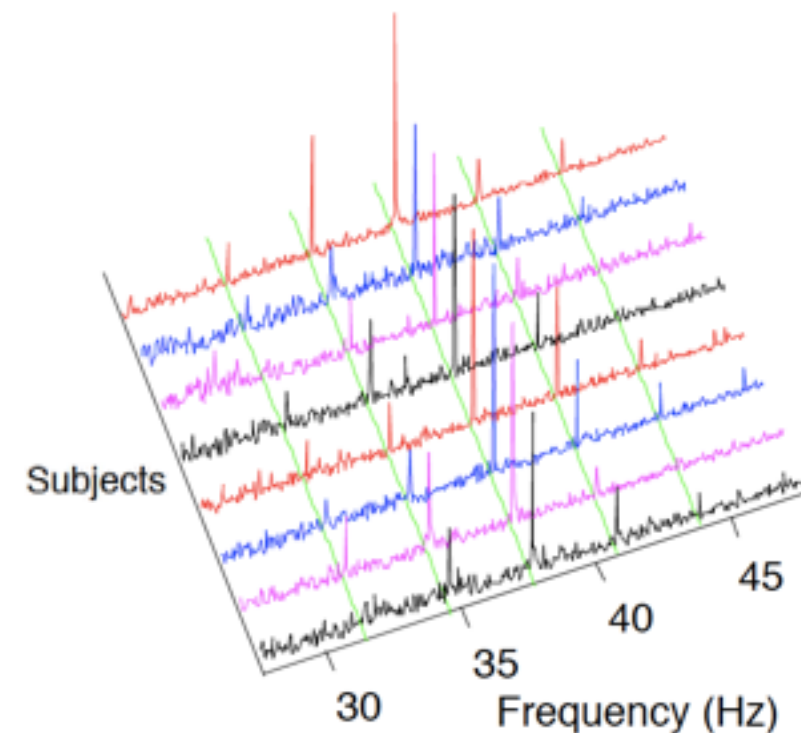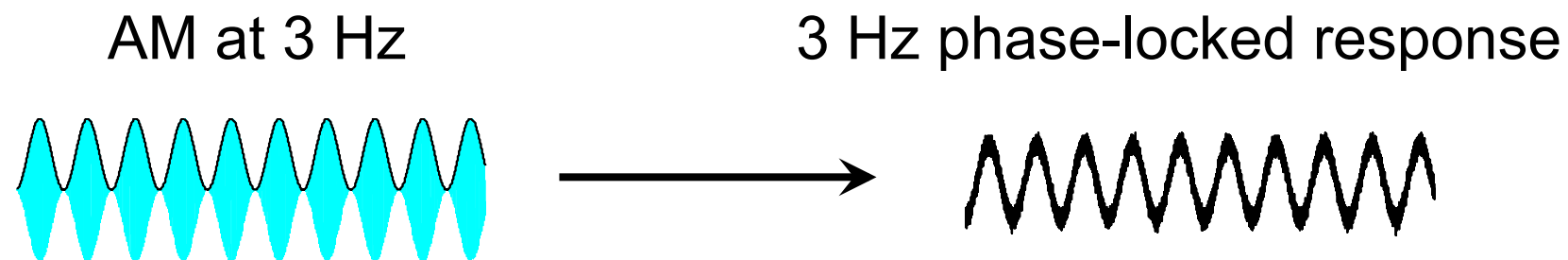- Filter can be applied to other signals, e.g. single trials

Särelä & Valpola (2005)
de Cheveigné & Simon, J. Neurosci. Methods (2008)

# DSS Example: Spectral

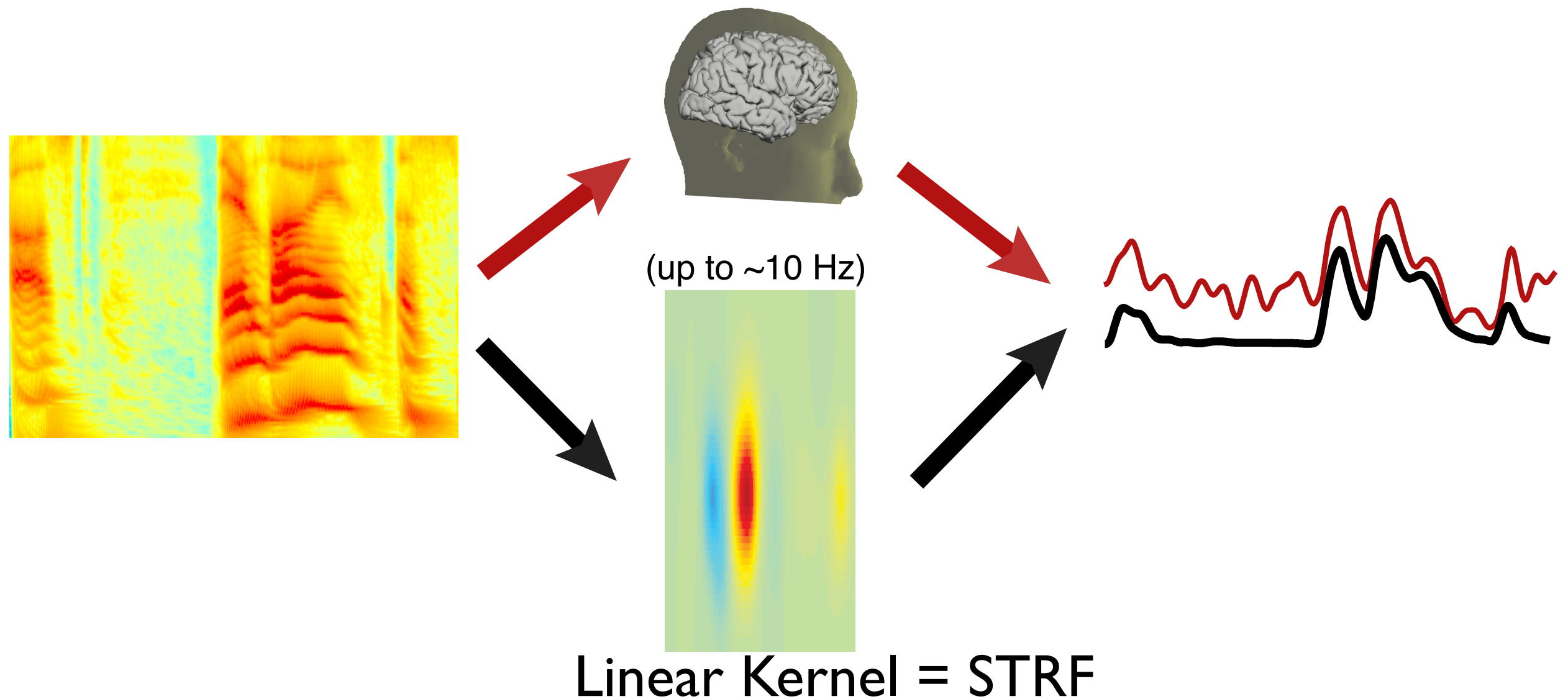Frequency Spectrum before DSS



Frequency Spectrum after DSS



Ding & Simon, J. Neurophysiol (2009)

# Phase-Locking in MEG to Slow Acoustic Modulations

AM at 3 Hz

3 Hz phase-locked response

Ding & Simon, J Neurophysiol (2009)
Wang et al., J Neurophysiol (2012)

# MEG Responses
# Predicted by STRF Model



(up to ~10 Hz)

Linear Kernel = STRF

"Spectro-Temporal Response Function"

# Neural Reconstruction of Speech Envelope
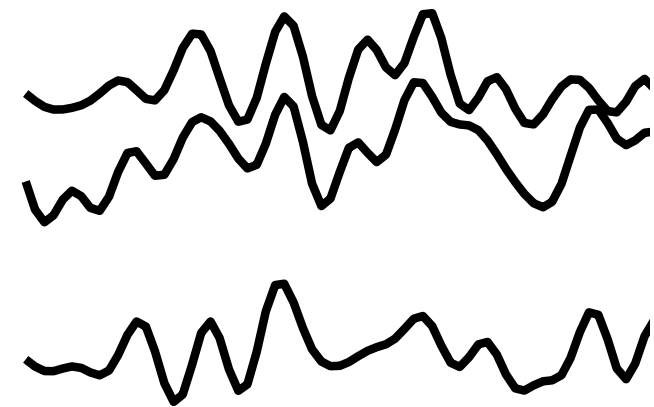
# Auditory Objects

- What is an auditory object?

