

Neural Representations of Cocktail Party Speech in Human Auditory Cortex

Jonathan Z. Simon

Department of Biology

Department of Electrical & Computer Engineering

Institute for Systems Research

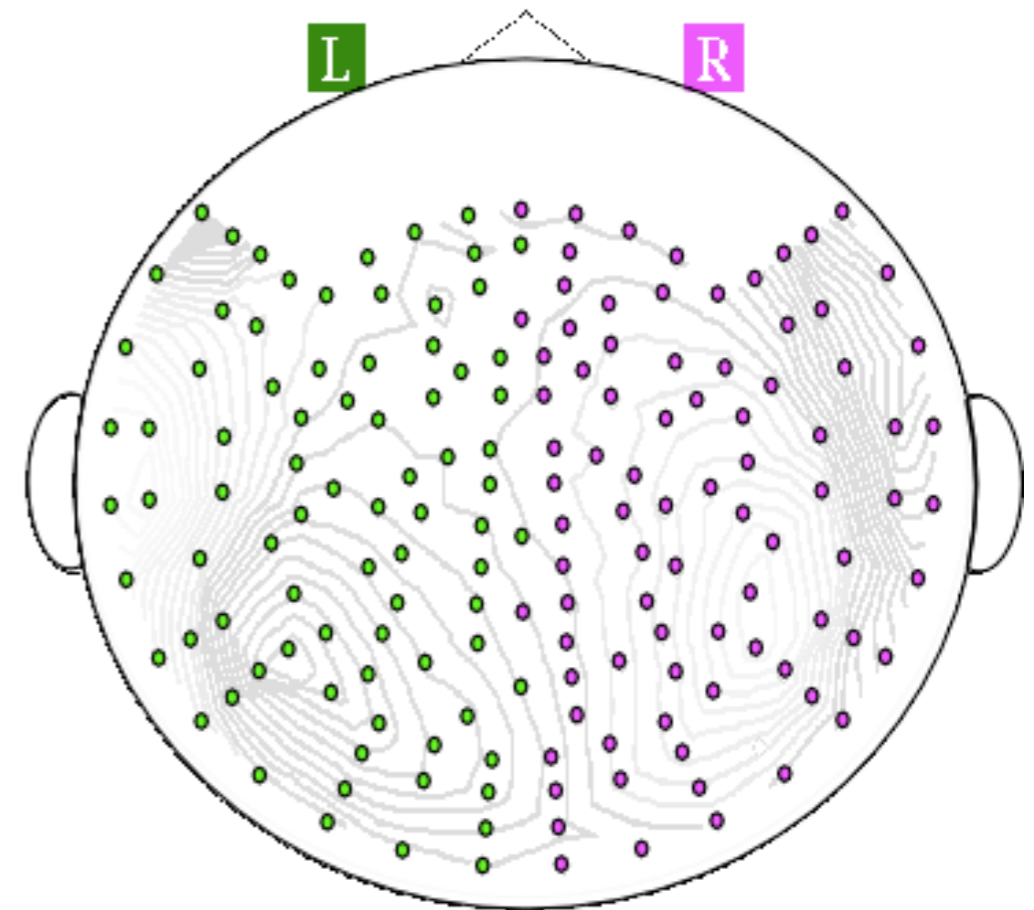
University of Maryland

Outline

- Cortical Representations of Speech (via MEG)
 - ▶ *Encoding vs. Decoding*
- Cortical Representations of Speech in Noise
- Cortical Representations of “Cocktail Party” Speech
- Effects of Aging on Cortical Representations of “Cocktail Party” Speech
- Cortical Representations of Internal Speech

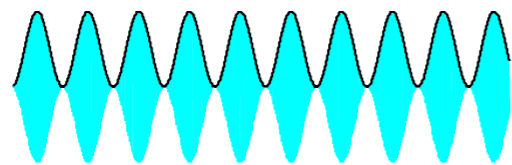
Magnetoencephalography (MEG)

- Non-invasive, Passive, Silent Neural Recordings of Cortex
- 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



MEG Phase-Locked Responses to Slow Acoustic Modulations

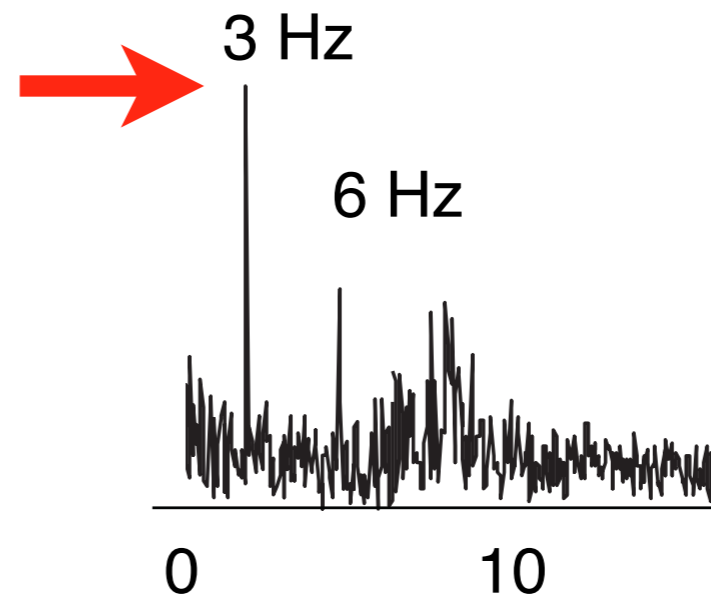
AM at 3 Hz



3 Hz phase-locked response

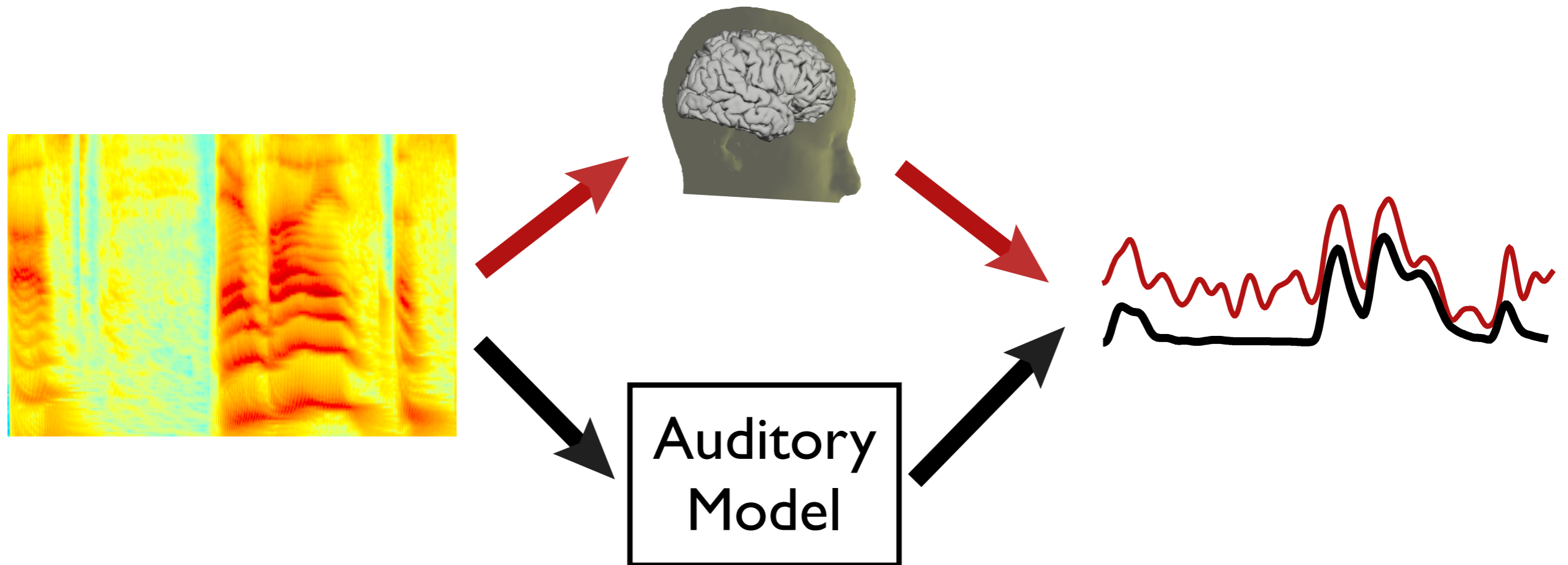


response spectrum (*subject R0747*)

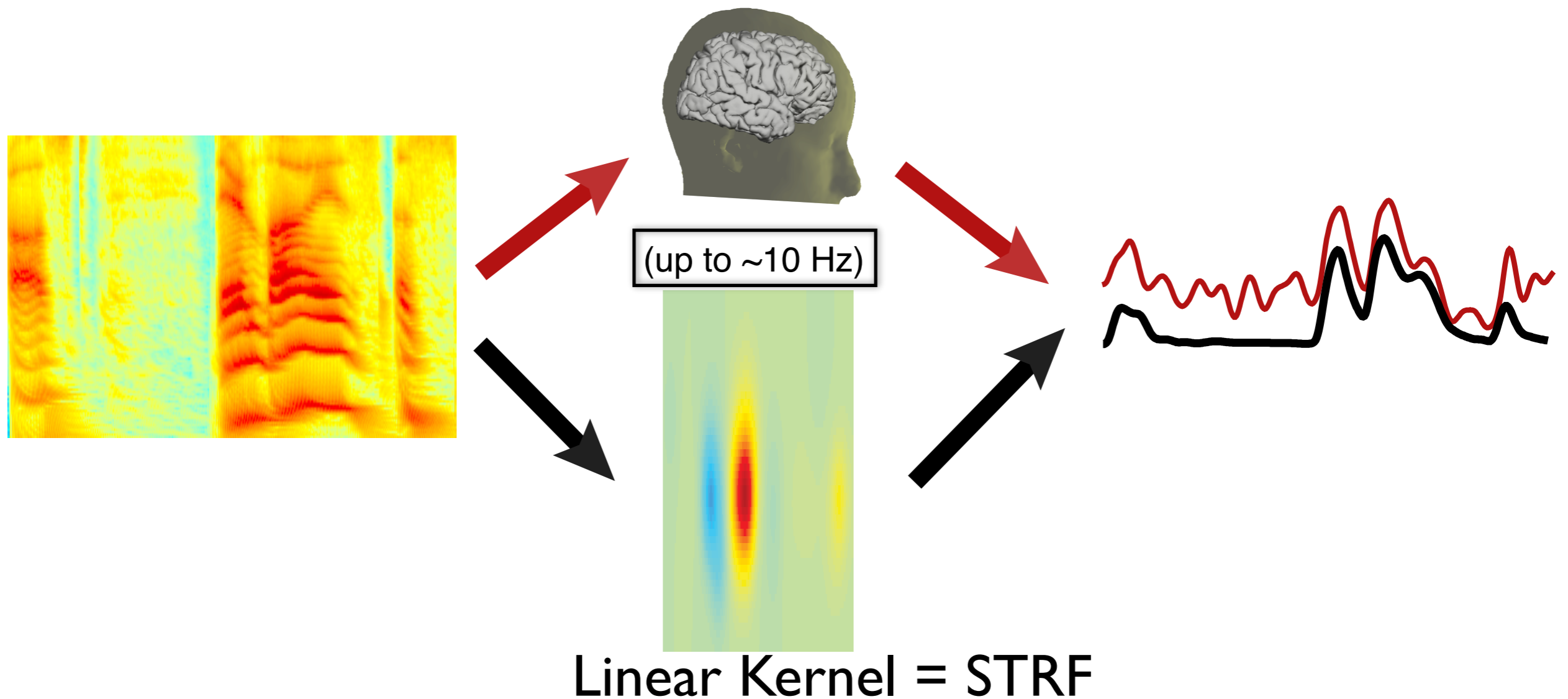


MEG activity is phase-locked to temporal modulations of sound

MEG Responses to Speech Modulations



MEG Responses Predicted by STRF Model

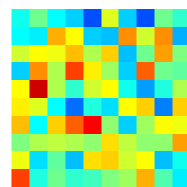


Neural Reconstruction of Speech Envelope

Speech Envelope

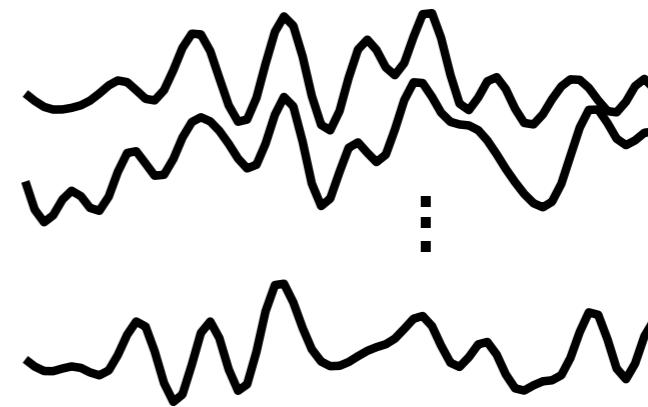


Decoder

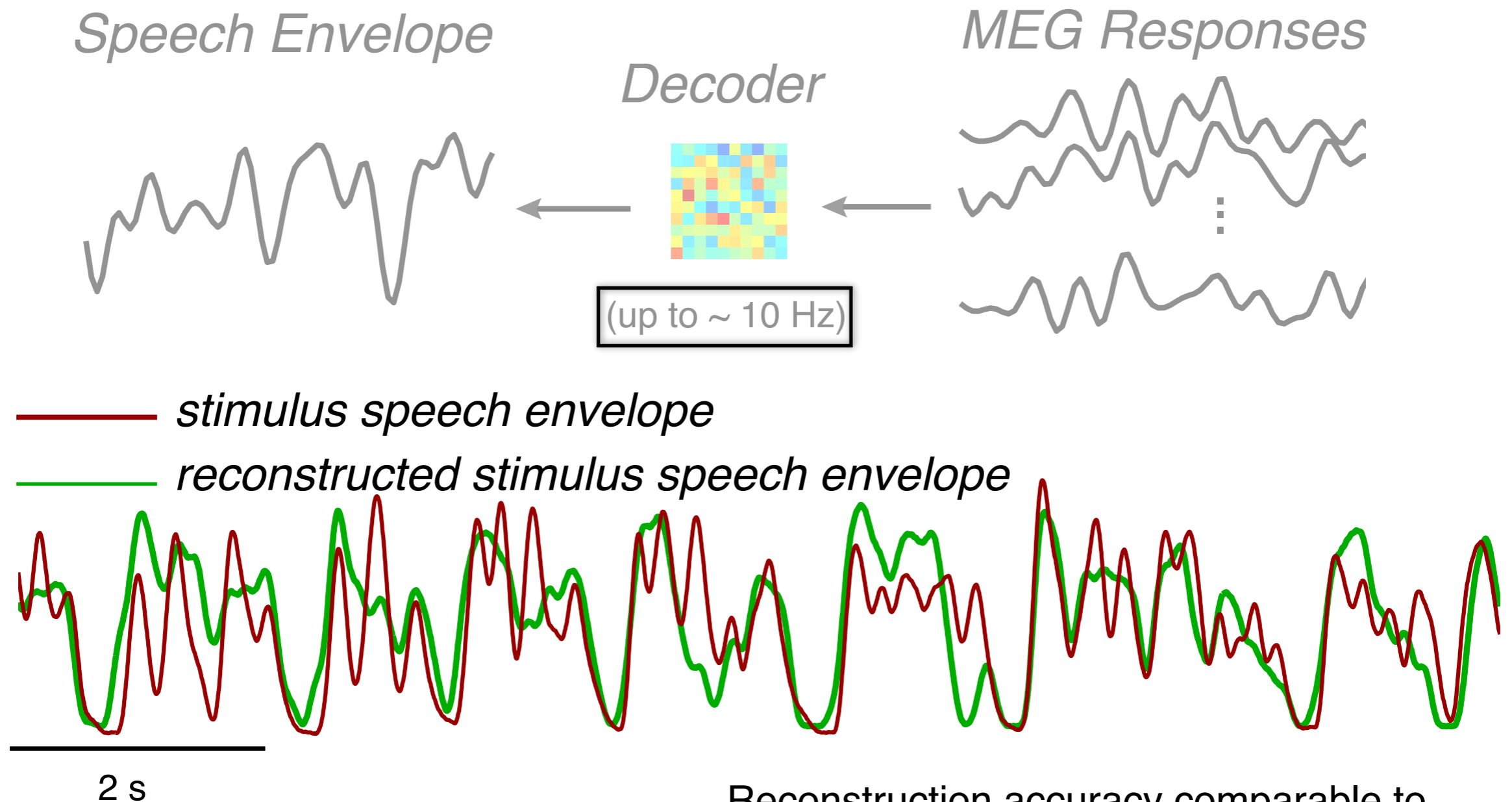


(up to ~ 10 Hz)