  - perceptual construct (not neural, not acoustic)

  - commonalities with visual objects

  - several potential formal definitions

# Auditory Object Definition

- Griffiths & Warren definition:

  - corresponds with *something* in the sensory world

  - object information *separate from* information of rest of sensory world

  - abstracted: object information *generalized over particular* sensory experiences

# Auditory Objects at the Cocktail Party
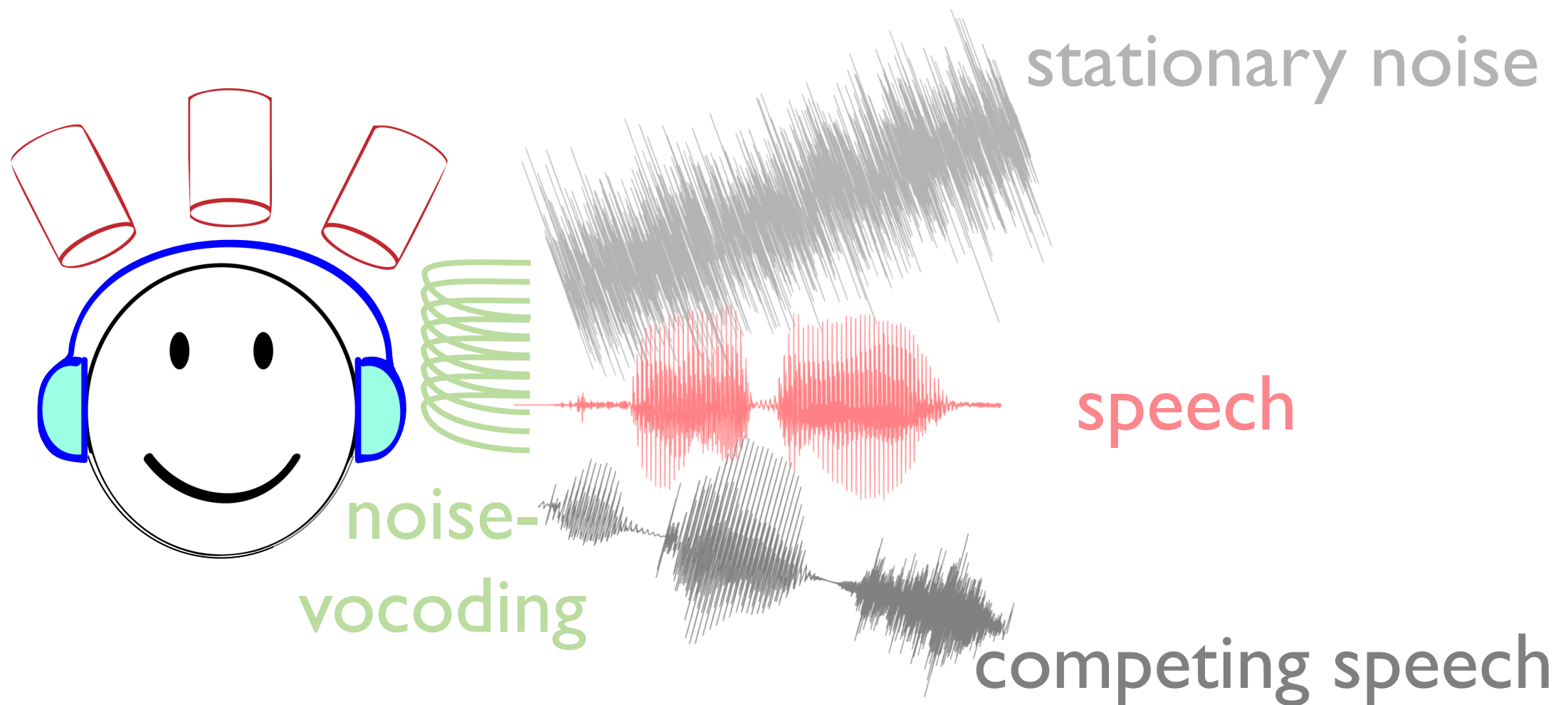


Alex Katz,
The Cocktail Party

# Auditory Objects at the Cocktail Party



Alex Katz,
The Cocktail Party

# Experiments



stationary noise

speech
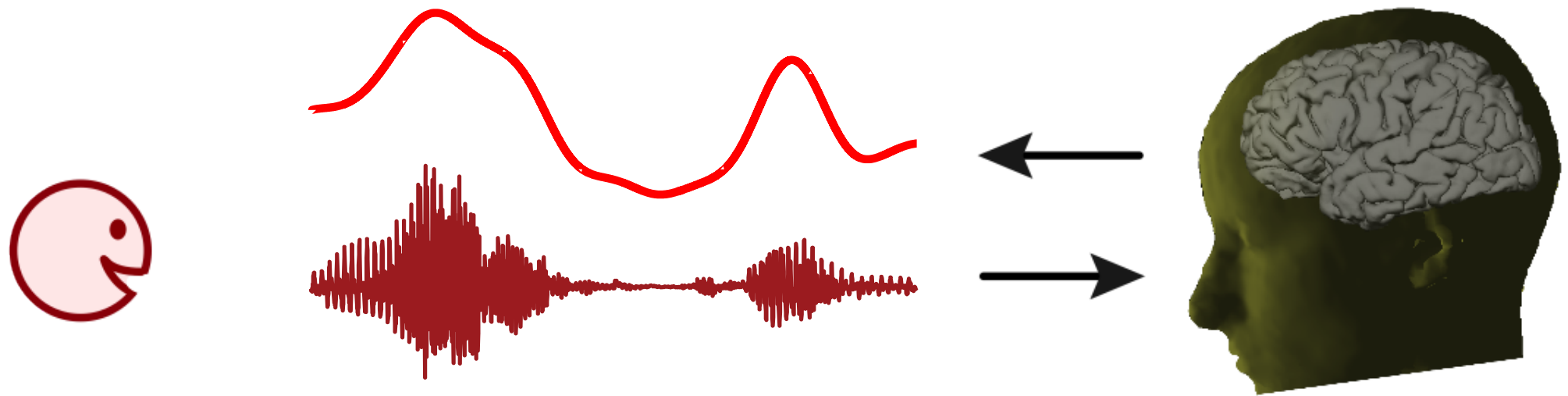
noise-vocoding

competing speech

# Speech Stream as an Auditory Object

- corresponds with something in the sensory world

- information *separate from* information of rest of sensory world
  e.g. other speech streams or noise

- abstracted: object information *generalized over particular* sensory experiences
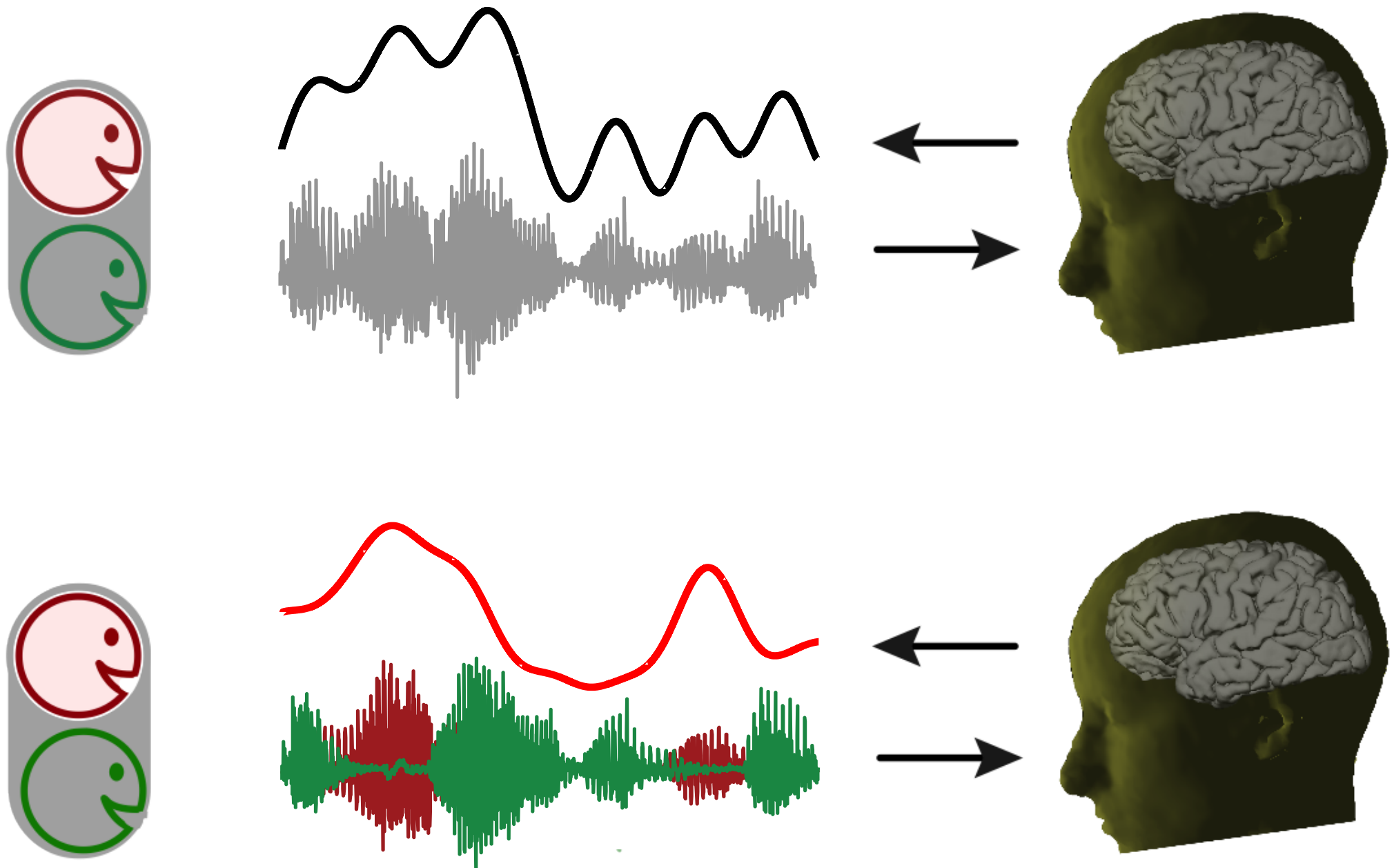  e.g. different sound mixtures

# *Neural Representation* of an Auditory Object

- neural representation is of something in sensory world

- when other sounds mixed in, neural representation is of that auditory object, not entire acoustic scene

- neural representation invariant under broad changes in specific acoustics
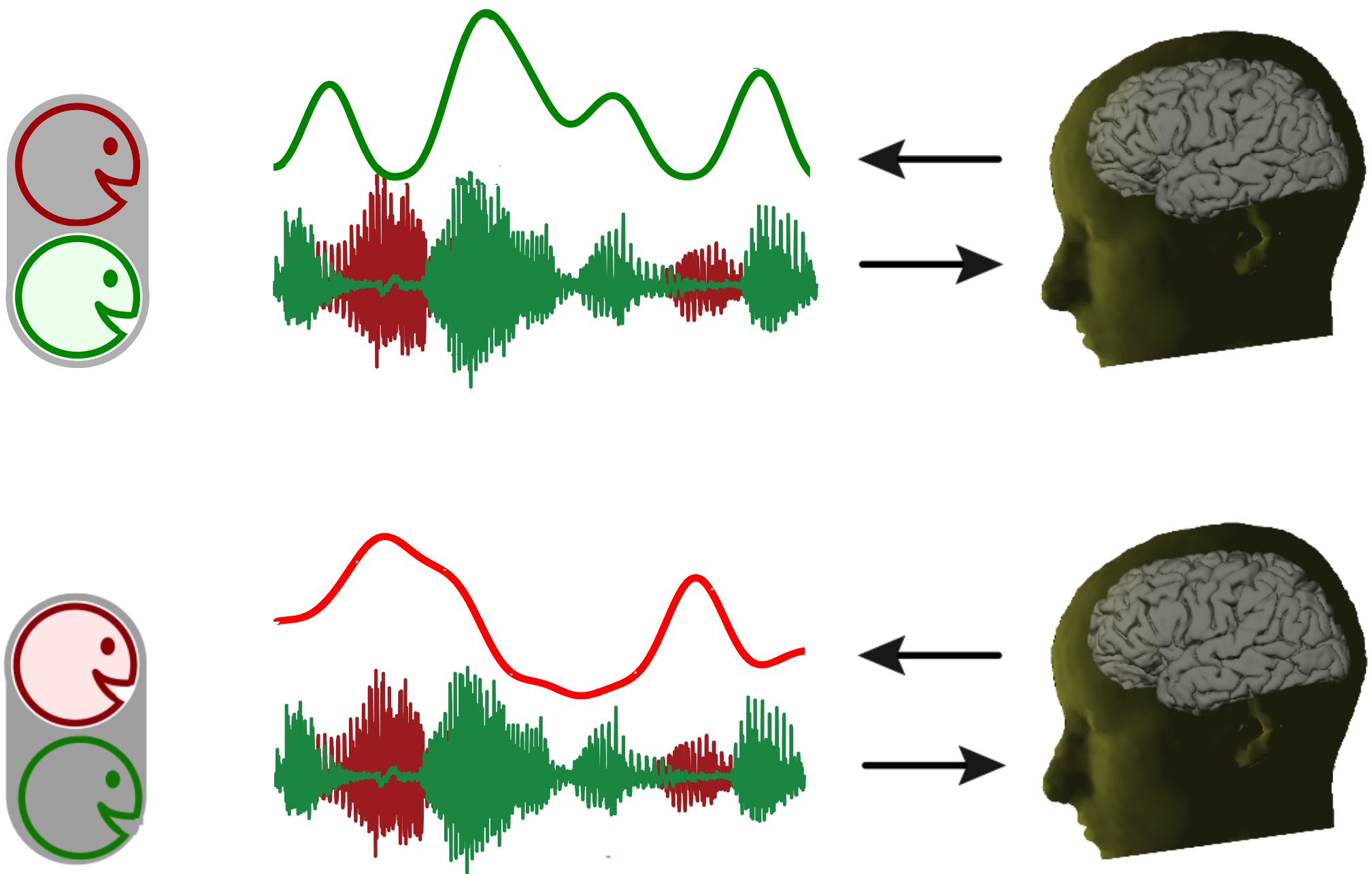
# Selective Neural Encoding

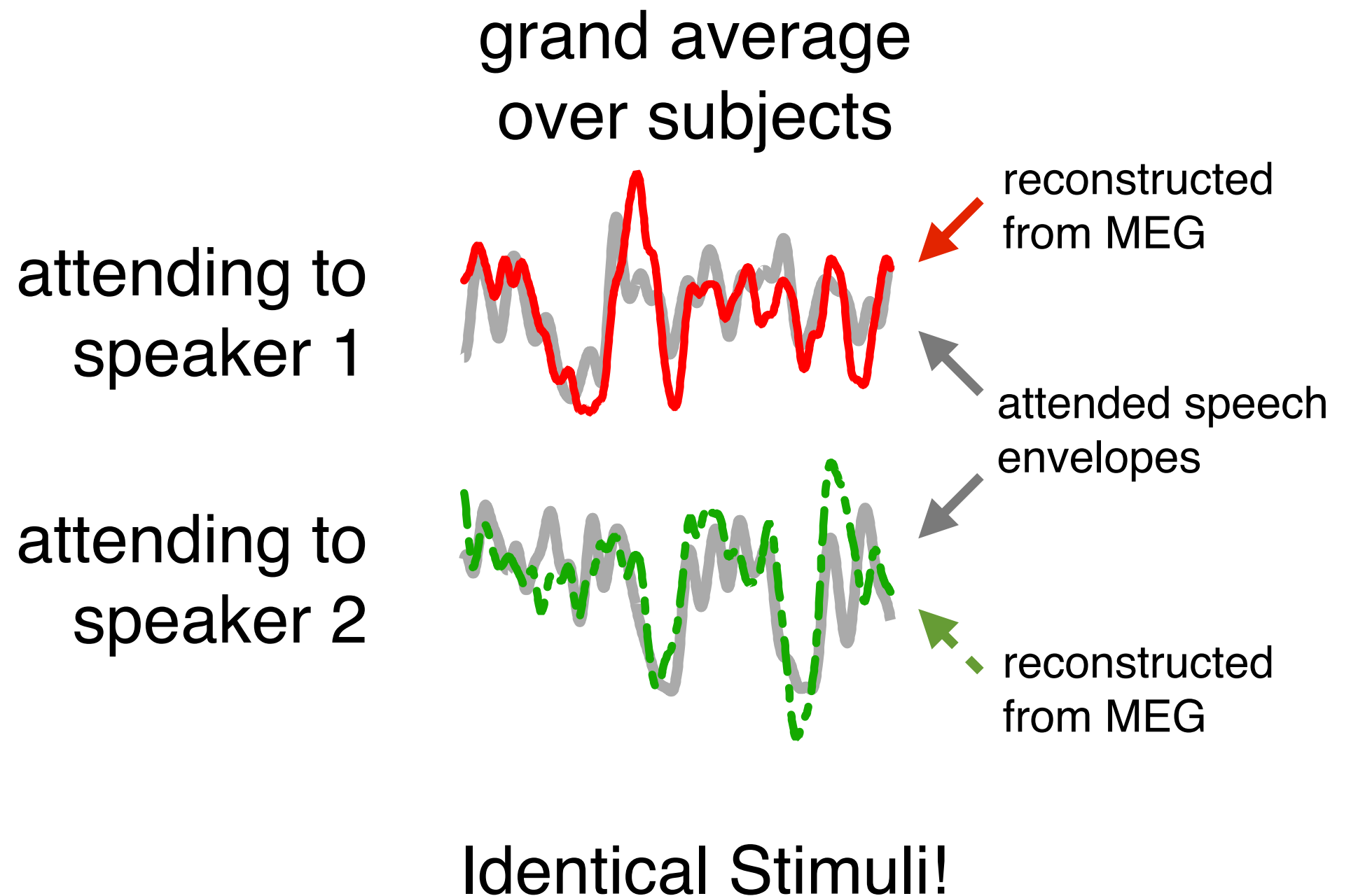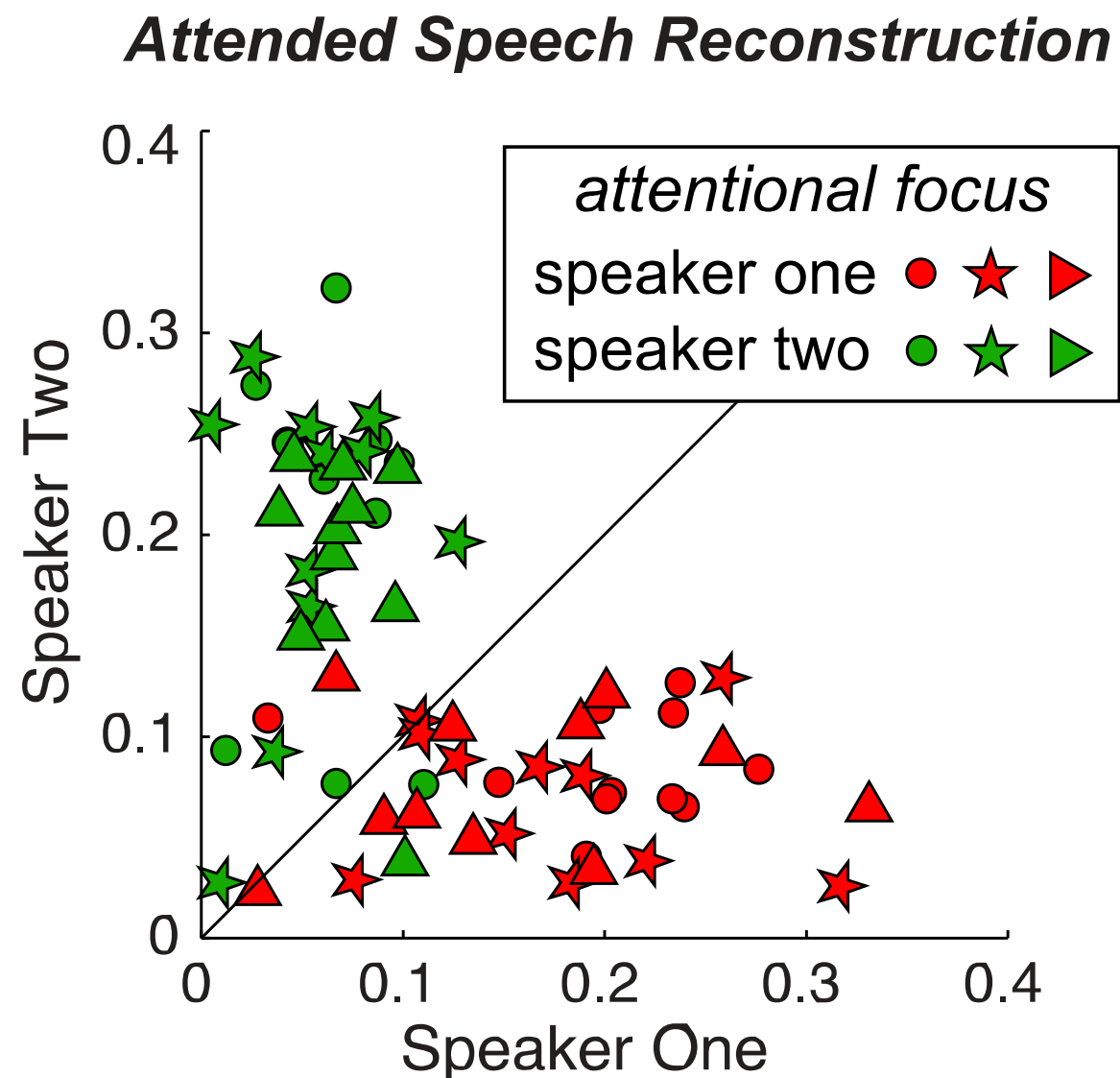# Unselective vs. Selective Neural Encoding