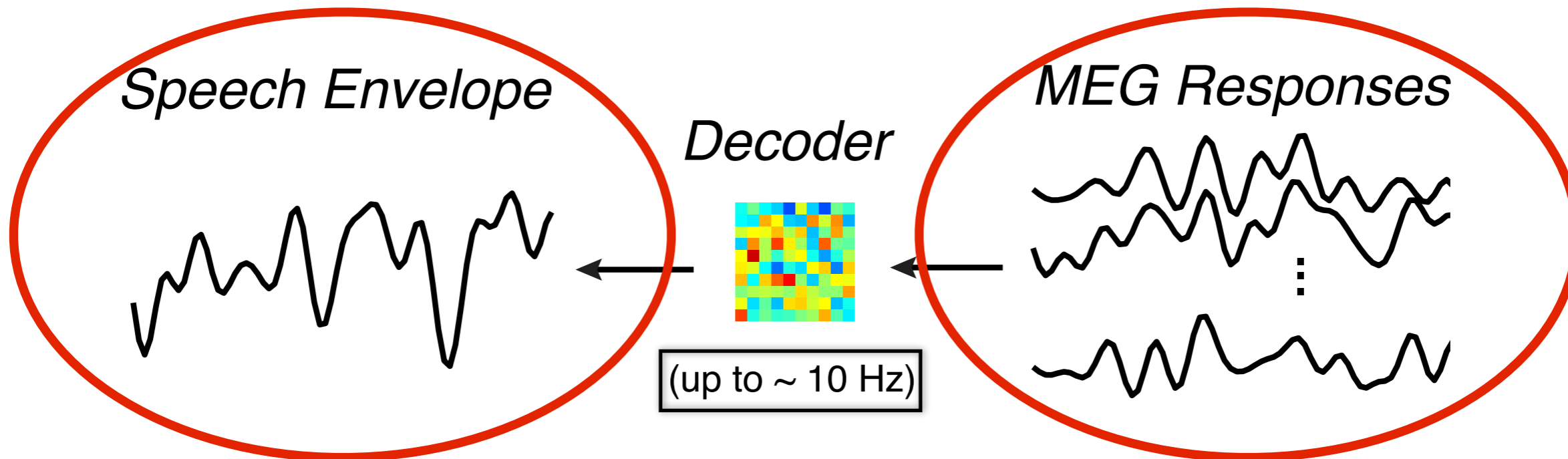
MEG Responses



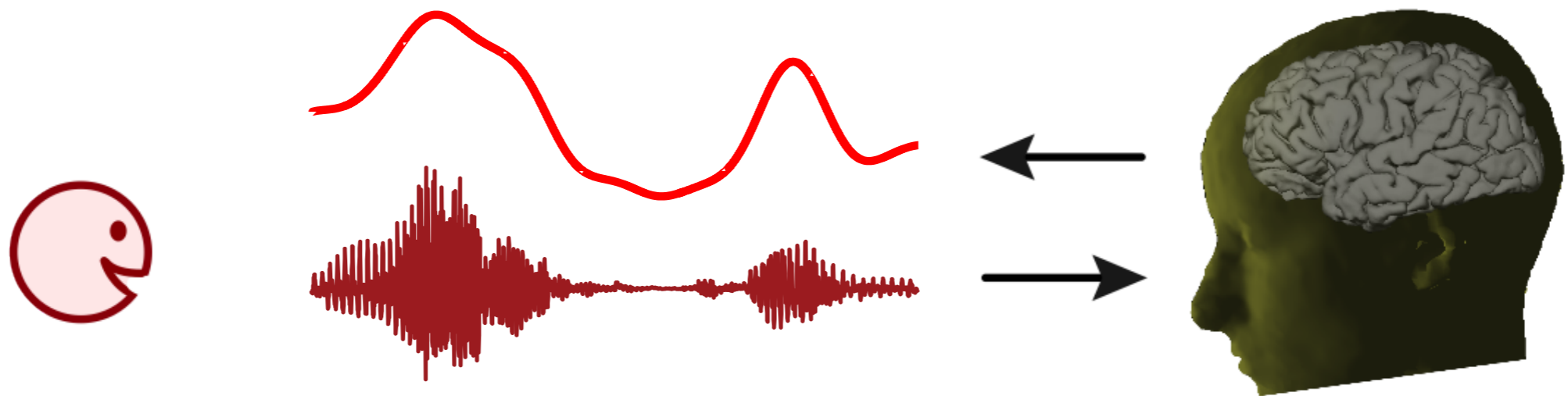
Neural Reconstruction of Speech Envelope



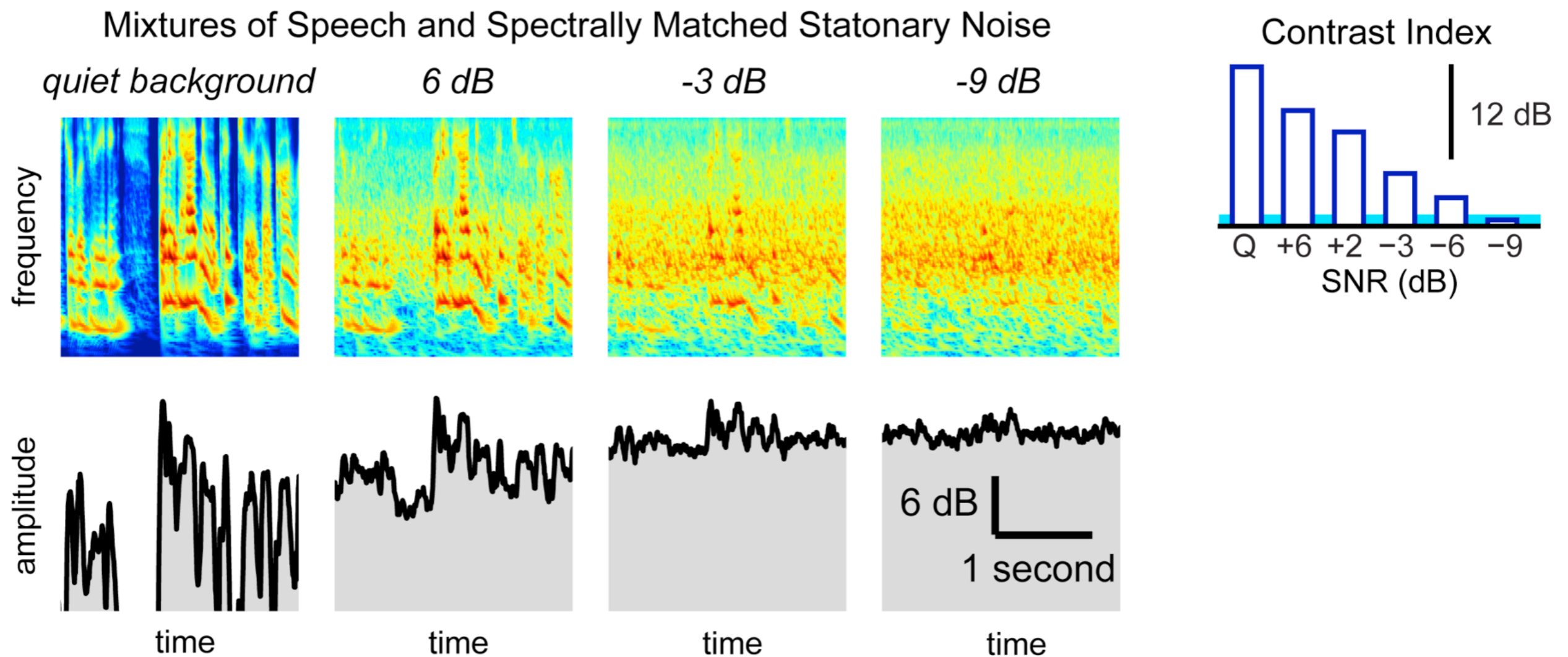
Reconstruction accuracy comparable to single unit & ECoG recordings



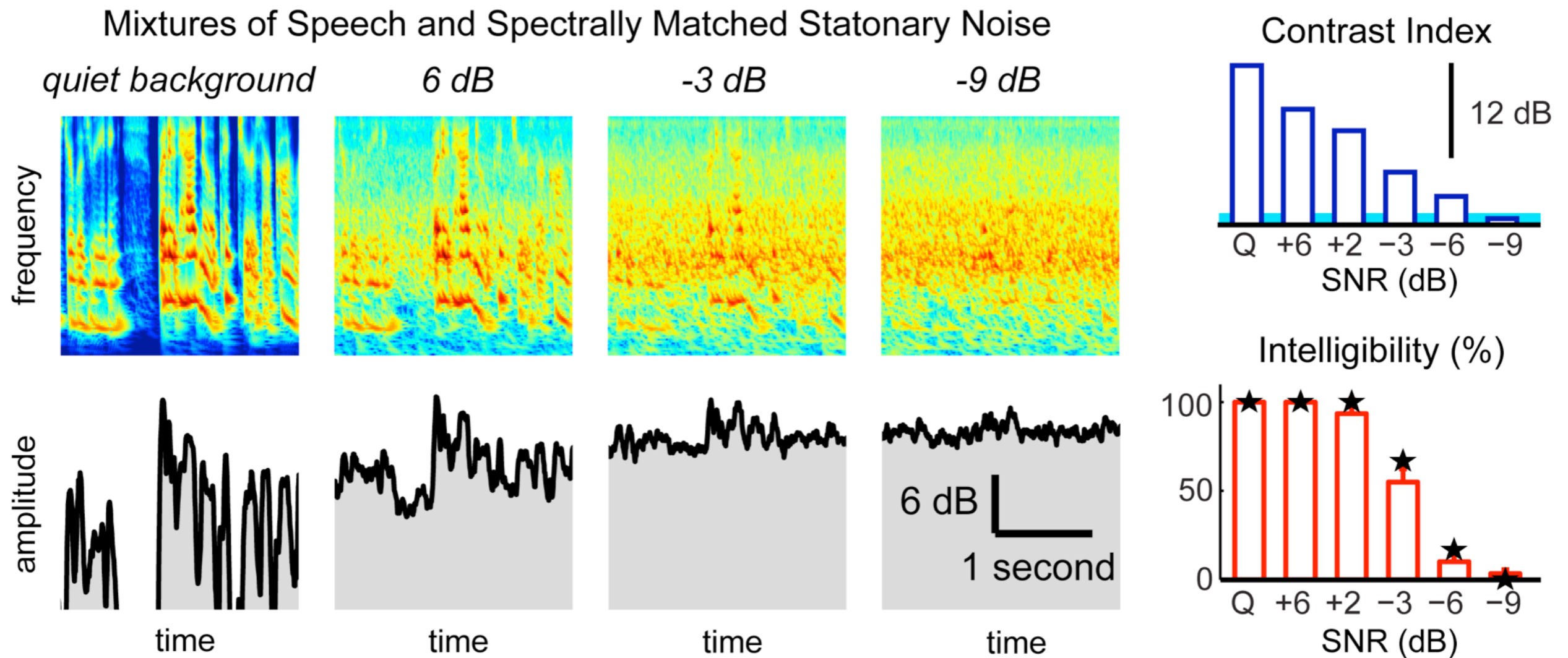
Neural Representation of Speech: Temporal



Speech in Stationary Noise

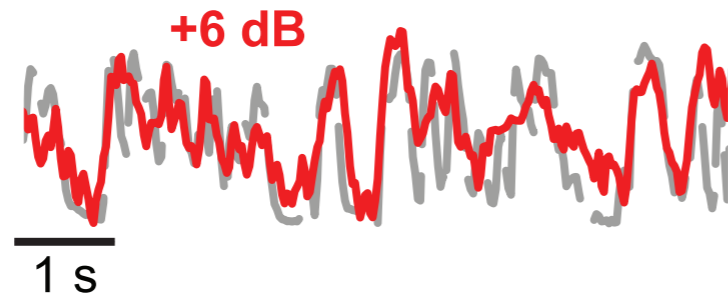


Speech in Stationary Noise



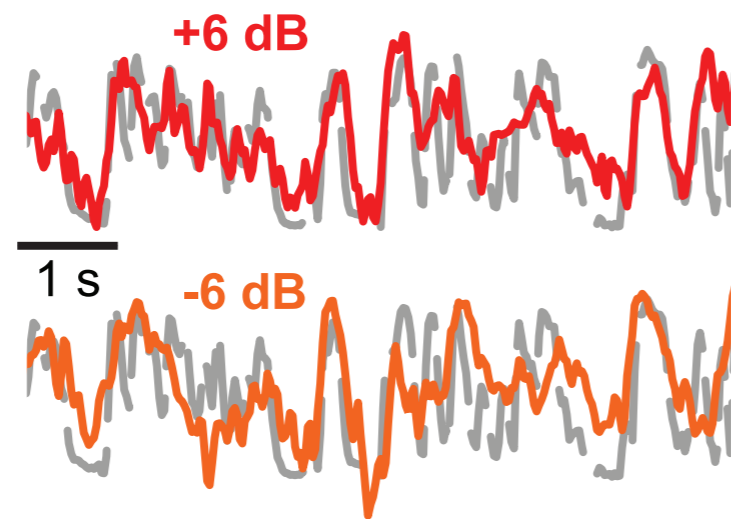
Speech in Noise: Results

Neural Reconstruction of
Underlying Speech Envelope



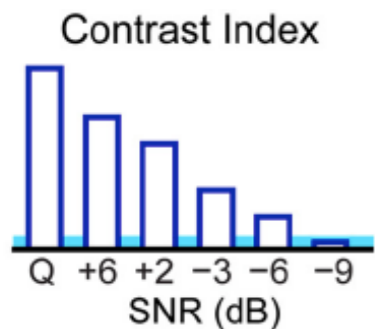
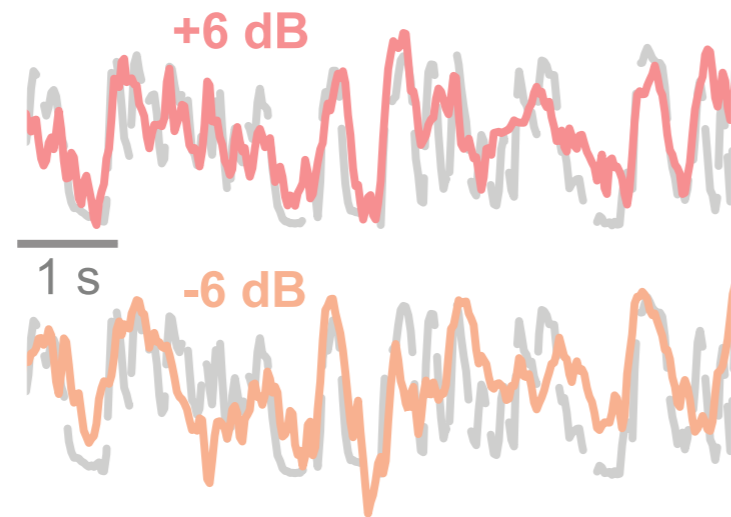
Speech in Noise: Results