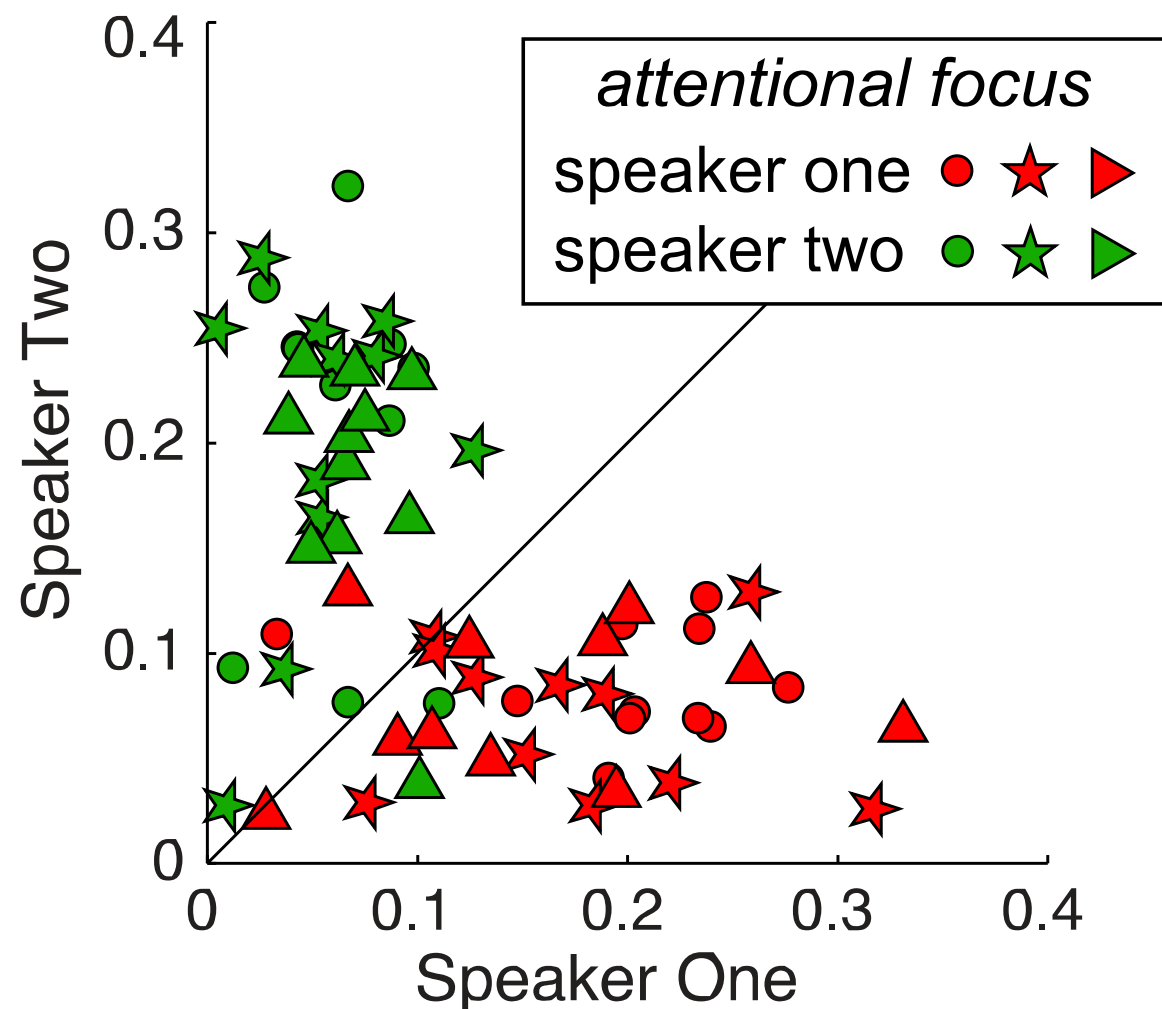# Selective Neural Encoding

# Stream-Specific Representation

grand average
over subjects



attending to speaker 1

reconstructed from MEG

attended speech envelopes

attending to speaker 2

reconstructed from MEG

Identical Stimuli!

# Single Trial Speech Reconstruction
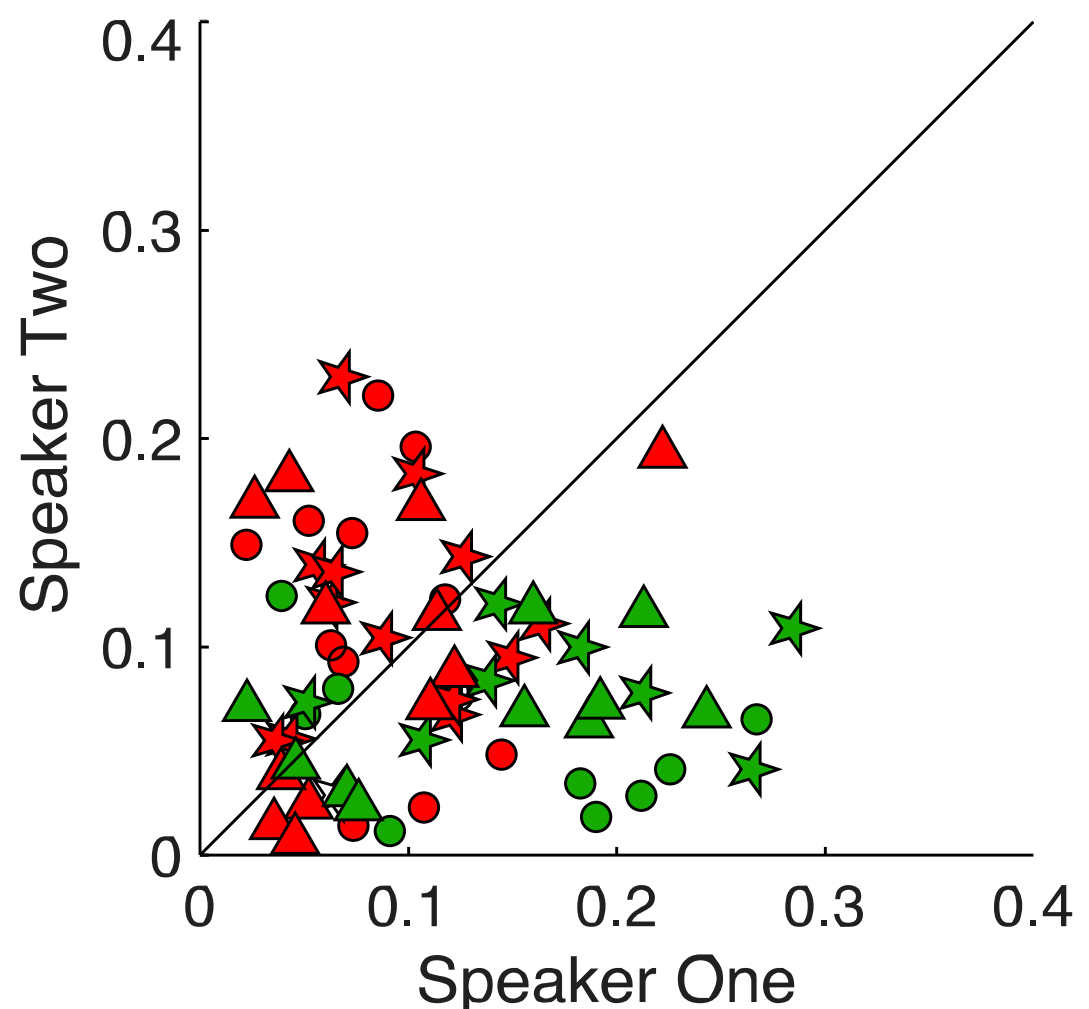


**Attended Speech Reconstruction**

Ding & Simon, PNAS (2013)