Neural Reconstruction of
Underlying Speech Envelope

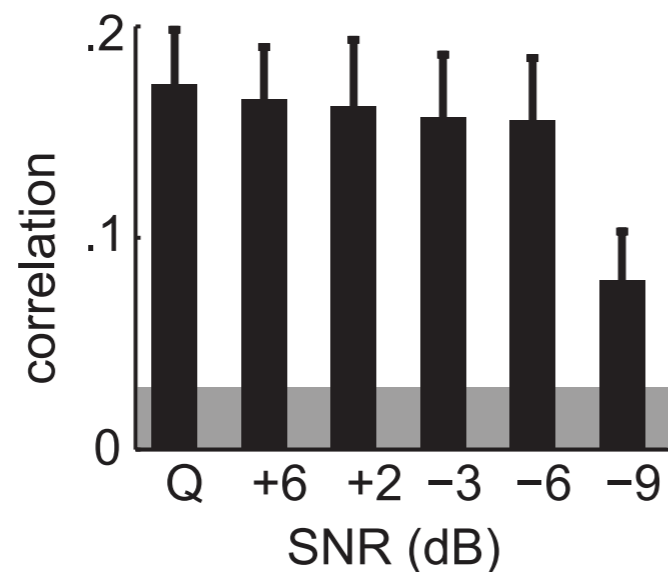


Speech in Noise: Results

Neural Reconstruction of Underlying Speech Envelope

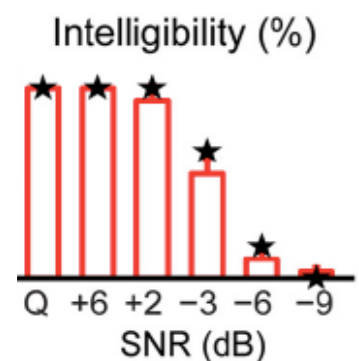
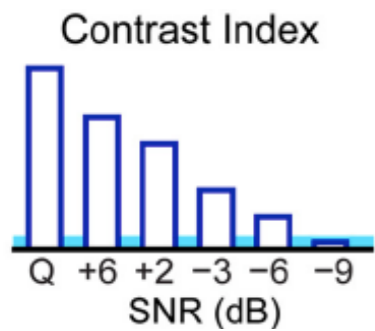
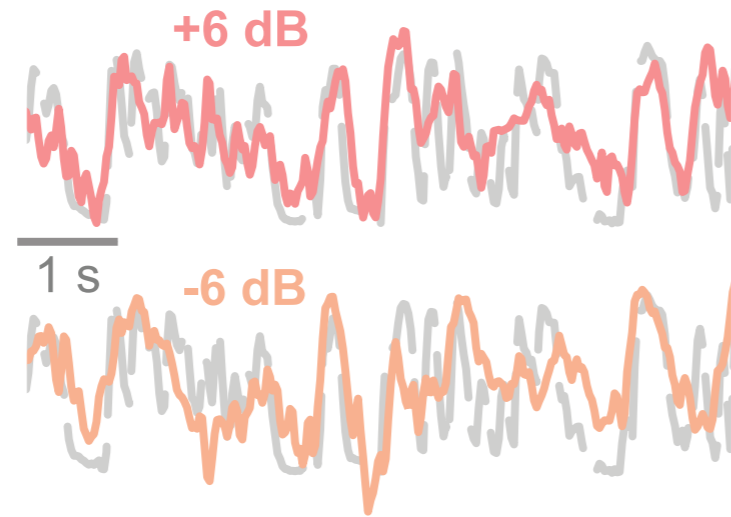


Reconstruction Accuracy

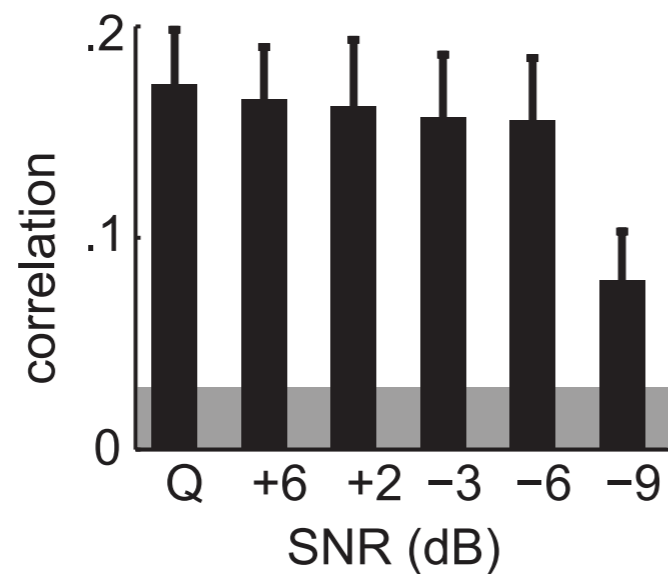


Speech in Noise: Results

Neural Reconstruction of Underlying Speech Envelope

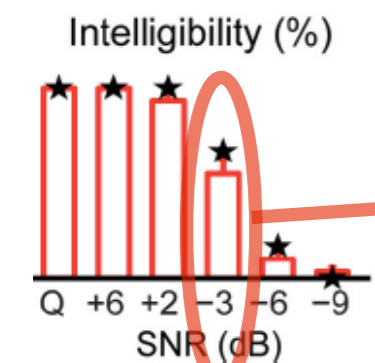
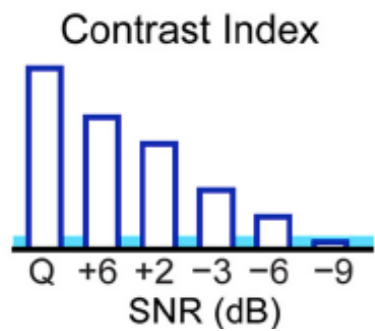
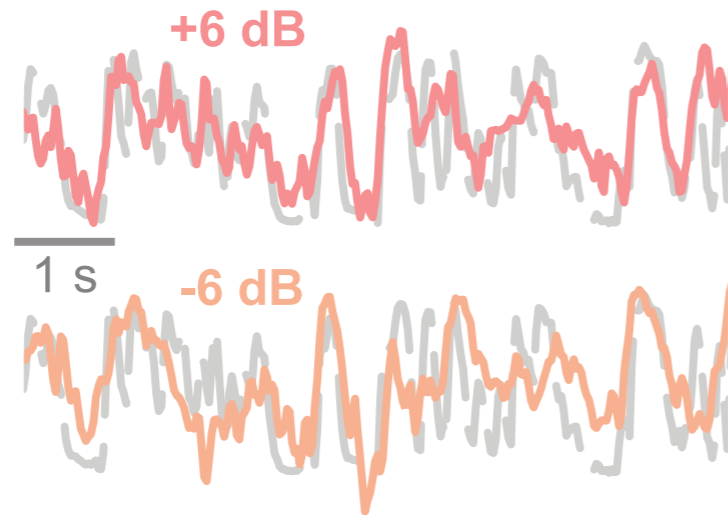


Reconstruction Accuracy

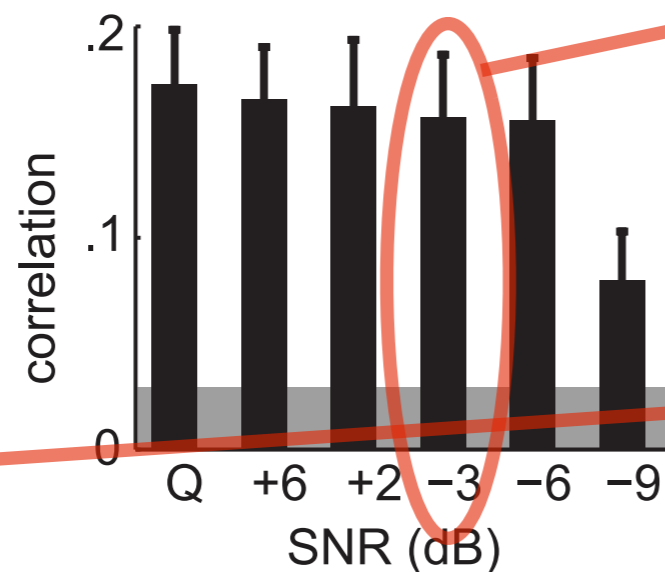


Speech in Noise: Results

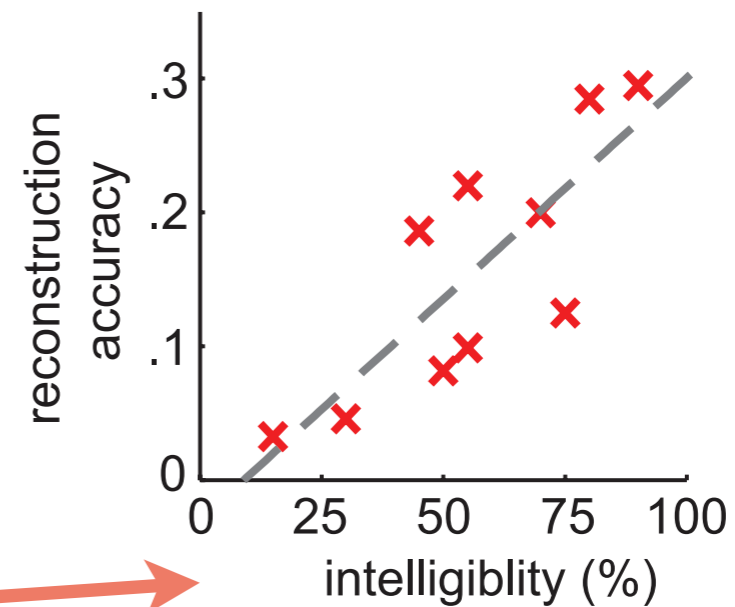
Neural Reconstruction of Underlying Speech Envelope



Reconstruction Accuracy

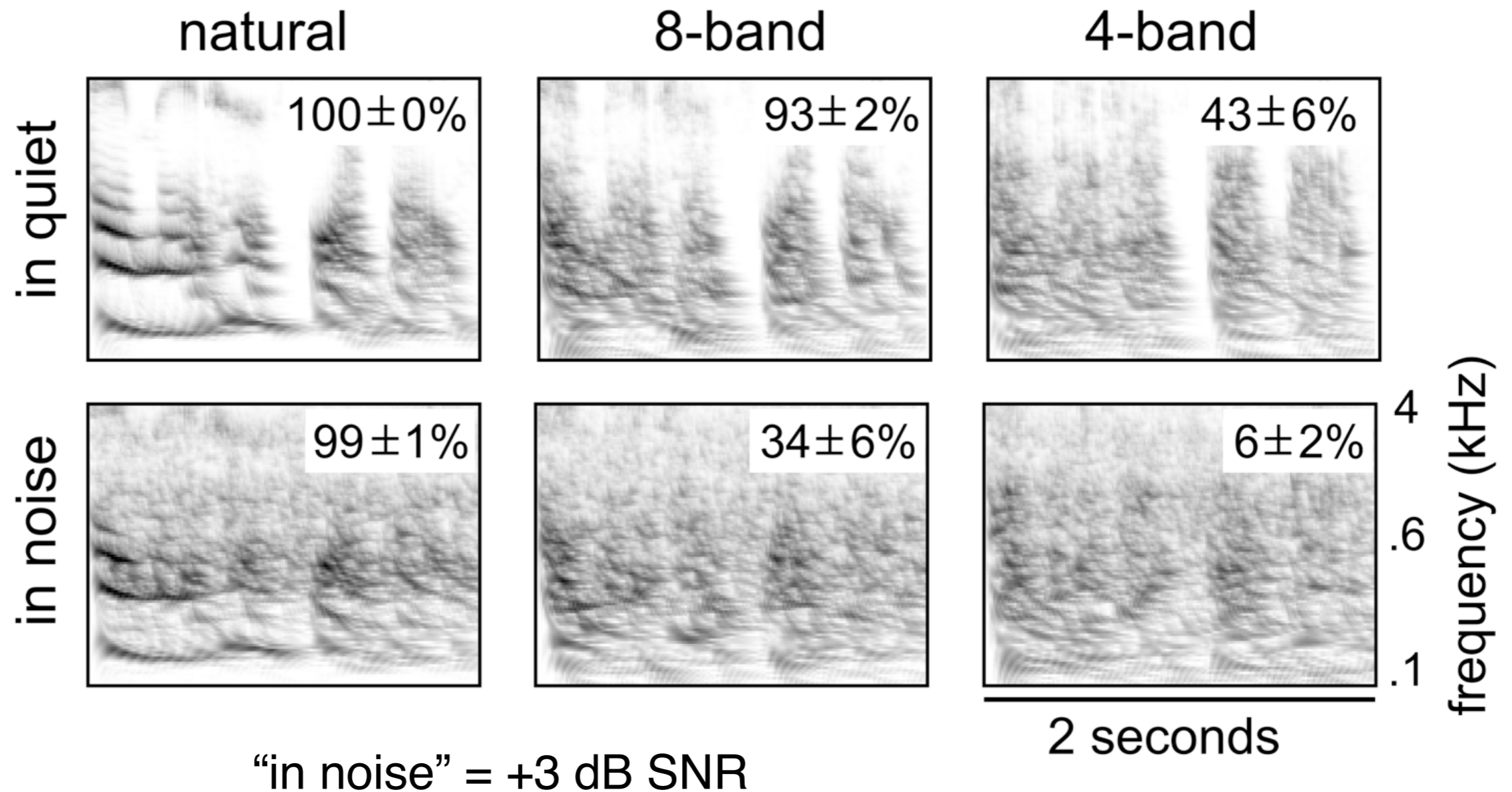


Correlation with Intelligibility

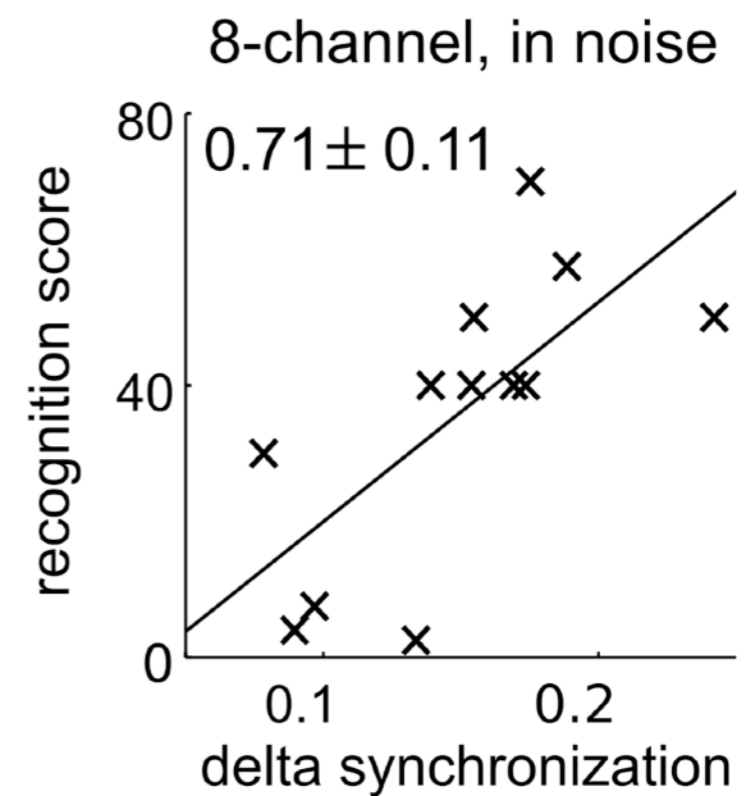
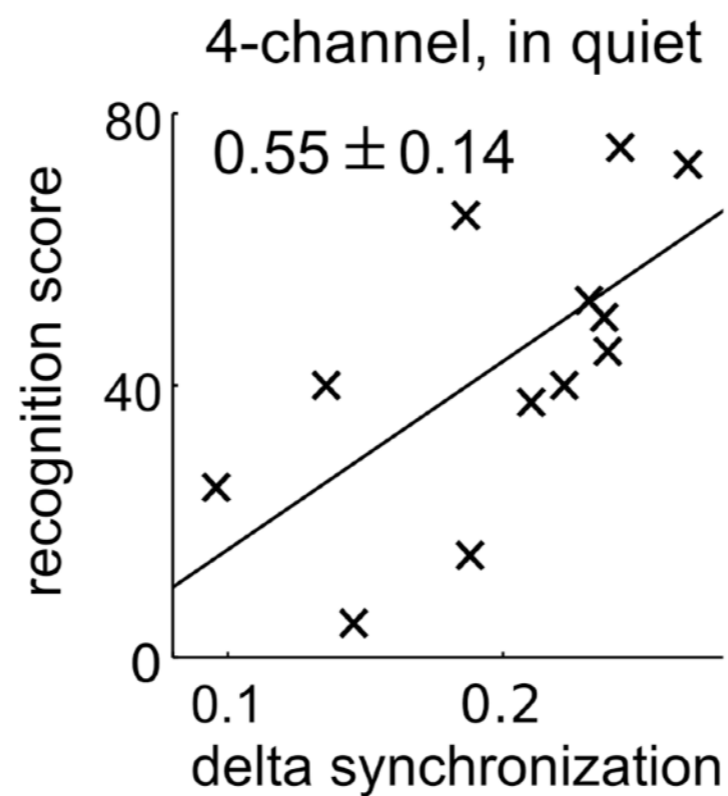
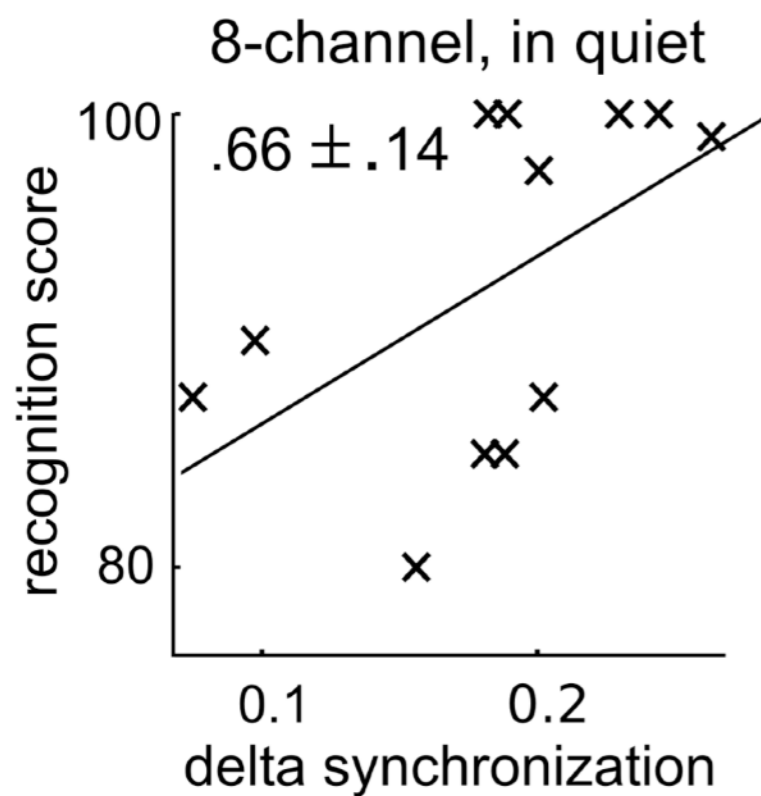


across Subjects

Noise-Vocoded Speech



Noise-Vocoded Speech: Results



- Intelligibility linked to response in Delta band

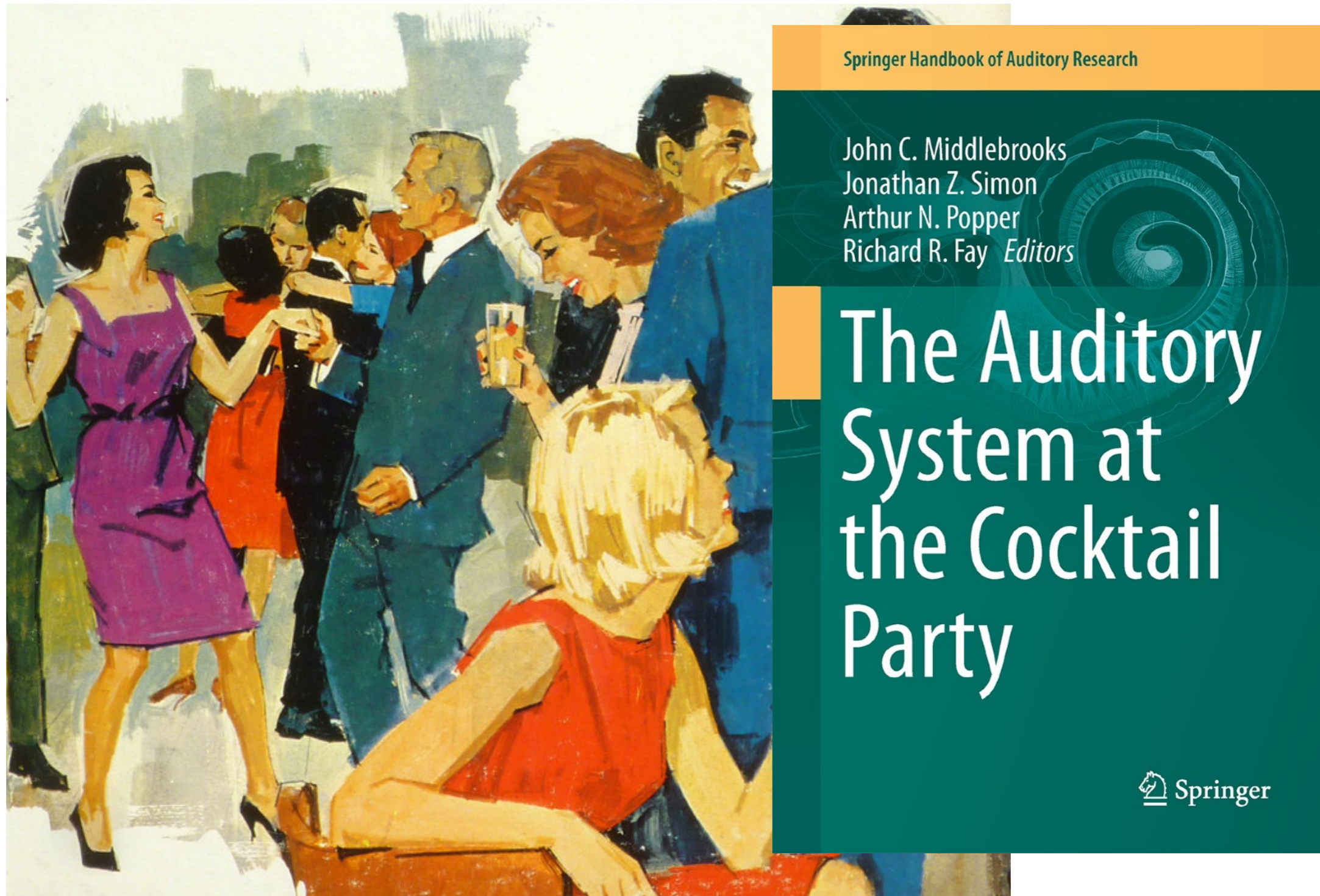
Cortical Speech Representations

- Neural Representations: Encoding & Decoding
- Linear models: Useful & Robust
- Speech **Envelope** only
- Envelope Rates: $\sim 1 - 10$ Hz
- Intelligibility linked to lower range of frequencies (Delta)

Listening to Speech at the Cocktail Party



Listening to Speech at the Cocktail Party



Springer Handbook of Auditory Research

John C. Middlebrooks
Jonathan Z. Simon
Arthur N. Popper
Richard R. Fay *Editors*

The Auditory System at the Cocktail Party

 Springer

Listening to Speech at the Cocktail Party



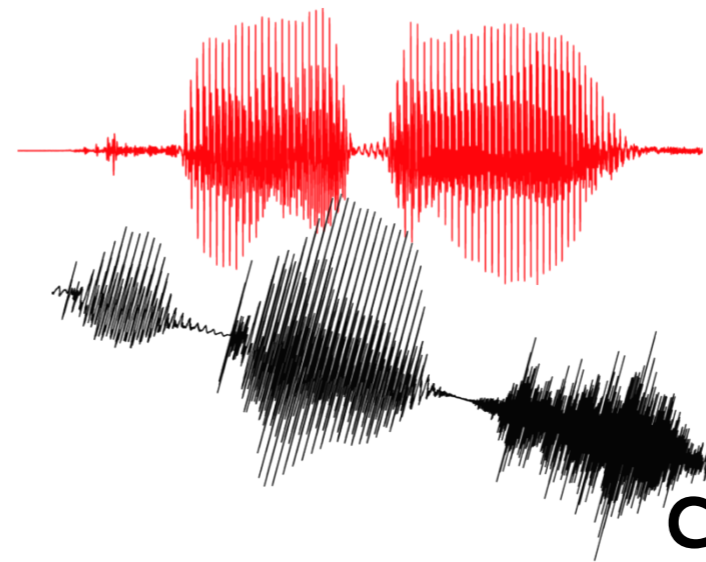
Listening to Speech at the Cocktail Party