# Single Trial Speech Reconstruction



Attended Speech Reconstruction

Background Speech Reconstruction

Ding & Simon, PNAS (2013)

# Invariance Under Acoustic Changes

# Invariance Under Acoustic Changes

# Stream-Based Gain Control?

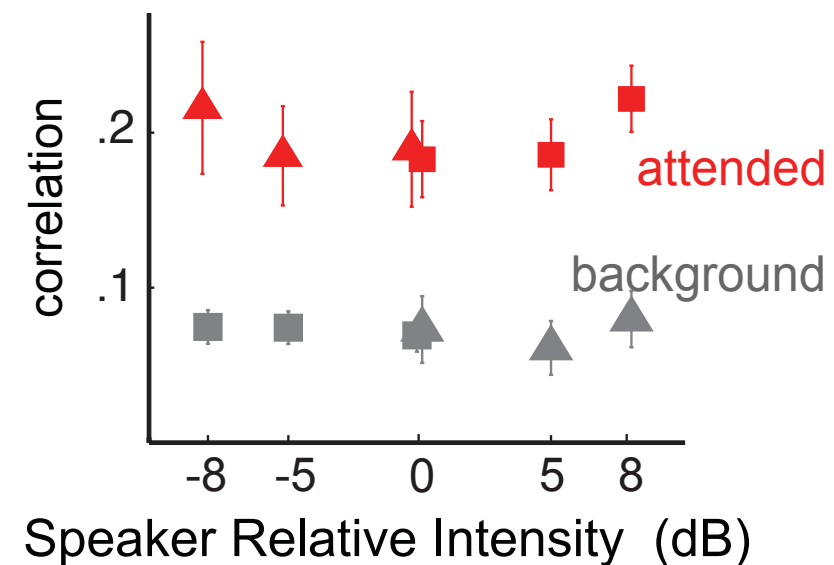## Gain-Control Models

**Object-Based**



correlation
.2 — attended
.1 — background
-8  -5  0  5  8
Speaker Relative Intensity  (dB)

**Stimulus- Based**



correlation
.2 — attended
.1 — background
-8  -5  0  5  8
Speaker Relative Intensity  (dB)

## Neural Results



correlation
.2 — attended
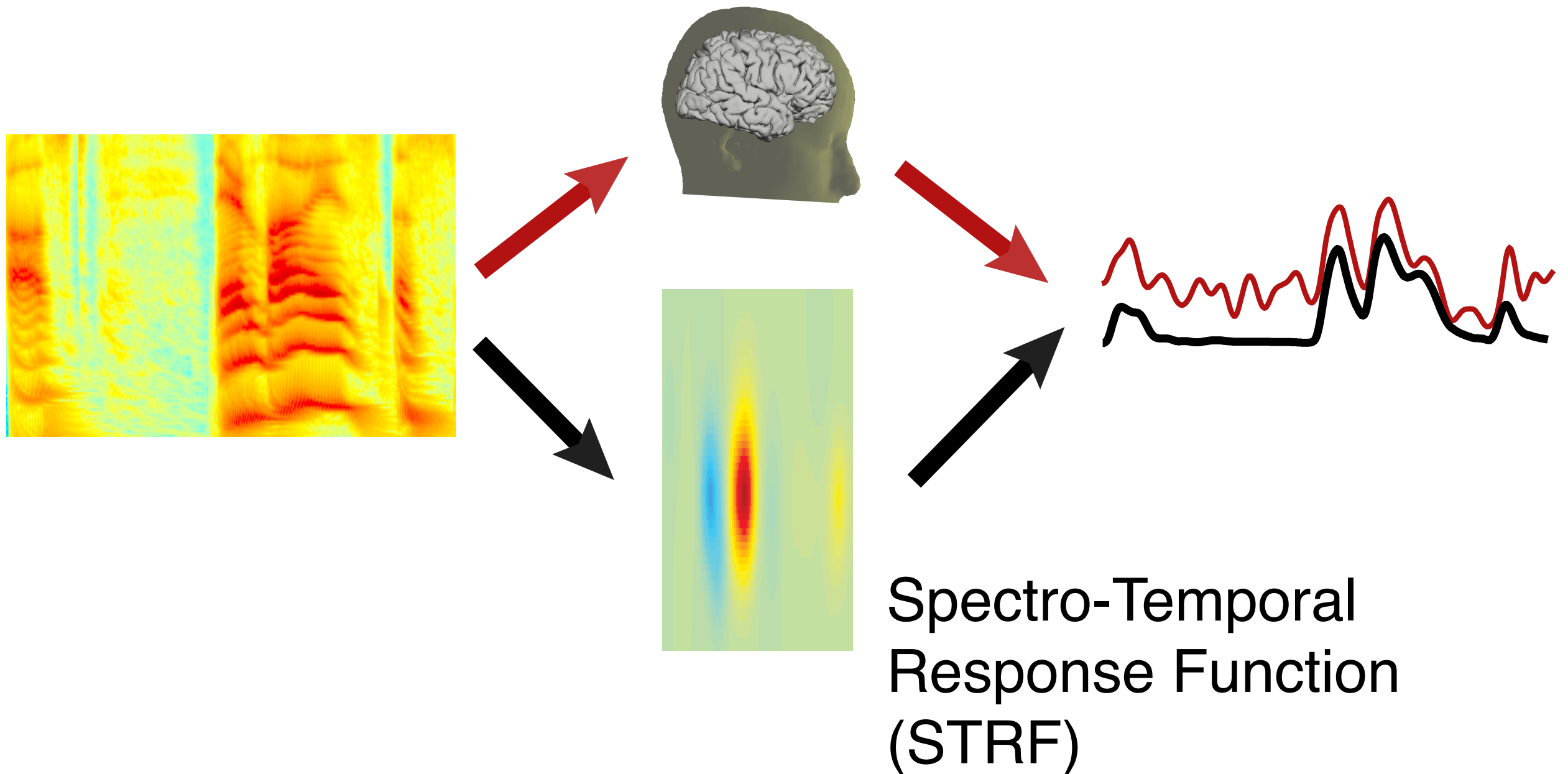.1 — background
-8  -5  0  5  8
Speaker Relative Intensity  (dB)

- Stream-based not stimulus-based
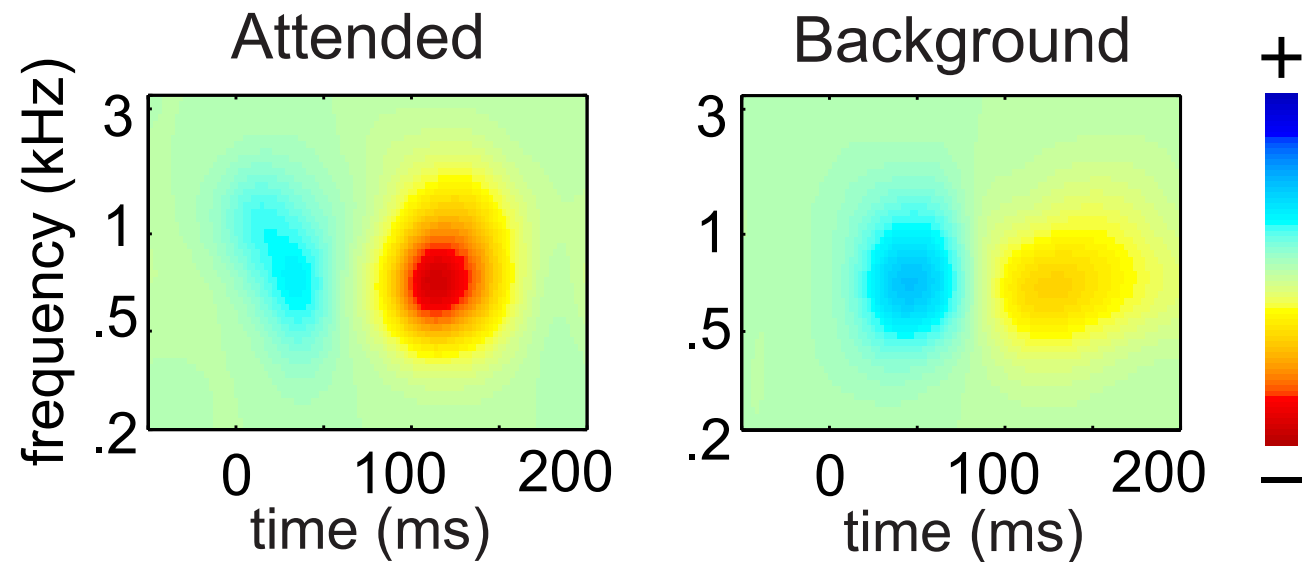- Neural representation is invariant to acoustic changes.