Listening to Speech at the Cocktail Party



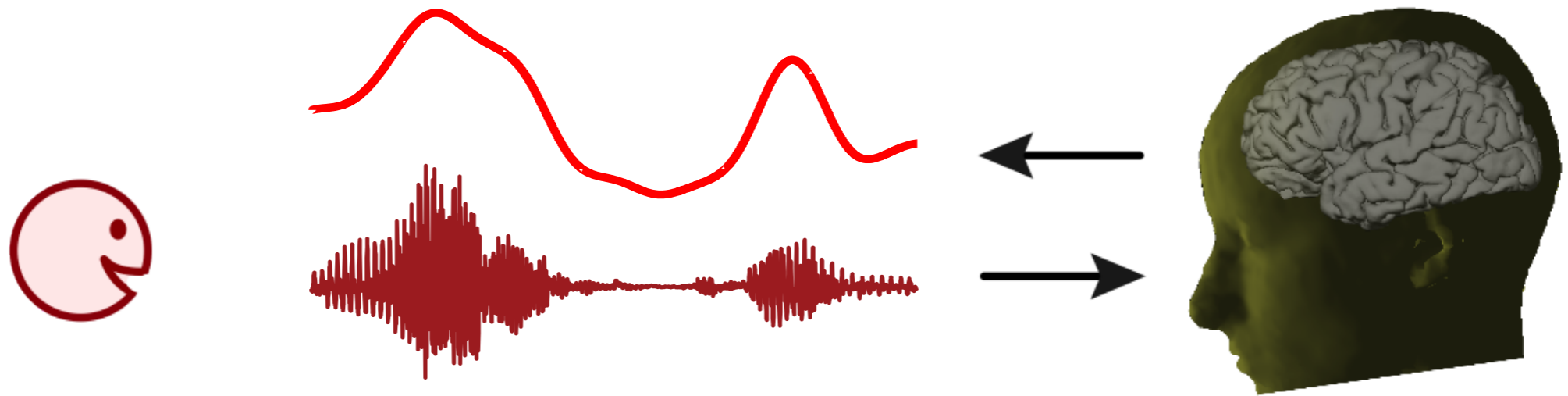
Two Competing Speakers



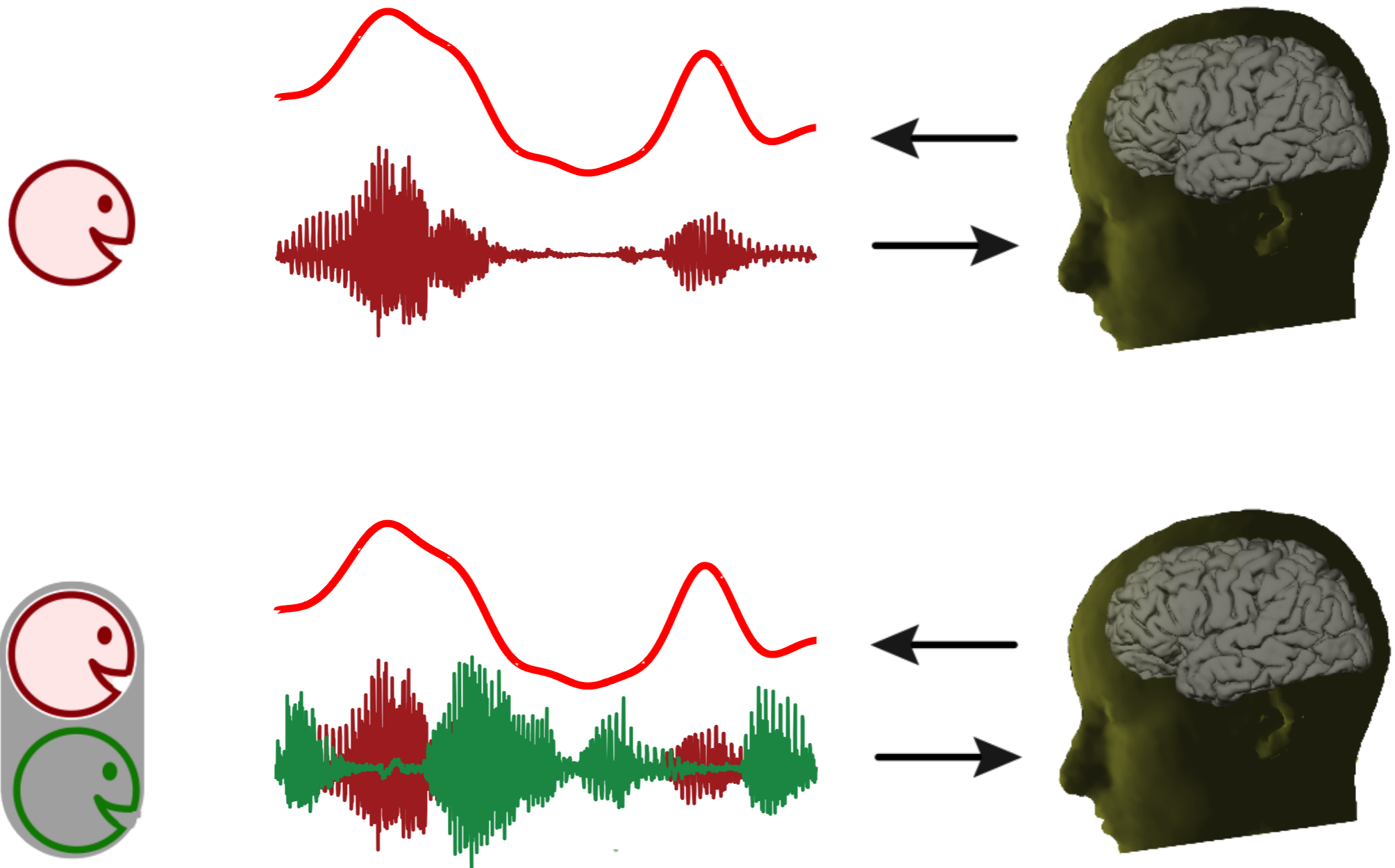
speech

competing speech

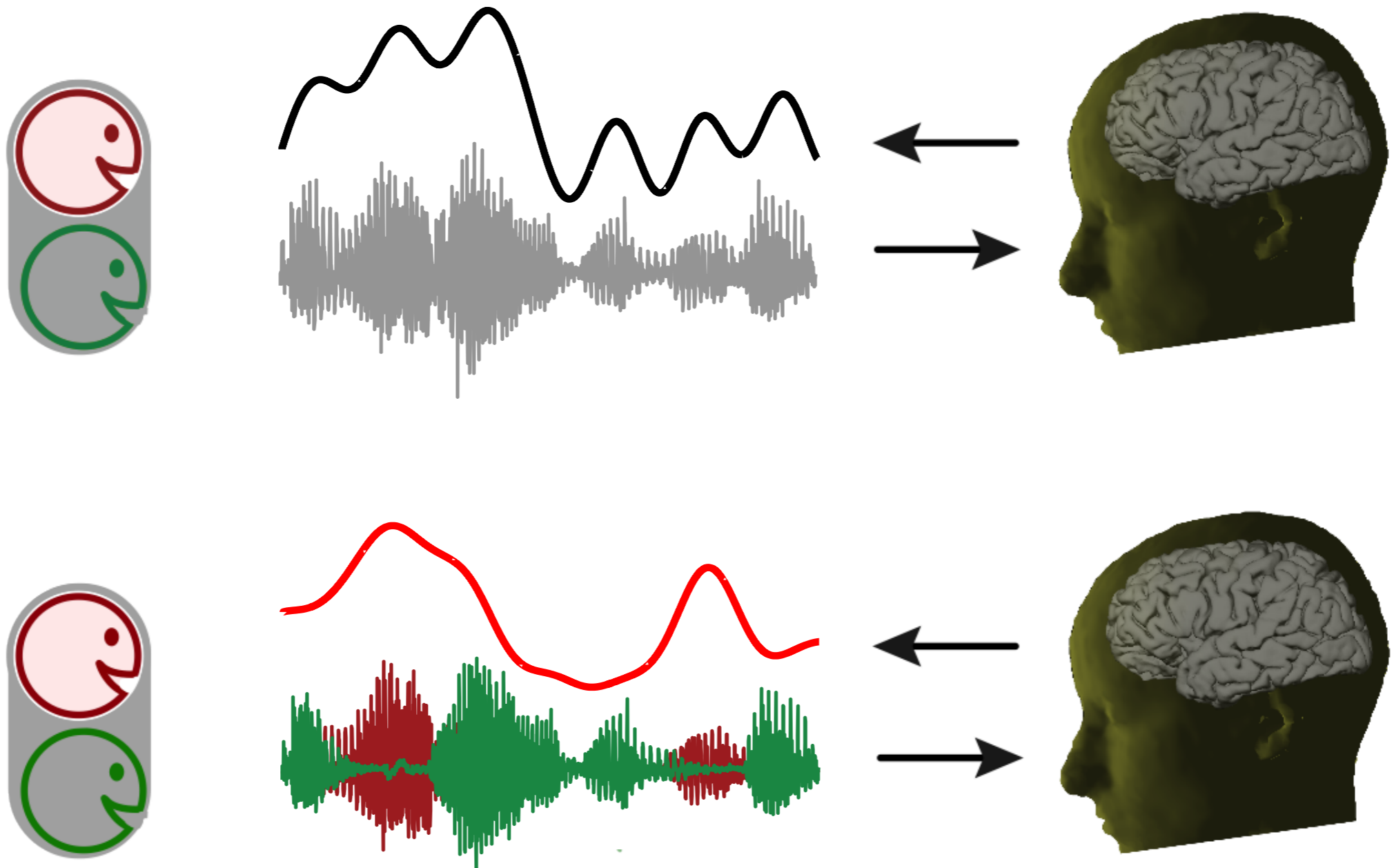
Selective Neural Encoding



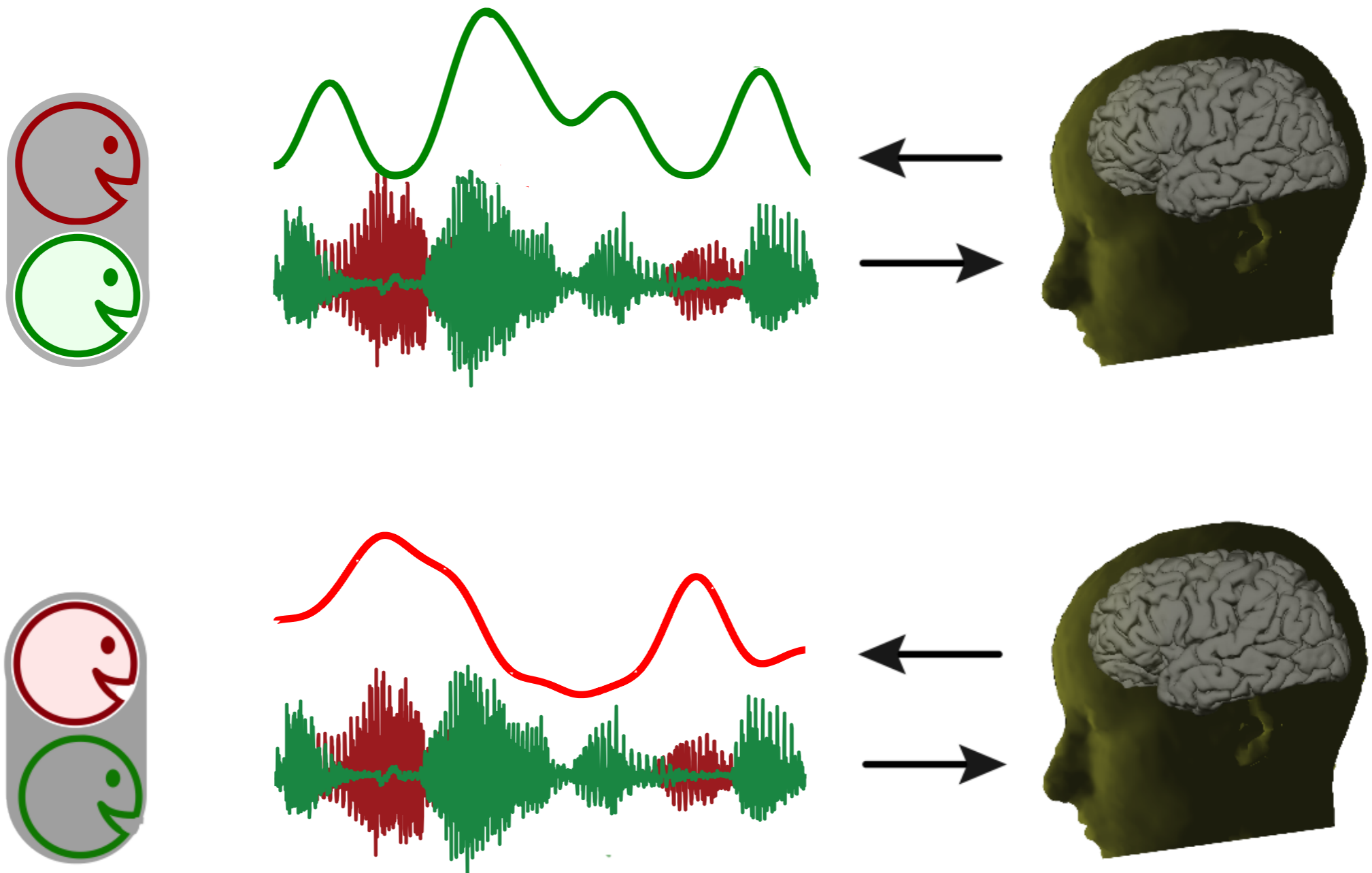
Selective Neural Encoding



Unselective vs. Selective Neural Encoding

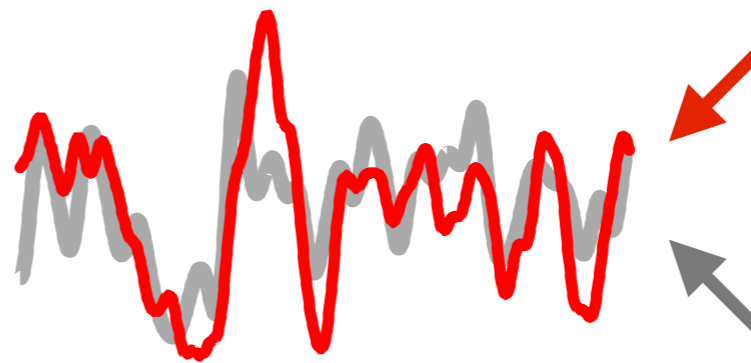


Selective Neural Encoding



Stream-Specific Representation

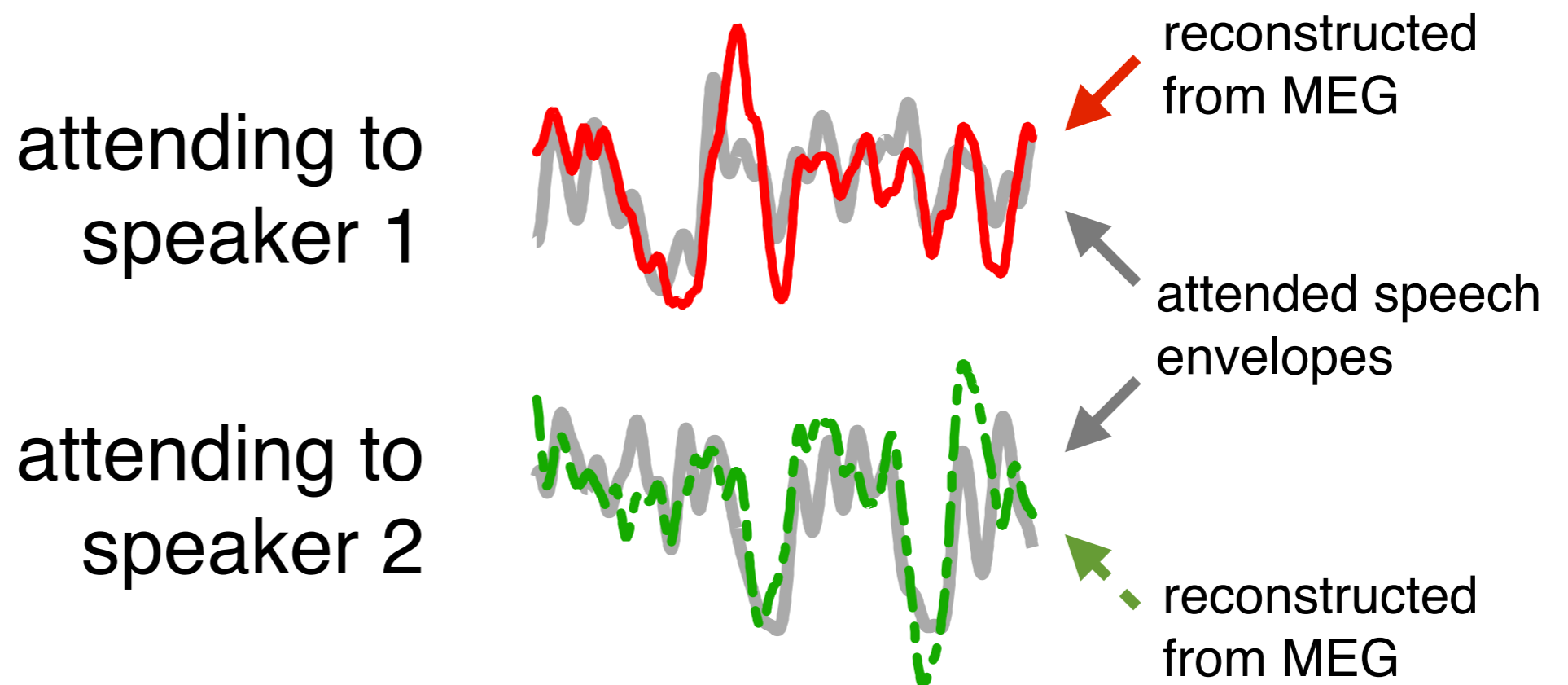
attending to
speaker 1



reconstructed
from MEG

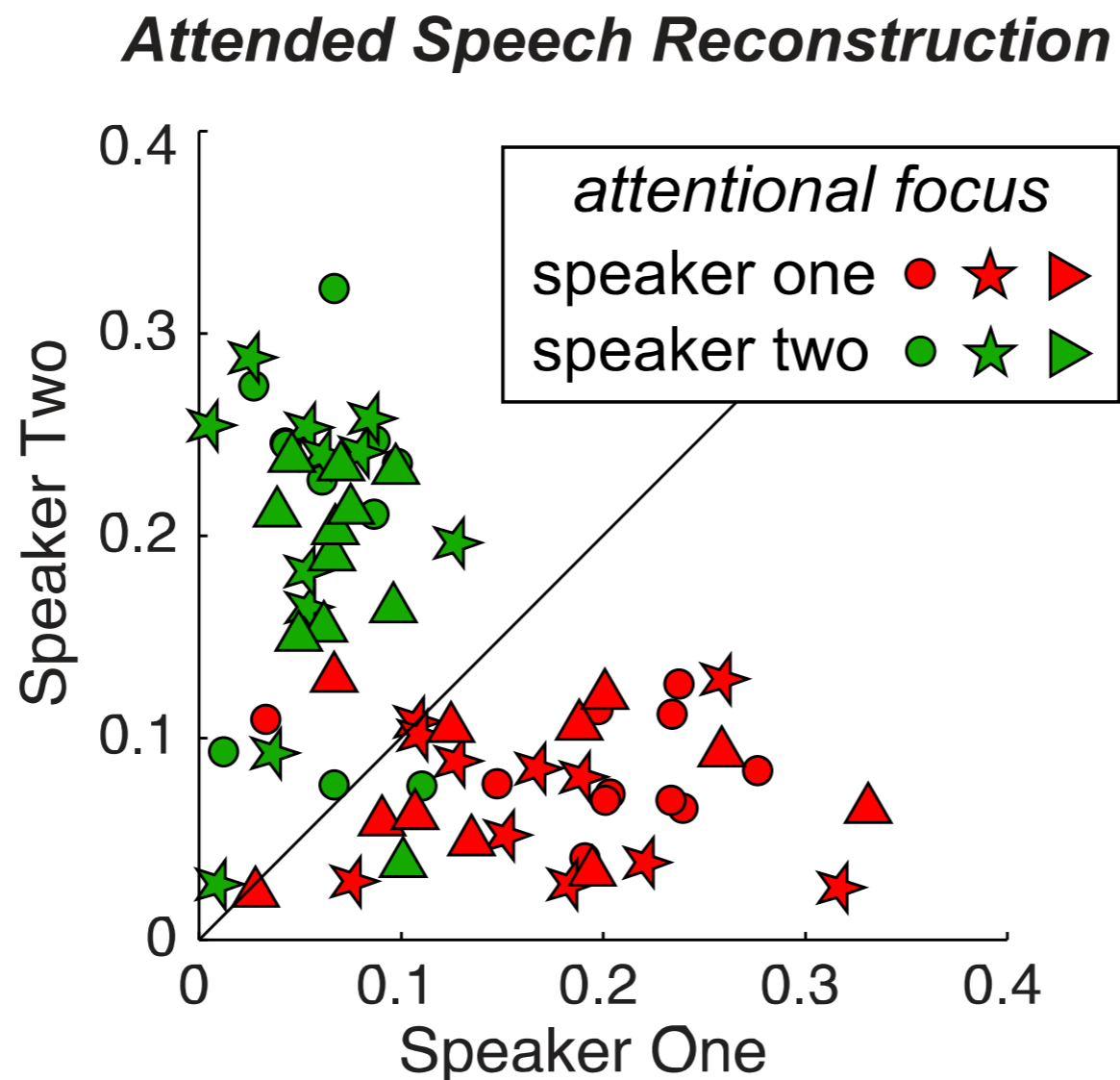
attended speech
envelopes

Stream-Specific Representation



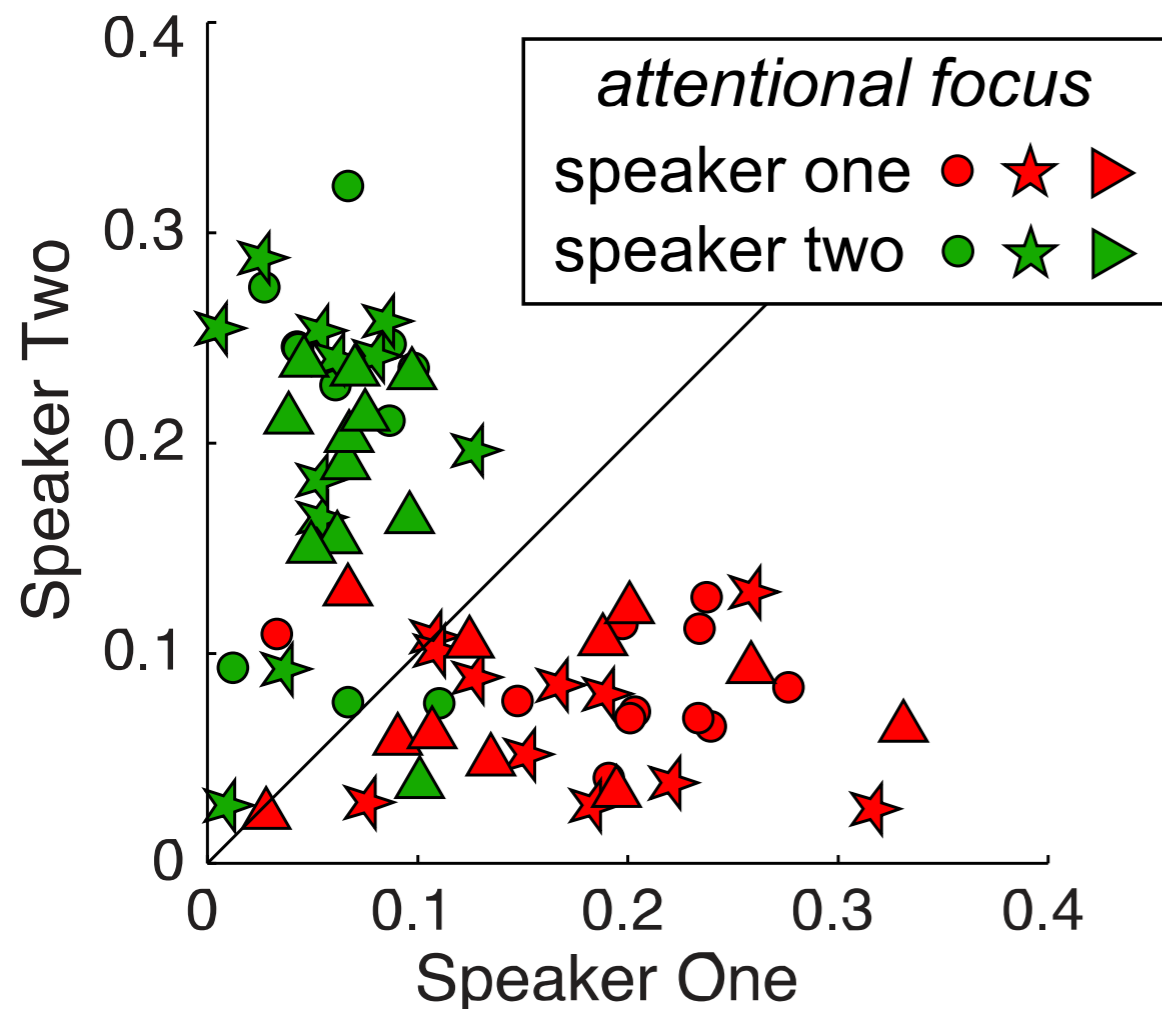
Identical Stimuli!

Single Trial Speech Reconstruction

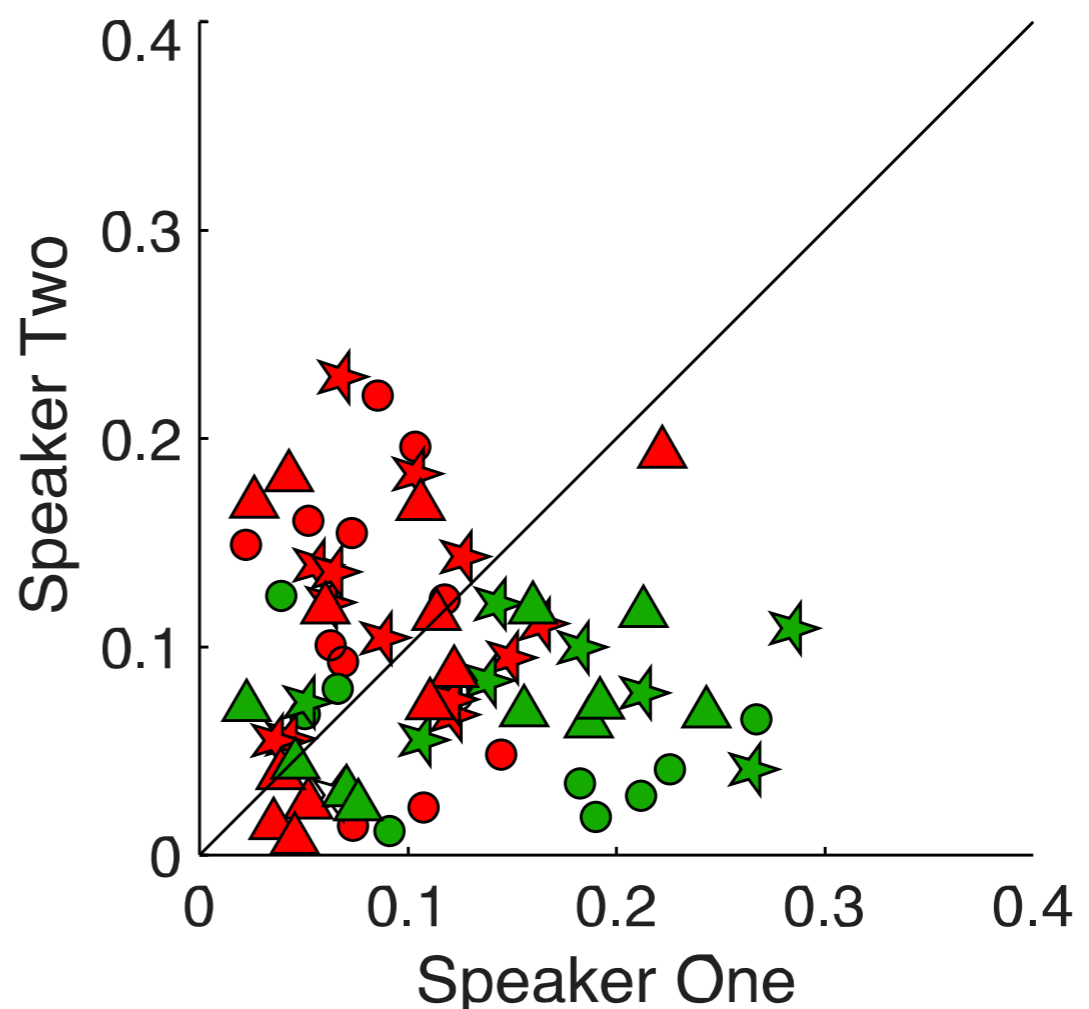


Single Trial Speech Reconstruction

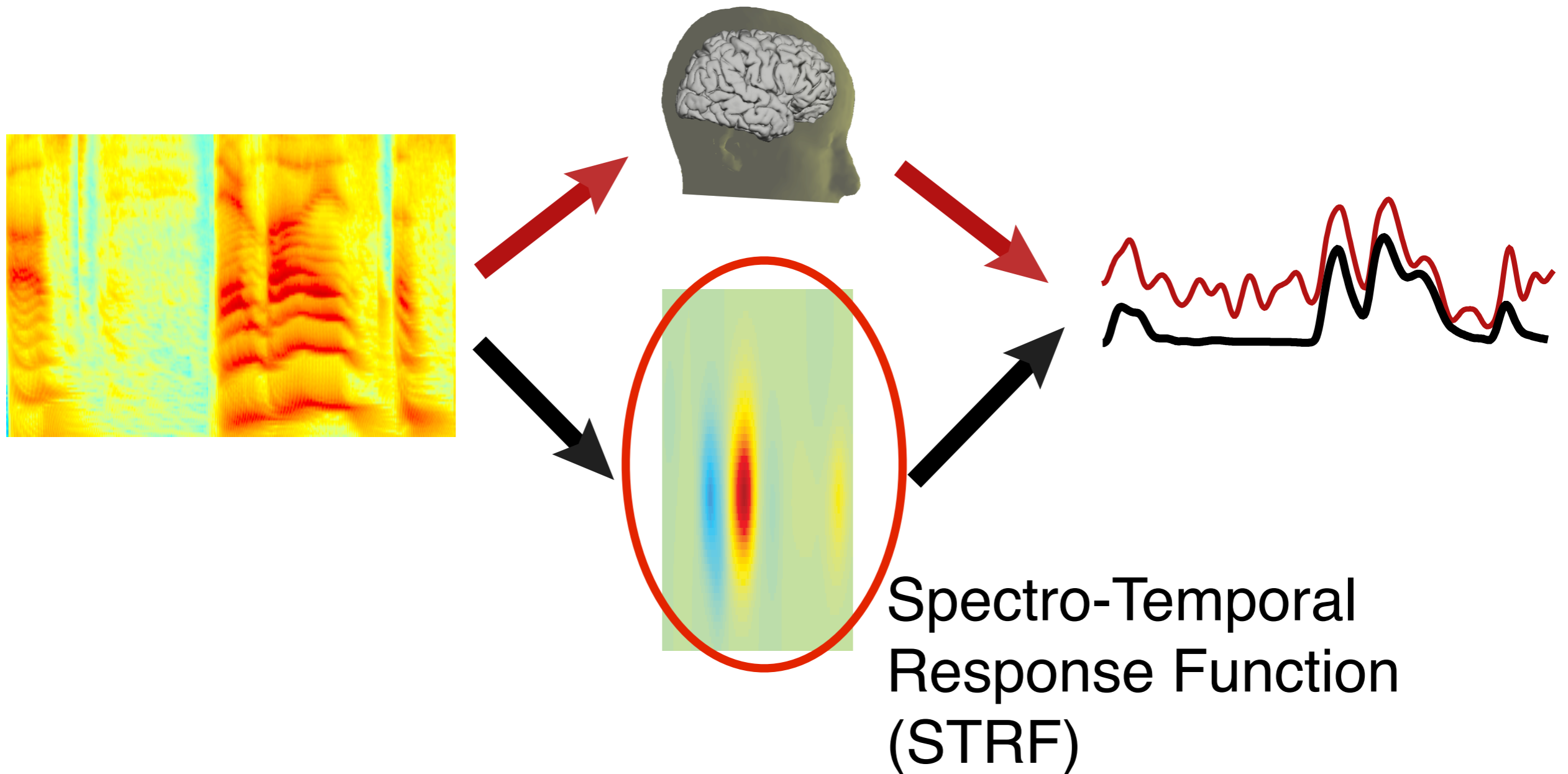
Attended Speech Reconstruction



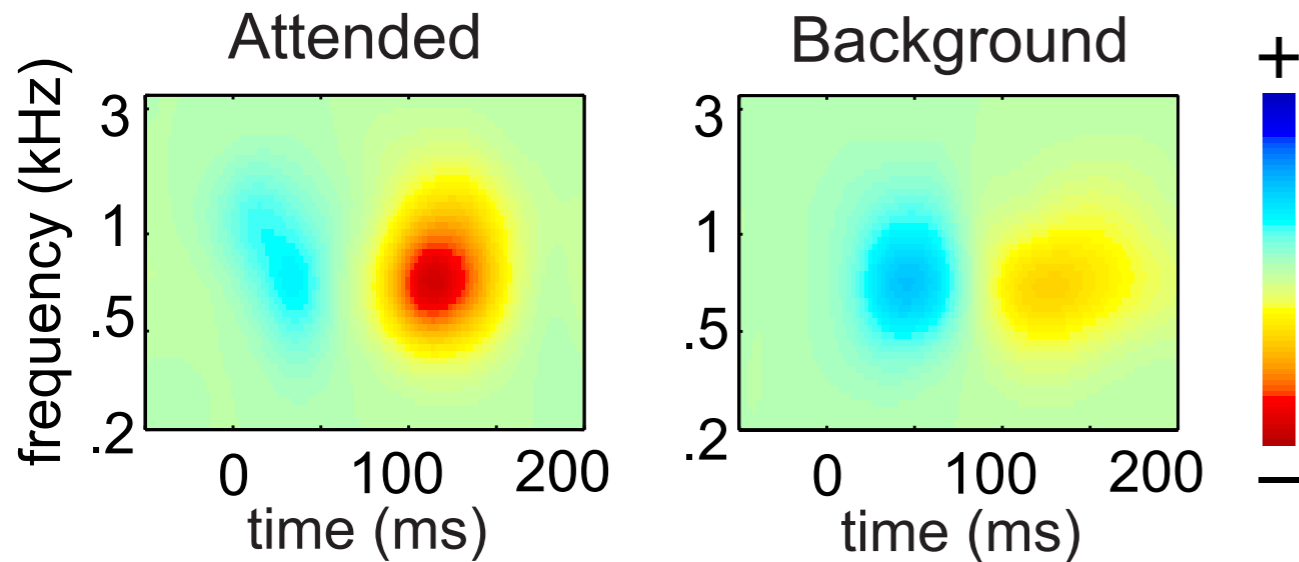
Background Speech Reconstruction



Forward STRF Model

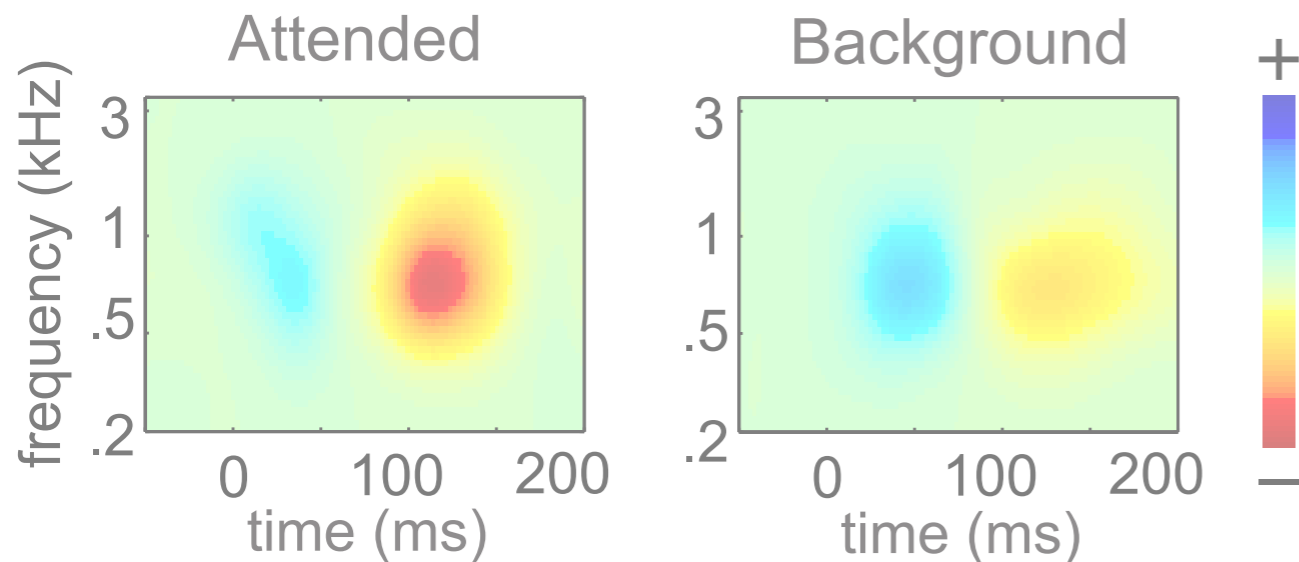


STRF Results

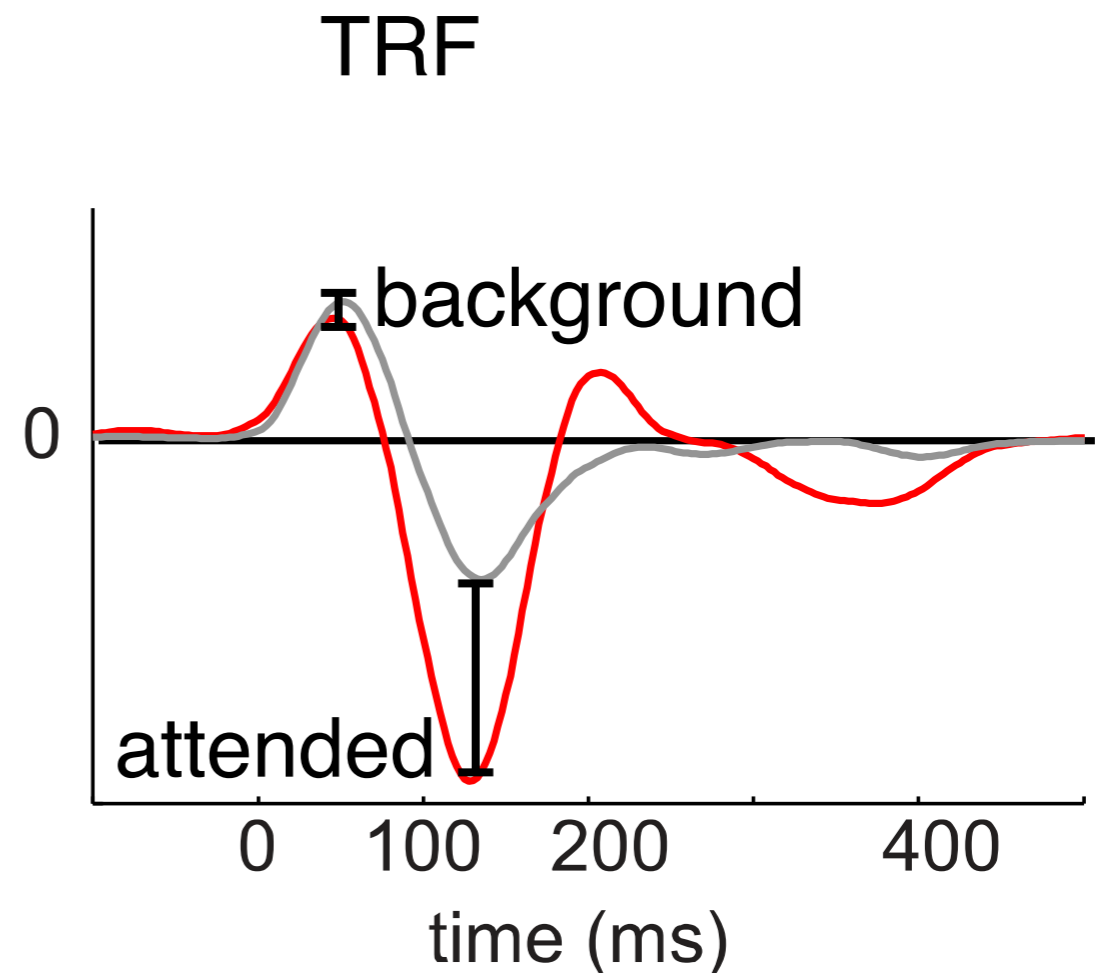


- STRF separable (time, frequency)
- 300 Hz - 2 kHz dominant carriers
- $M50_{\text{STRF}}$ positive peak
- $M100_{\text{STRF}}$ negative peak