# Neural Representation of an Auditory Object

✓ neural representation is of something in sensory world

✓ when other sounds mixed in,
neural representation is of auditory object,
not entire acoustic scene

✓ neural representation invariant
under broad changes in specific acoustics

# Forward STRF Model



Spectro-Temporal
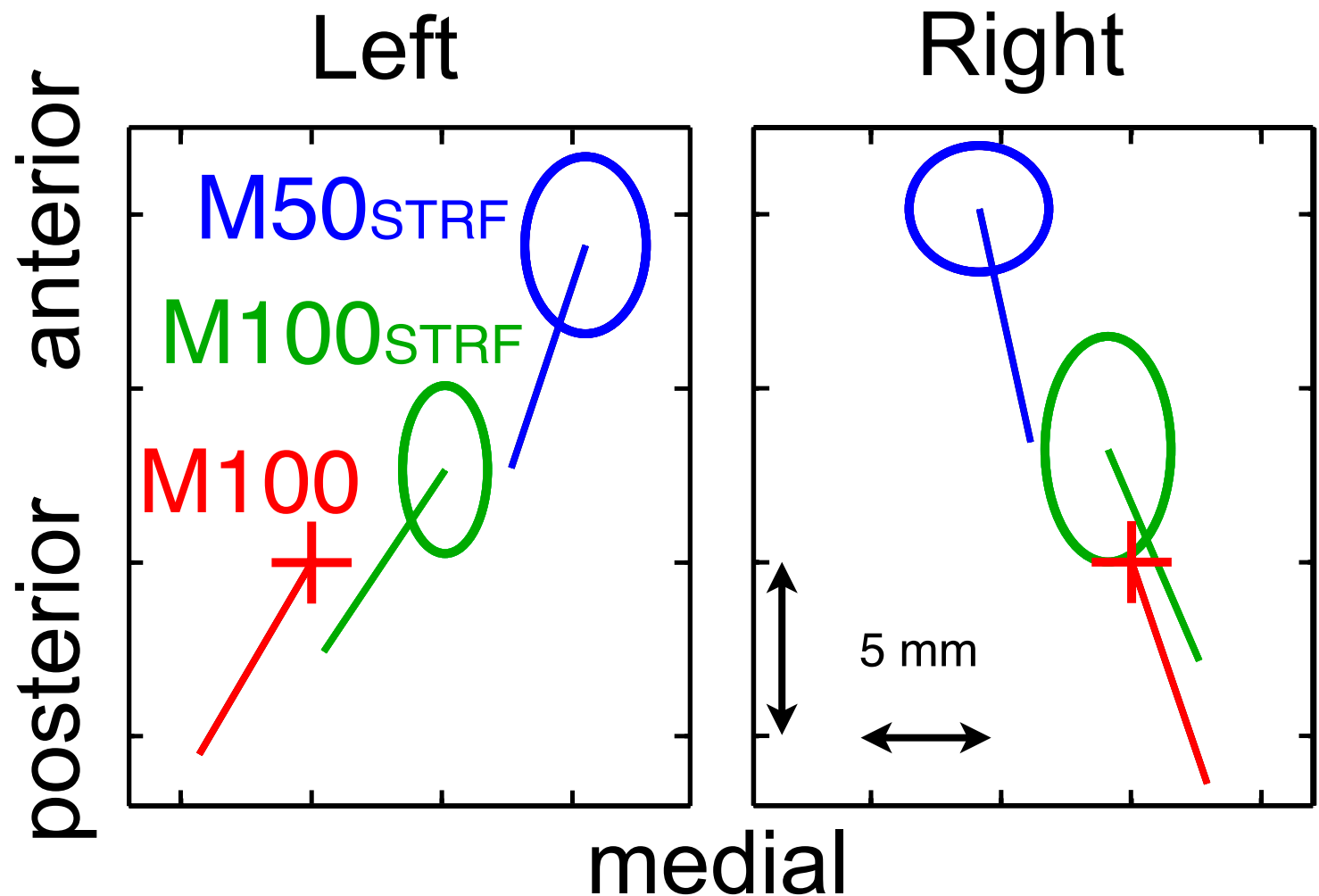Response Function
(STRF)

# STRF Results



- STRF separable (time, frequency)
- 300 Hz - 2 kHz dominant carriers
- M50$_{STRF}$ positive peak
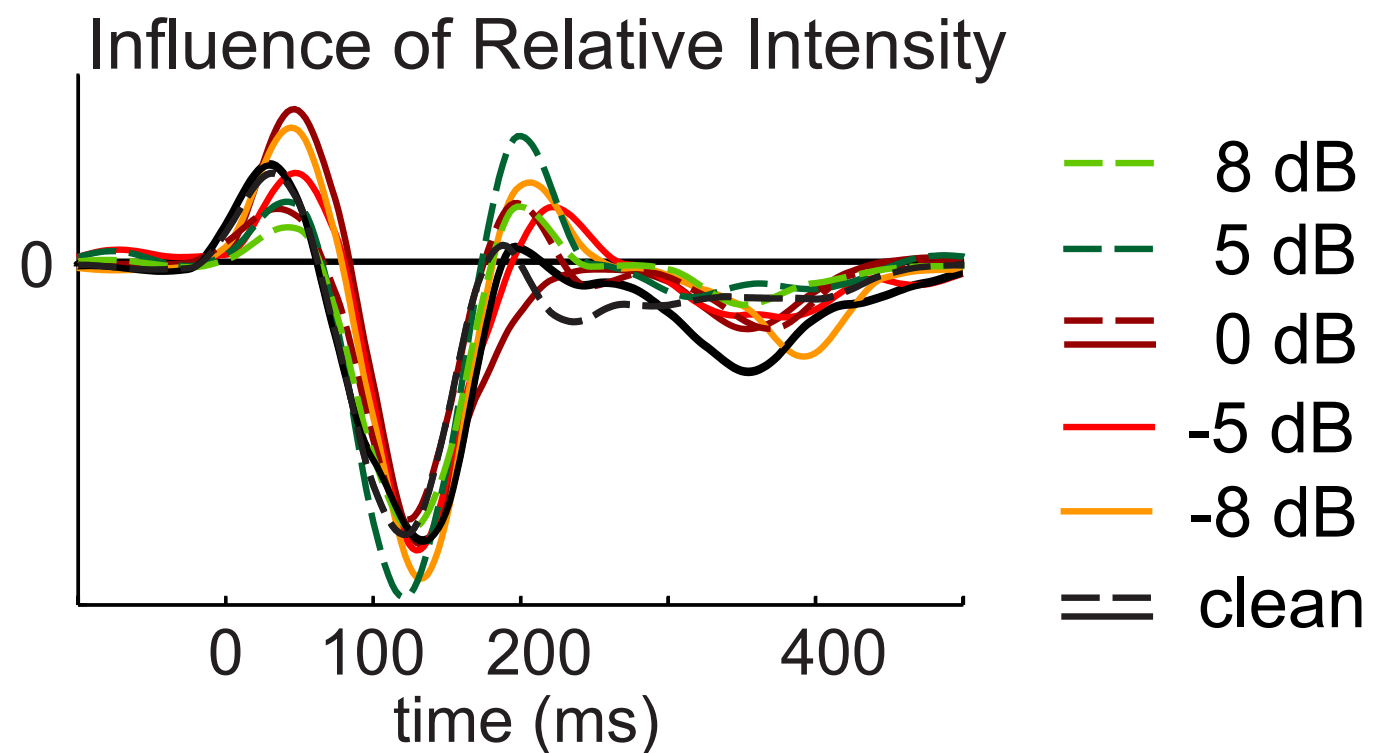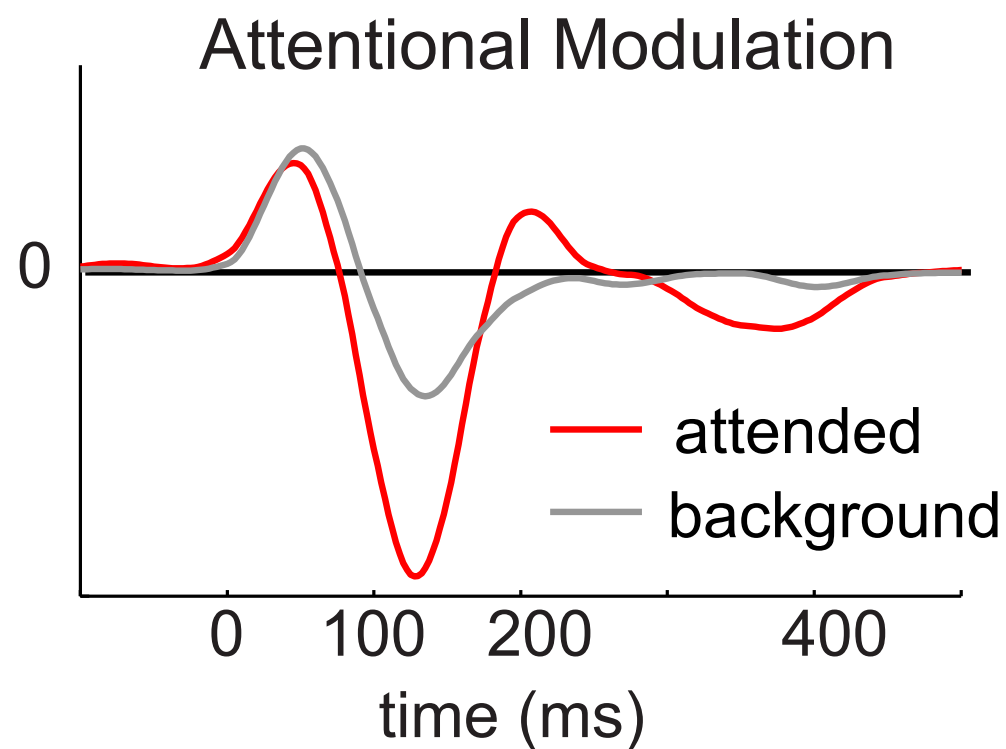- M100$_{STRF}$ negative peak

# Neural Sources

• M100$_{STRF}$ source near (same as?) M100 source:
Planum Temporale

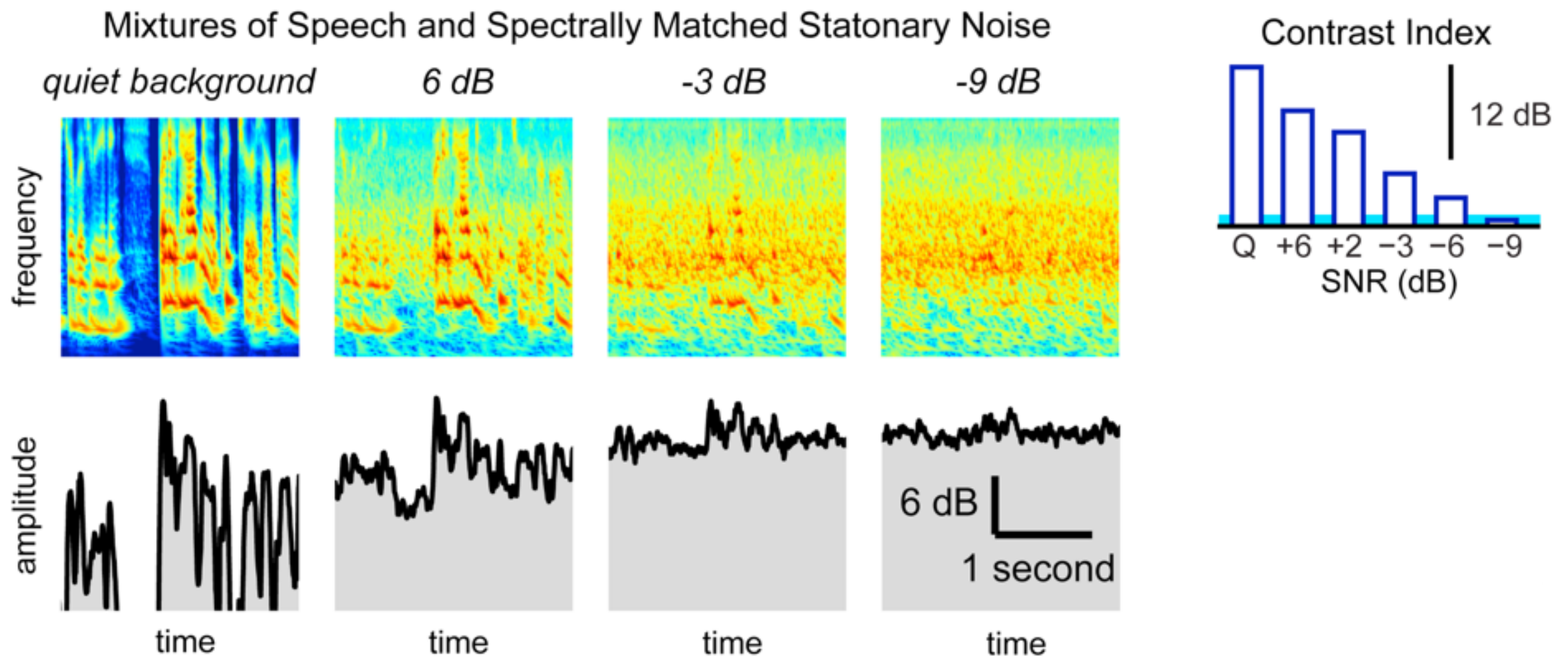• M50$_{STRF}$ source is anterior and medial to M100 (same as M50?):
Heschl's Gyrus

# Cortical Object-Processing Hierarchy



Attentional Modulation

- attended
- background

Influence of Relative Intensity

- 8 dB
- 5 dB
- 0 dB
- -5 dB
- -8 dB
- clean

time (ms)

- M100$_{STRF}$ strongly modulated by attention, but not M50$_{STRF}$.
- M100$_{STRF}$ invariant against acoustic changes.
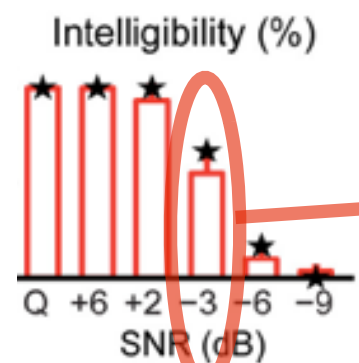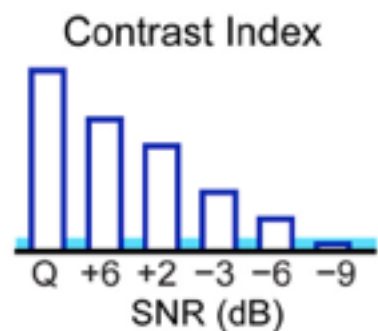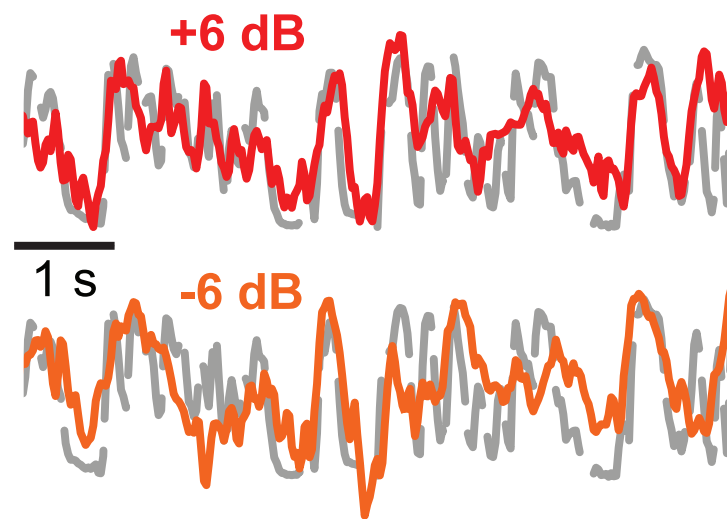- Objects well-neurally represented at 100 ms, but not 50 ms.
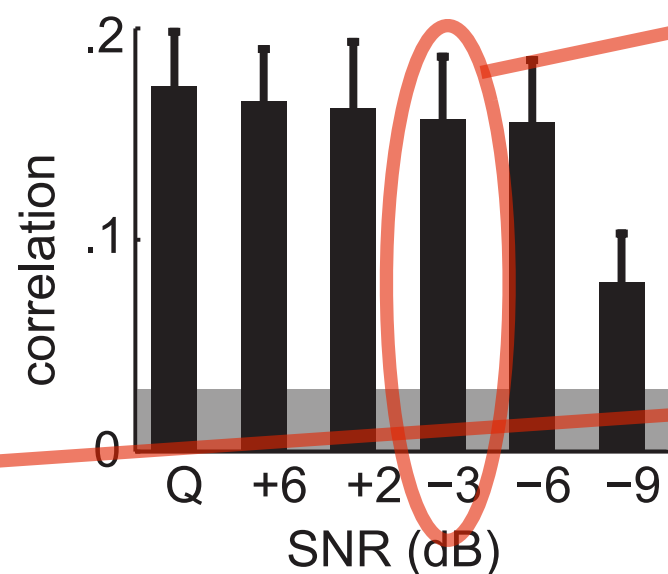
# Speech in Noise



Mixtures of Speech and Spectrally Matched Statonary Noise

quiet background     6 dB     -3 dB     -9 dB

Contrast Index

12 dB

Q  +6  +2  −3  −6  −9
SNR (dB)

# Speech in Noise: Results



Neural Reconstruction of Underlying Speech Envelope

+6 dB

1 s

-6 dB

Contrast Index

Reconstruction Accuracy

Correlation with Intelligiblity

Intelligibility (%)