STRF Results

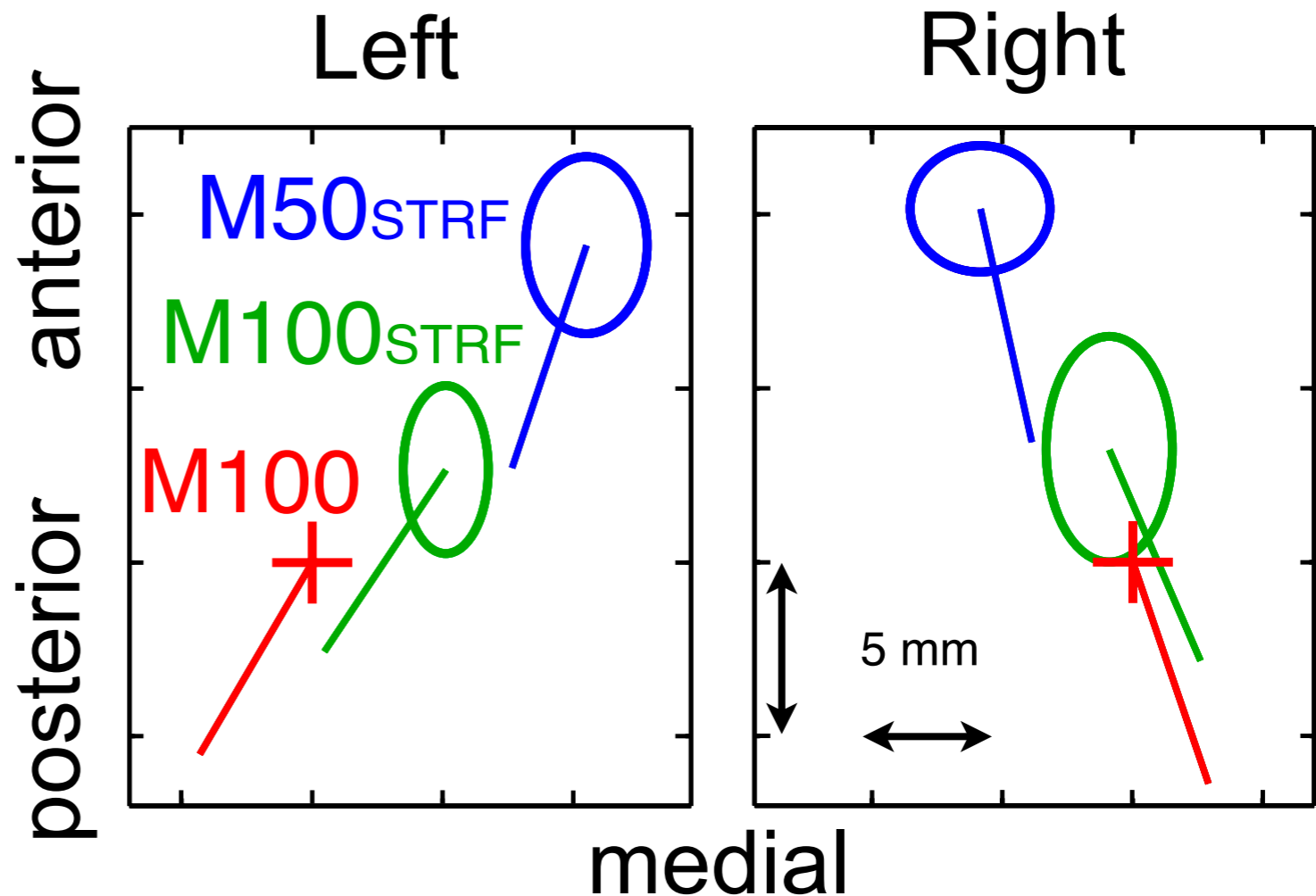


- STRF separable (time, frequency)
- 300 Hz - 2 kHz dominant carriers
- $M50_{STRF}$ positive peak
- $M100_{STRF}$ negative peak
- **$M100_{STRF}$ strongly modulated by attention, *but not* $M50_{STRF}$**

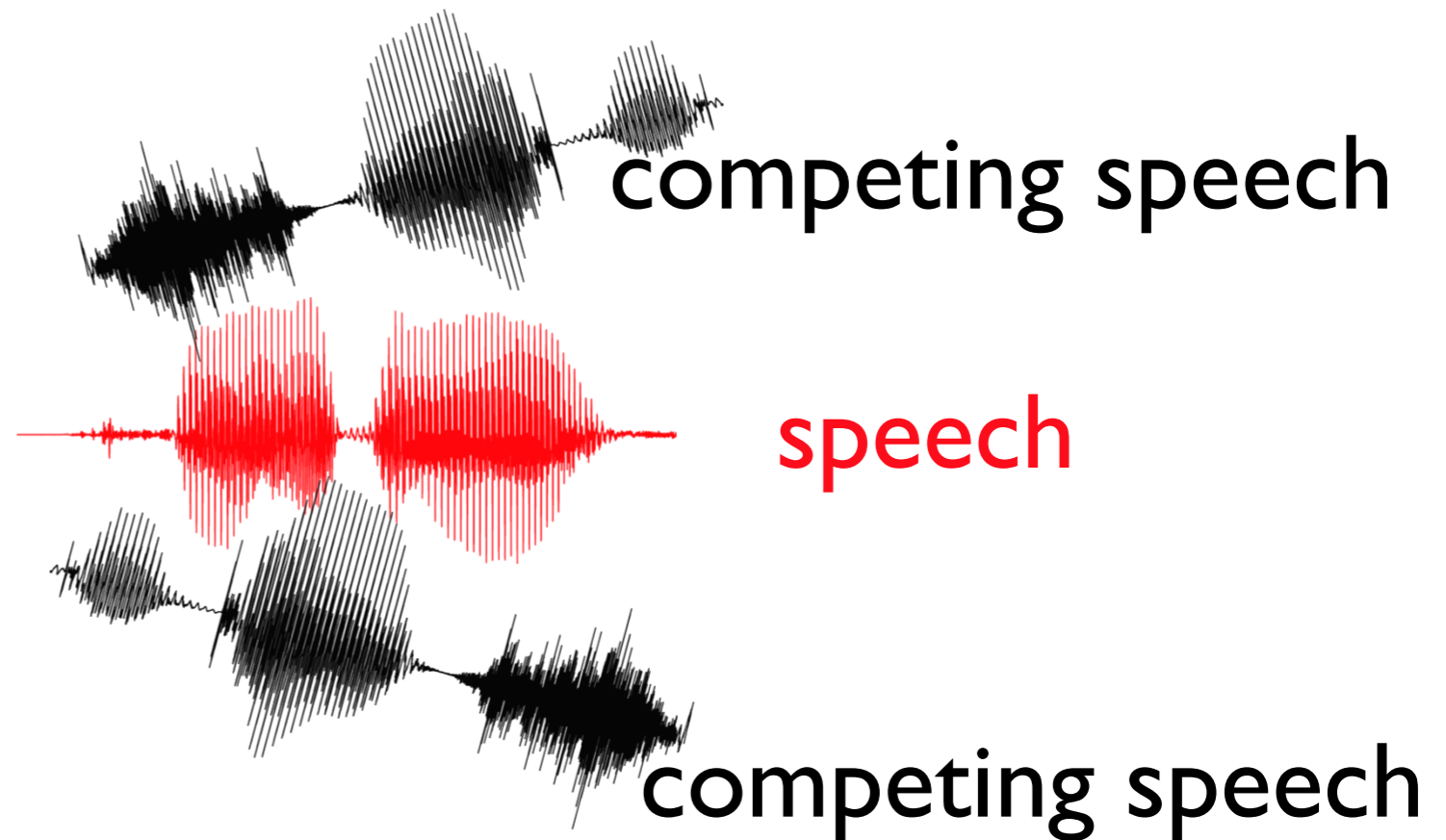


Neural Sources

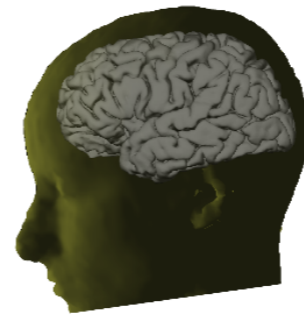
- **M100_{STRF}** source near (same as?) M100 source:
Planum Temporale
- **M50_{STRF}** source is anterior and medial to M100 (same as M50?):
Heschl's Gyrus
- **PT strongly affected by attention, *but not HG***



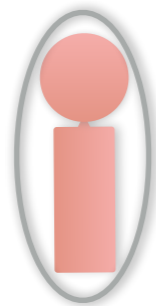
Three Competing Speakers



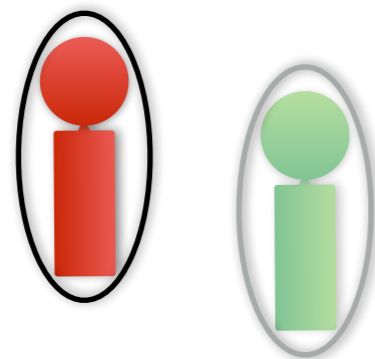
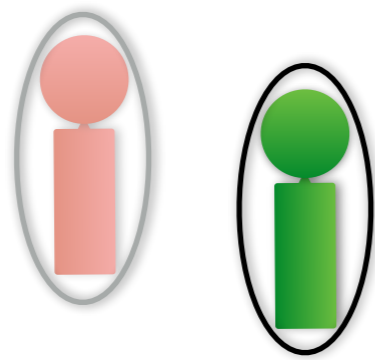
Foreground vs. Background



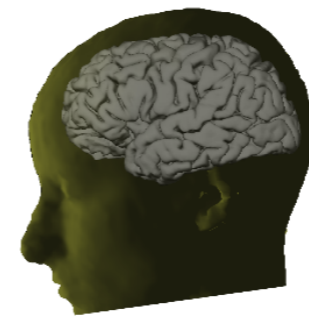
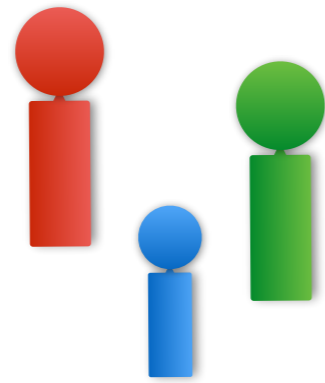
Foreground vs. Background



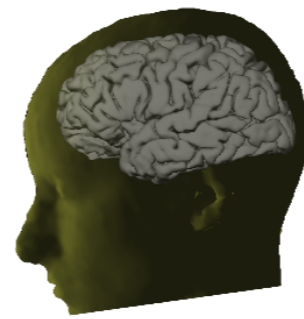
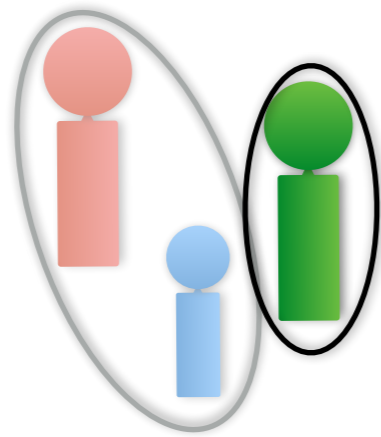
Foreground vs. Background



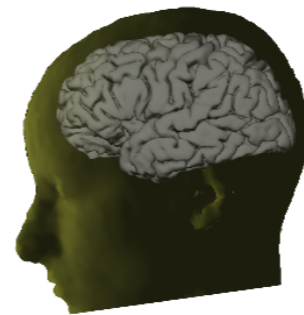
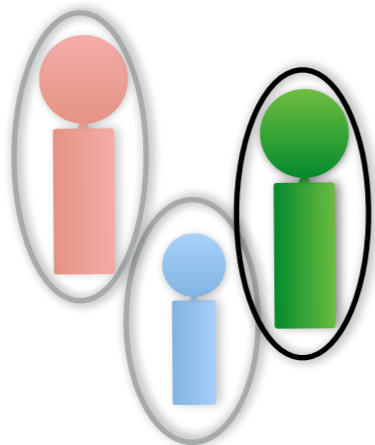
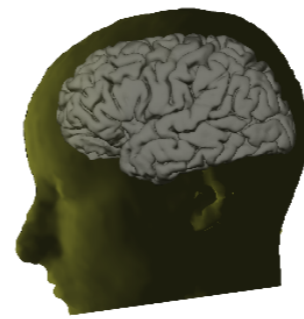
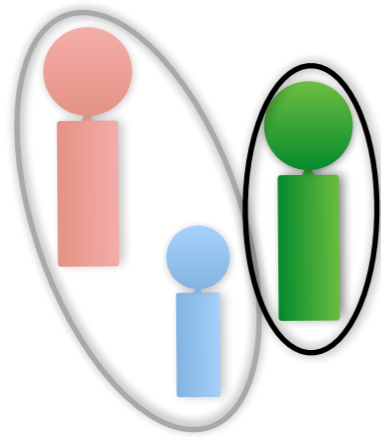
Foreground vs. Background