# Noise-Vocoded Speech



natural      8-band      4-band

in quiet: 100±0%   93±2%   43±6%

in noise: 99±1%   34±6%   6±2%

frequency (kHz): 4, .6, .1

2 seconds

"in noise" = +3 dB SNR
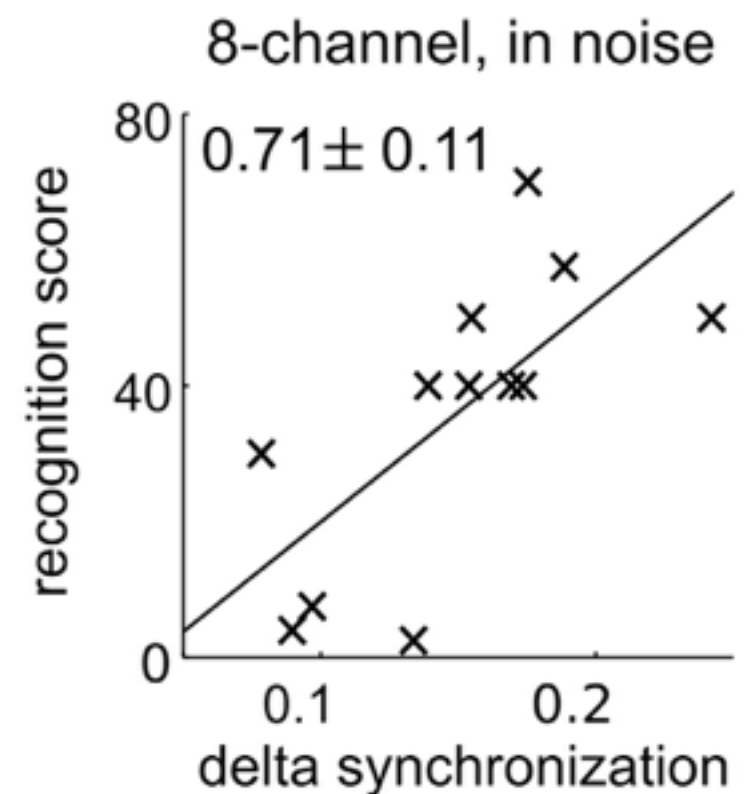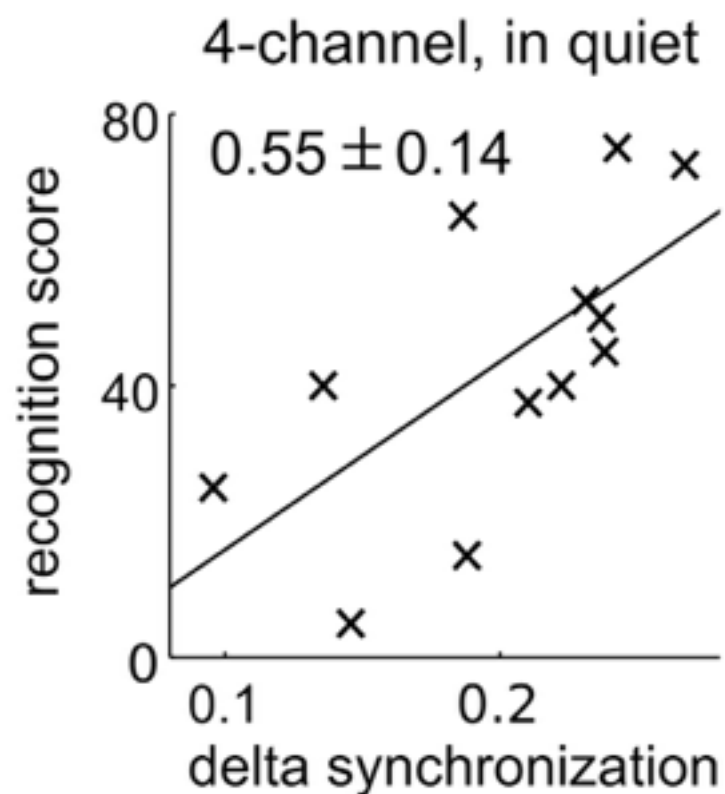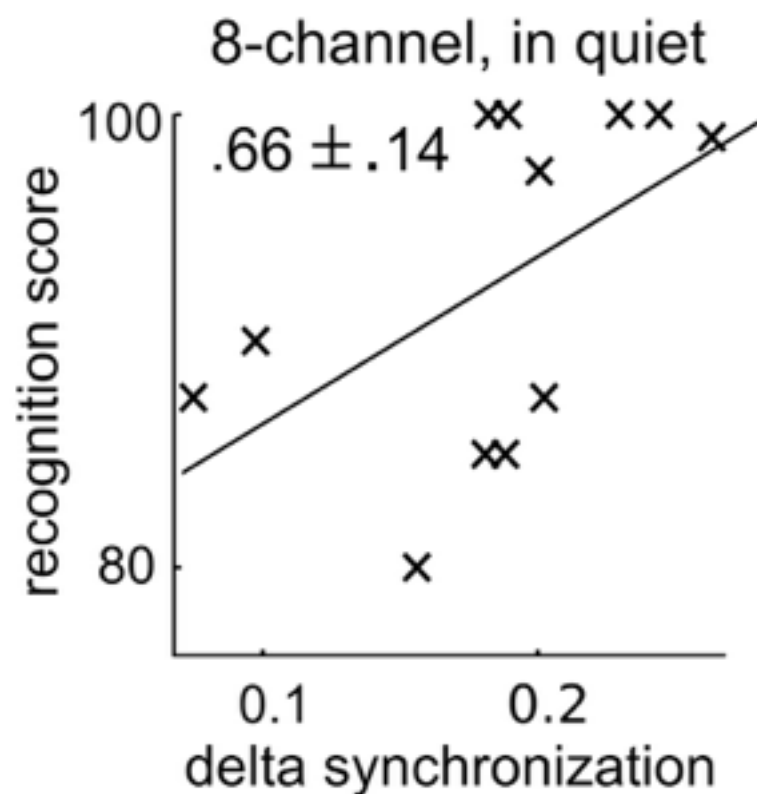
# Noise-Vocoded Speech: Results
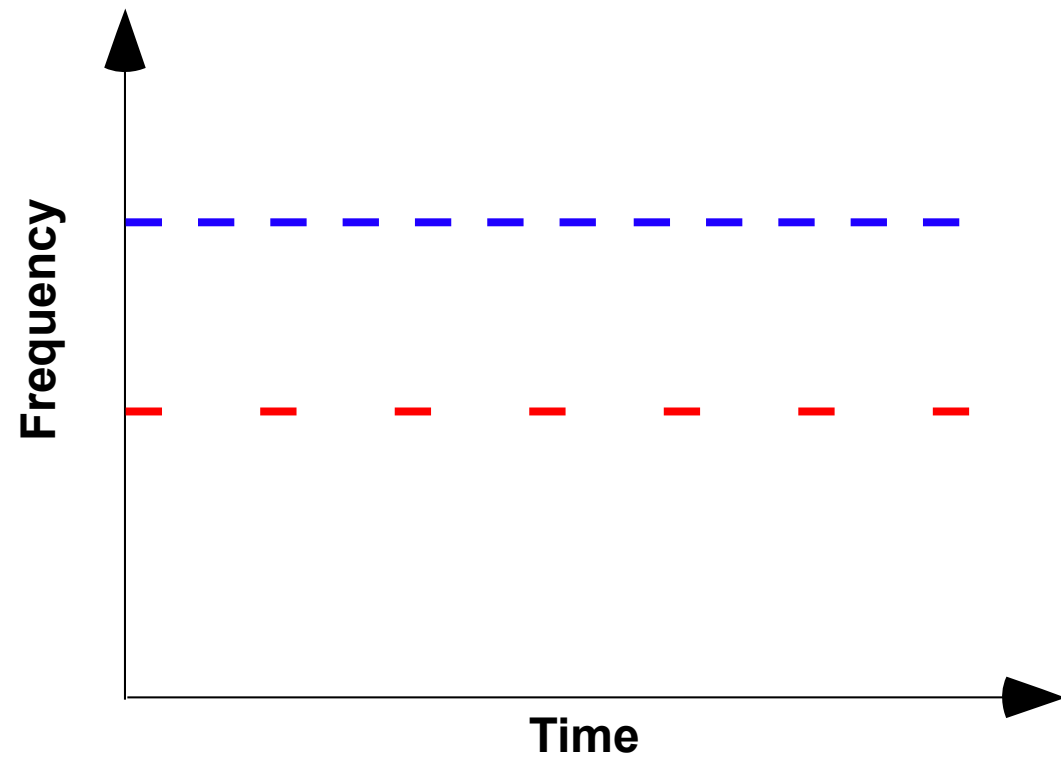


*Neural Synchronization Spectrum*

- 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
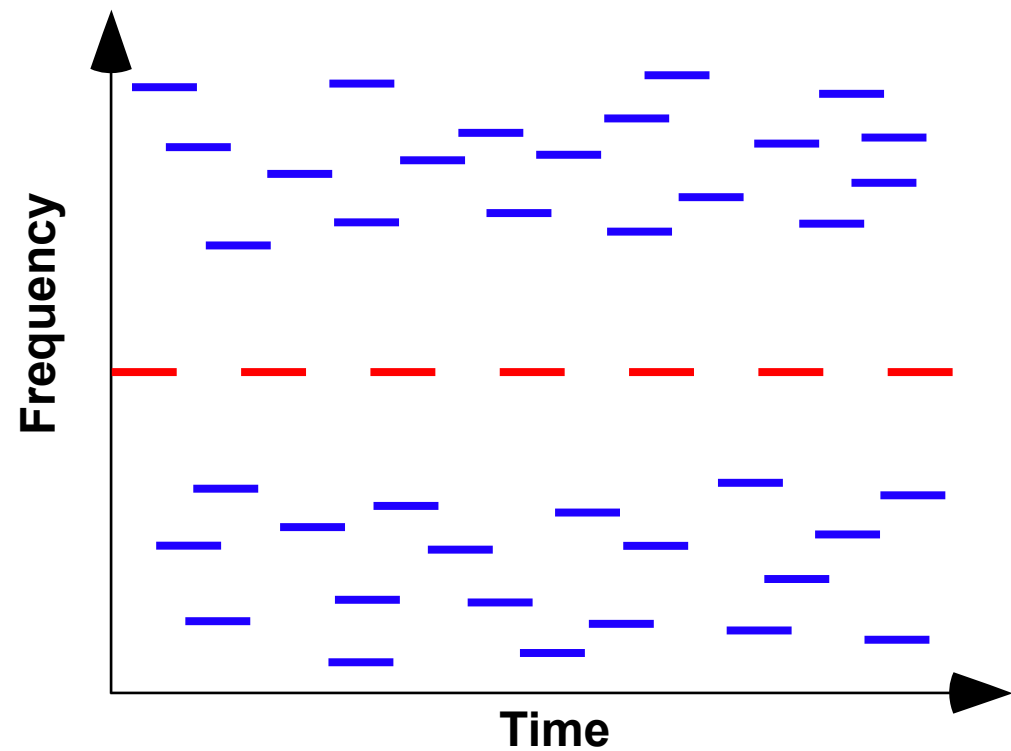
# Noise-Vocoded Speech: Results

# Not Just Speech

Competing Tone Streams

Tone Stream in Masker Cloud



Xiang et al., J Neuroscience (2010)

Elhilali et al., PLoS Biology (2009)

# Summary

- Cortical representations of speech found here:
  - ✓ consistent with being *neural* representations of auditory *perceptual* objects
  - ✓ very robust to noise (~intelligibility)
  - ✓ relies on *spectro*-temporal fine structure
  - ✓ explicitly temporal representation
- Object representation at 100 ms latency (PT), but not by 50 ms (HG)

# Thank You