Foreground vs. Background

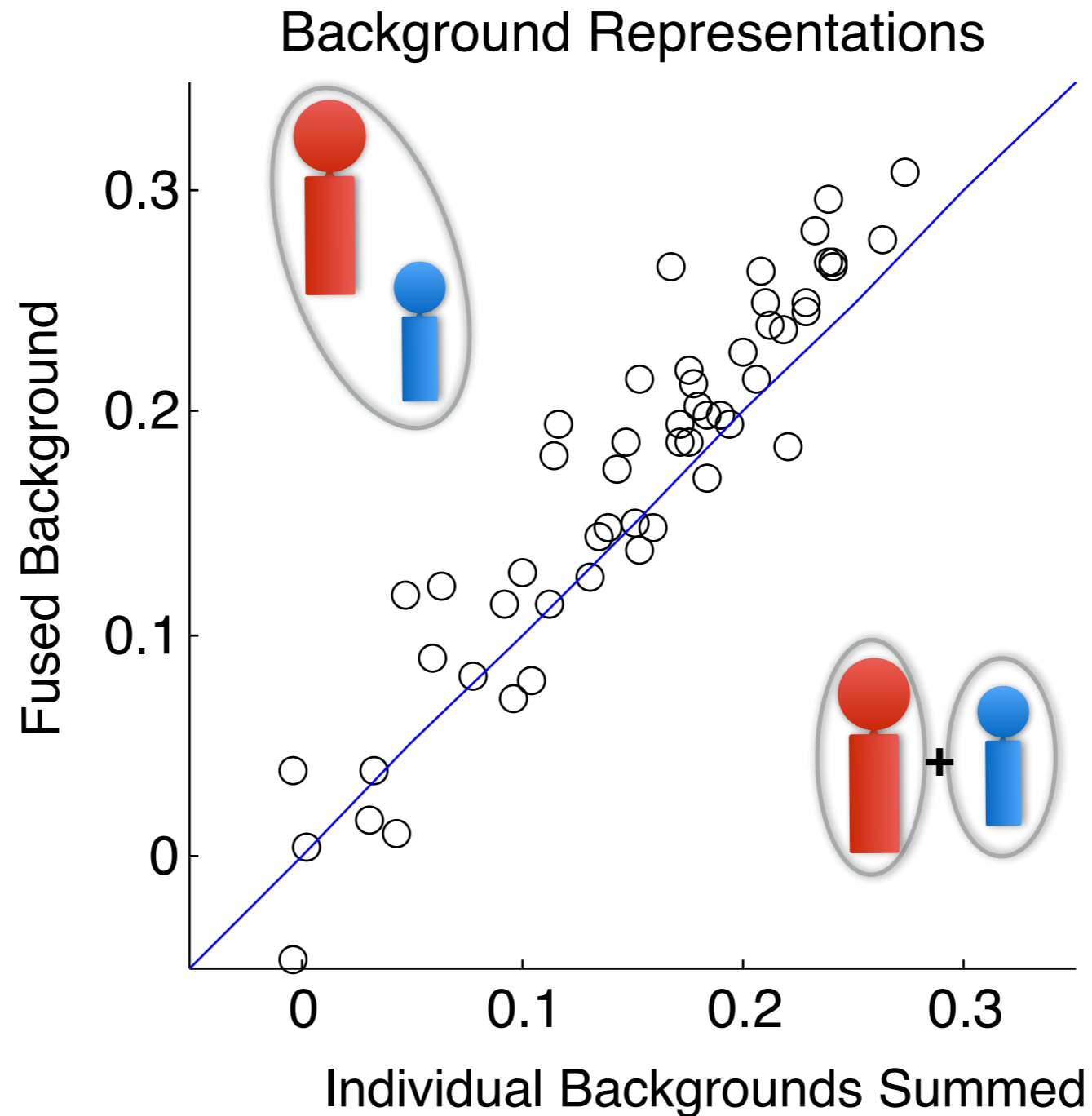


Foreground vs. Background



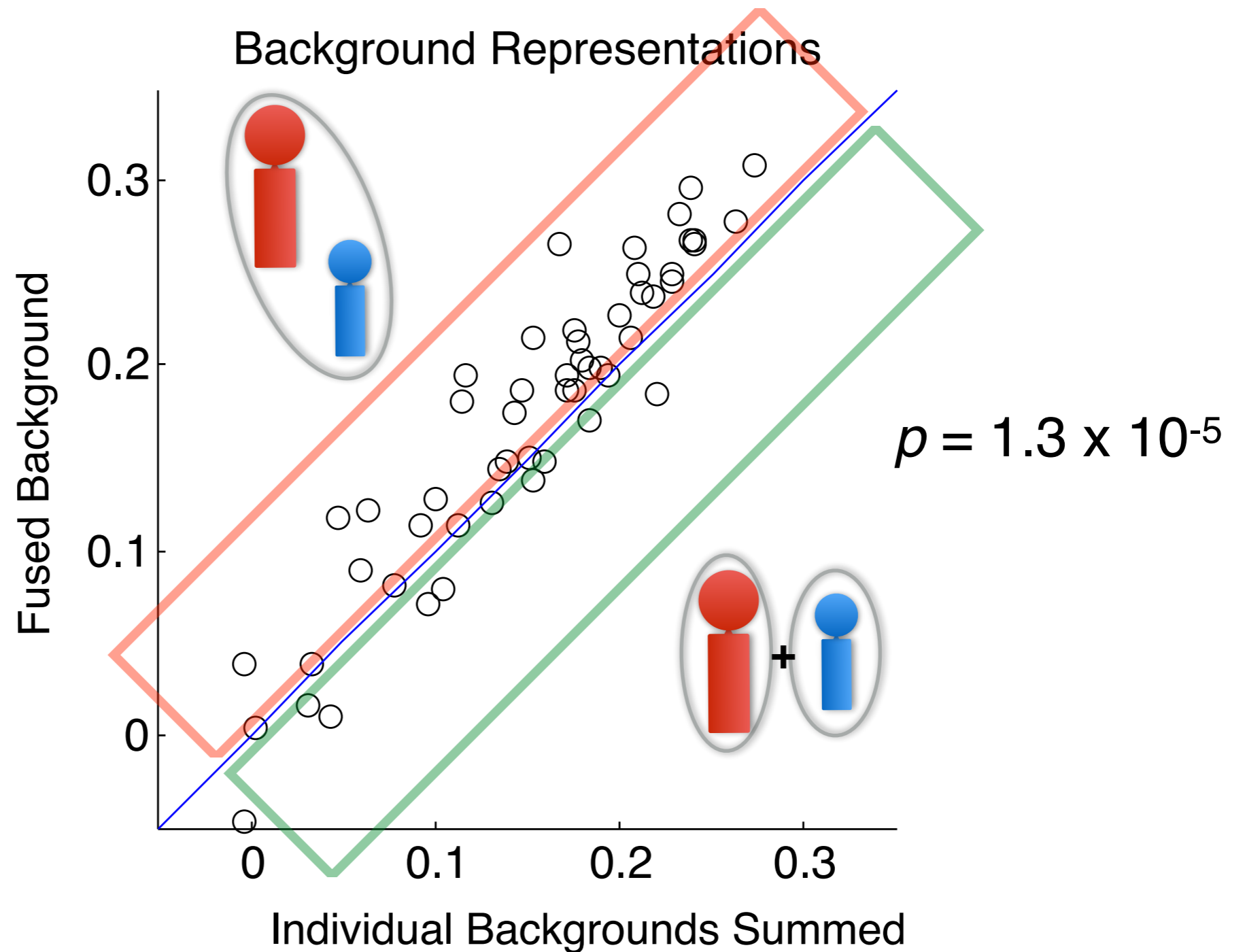
Backgrounds vs. Background

Integration Window
Late Times Only



Backgrounds vs. Background

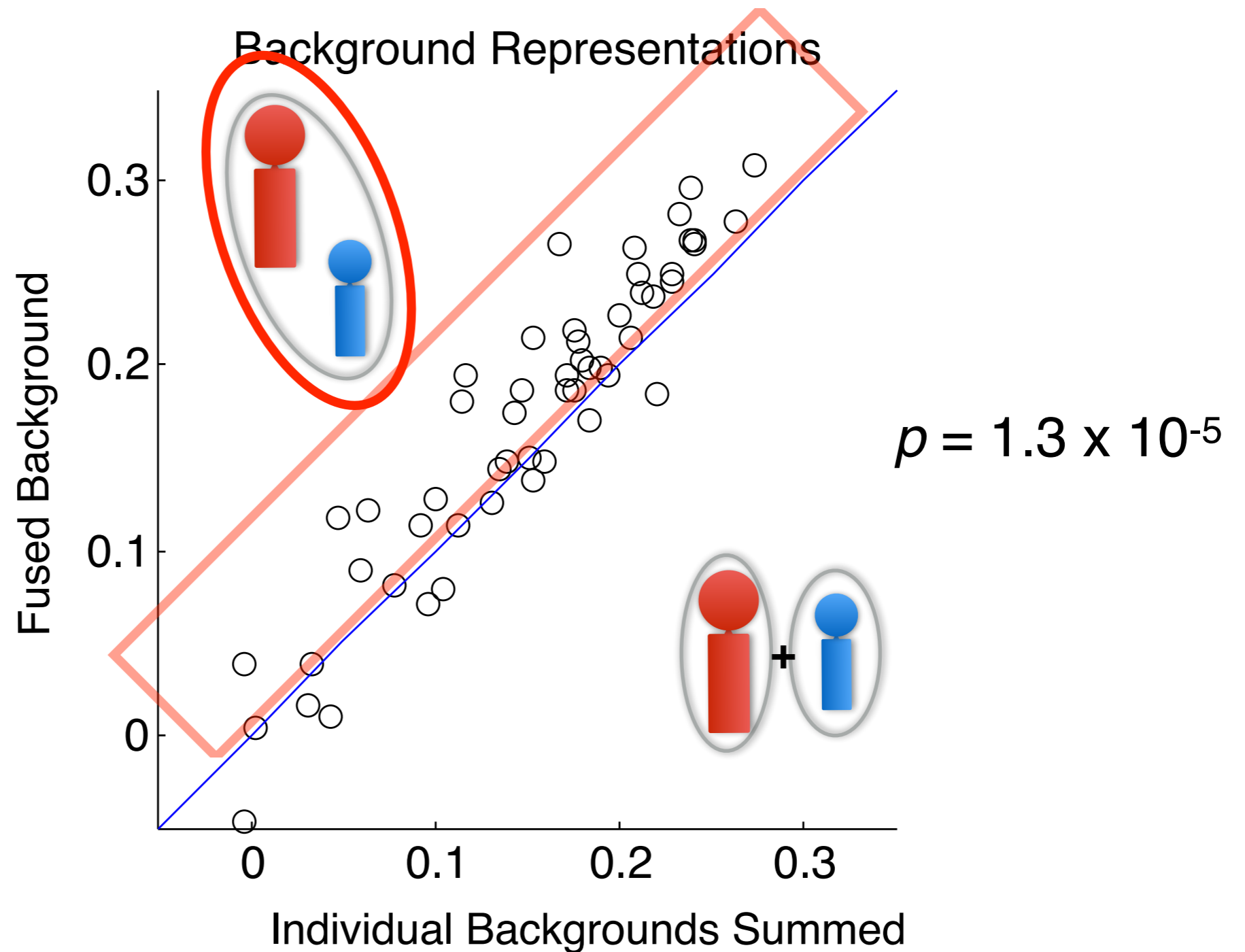
Integration Window
Late Times Only



Backgrounds vs. Background

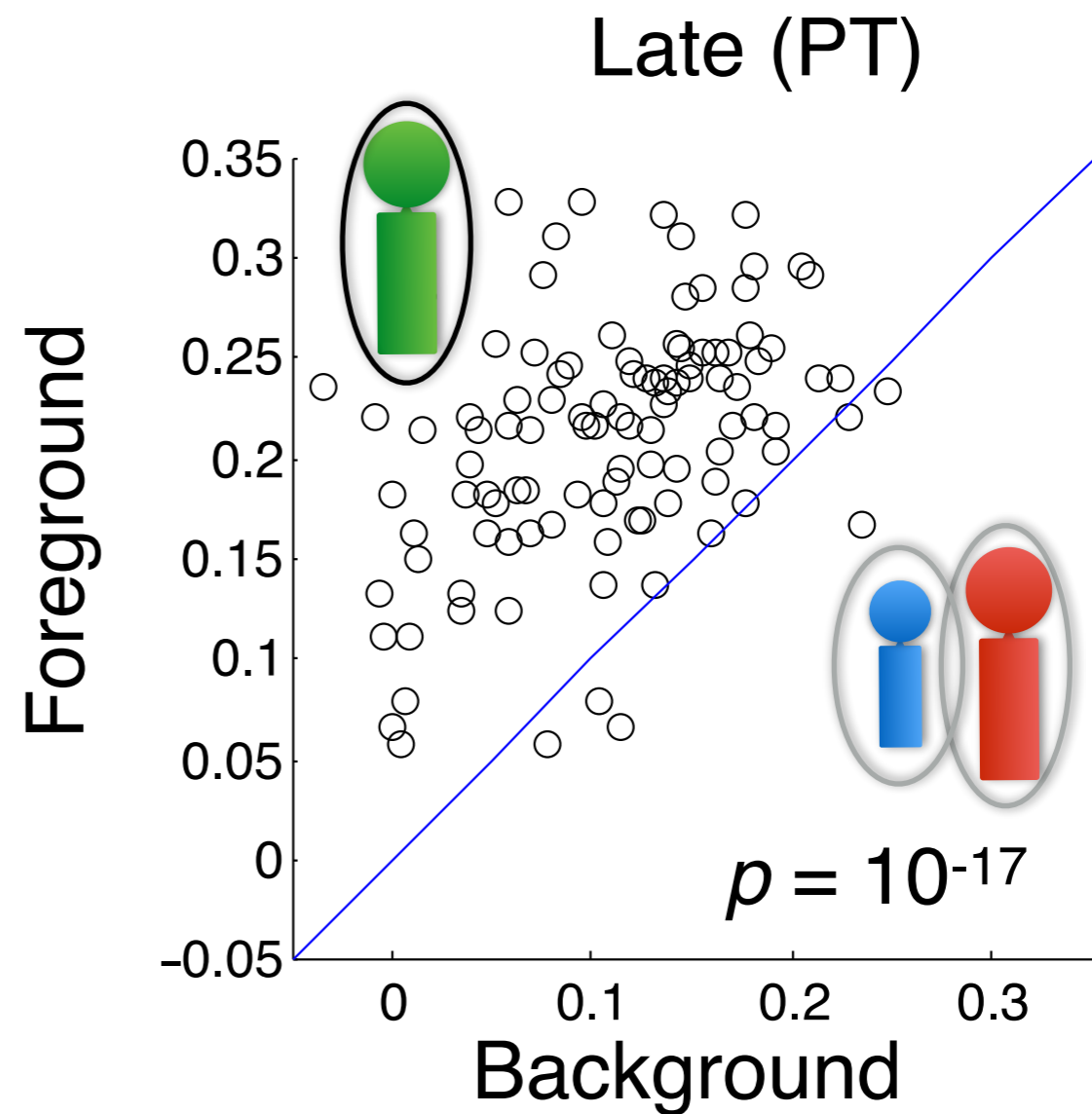
Integration Window
Late Times Only

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



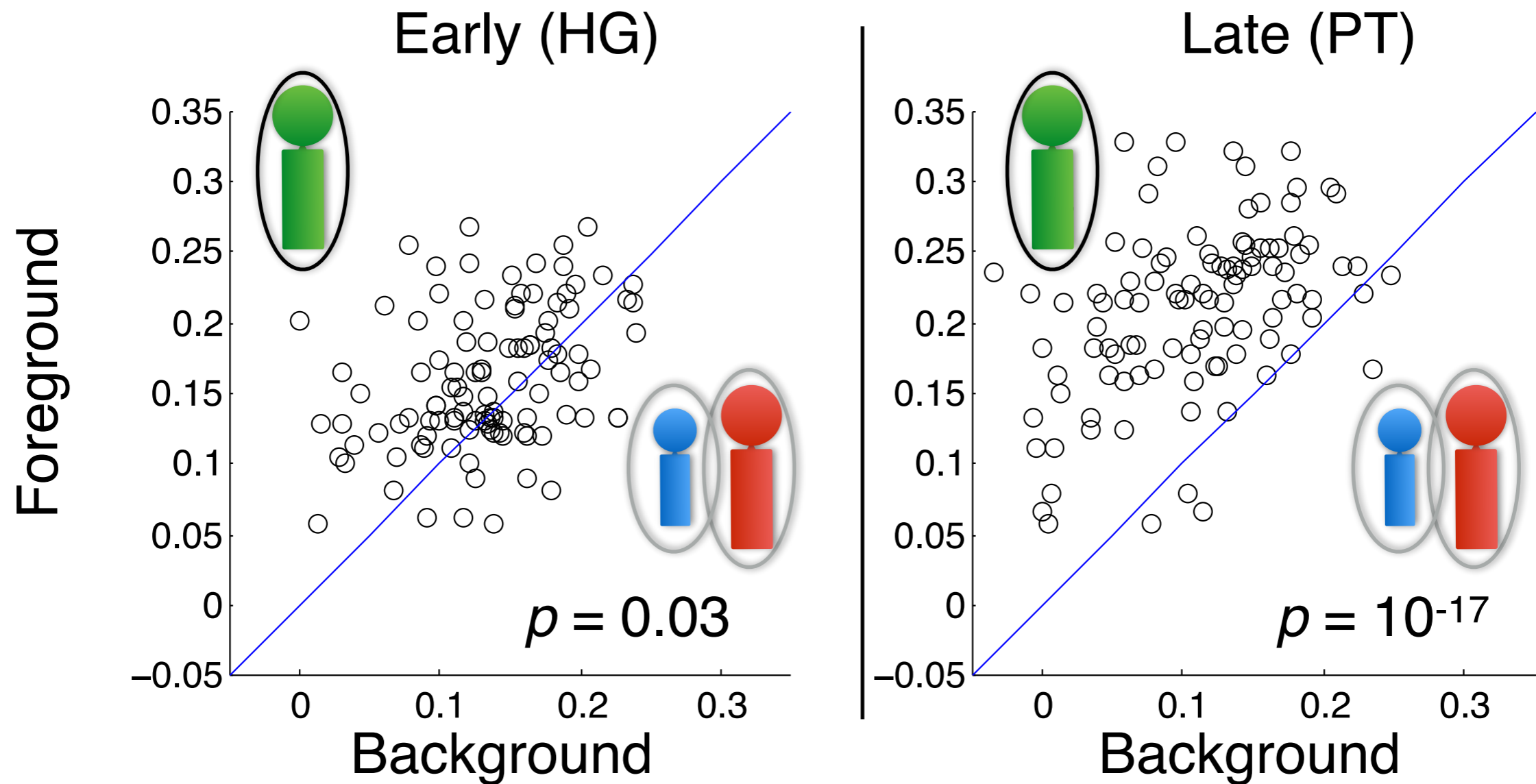
Foreground vs. Background

Early vs. Late



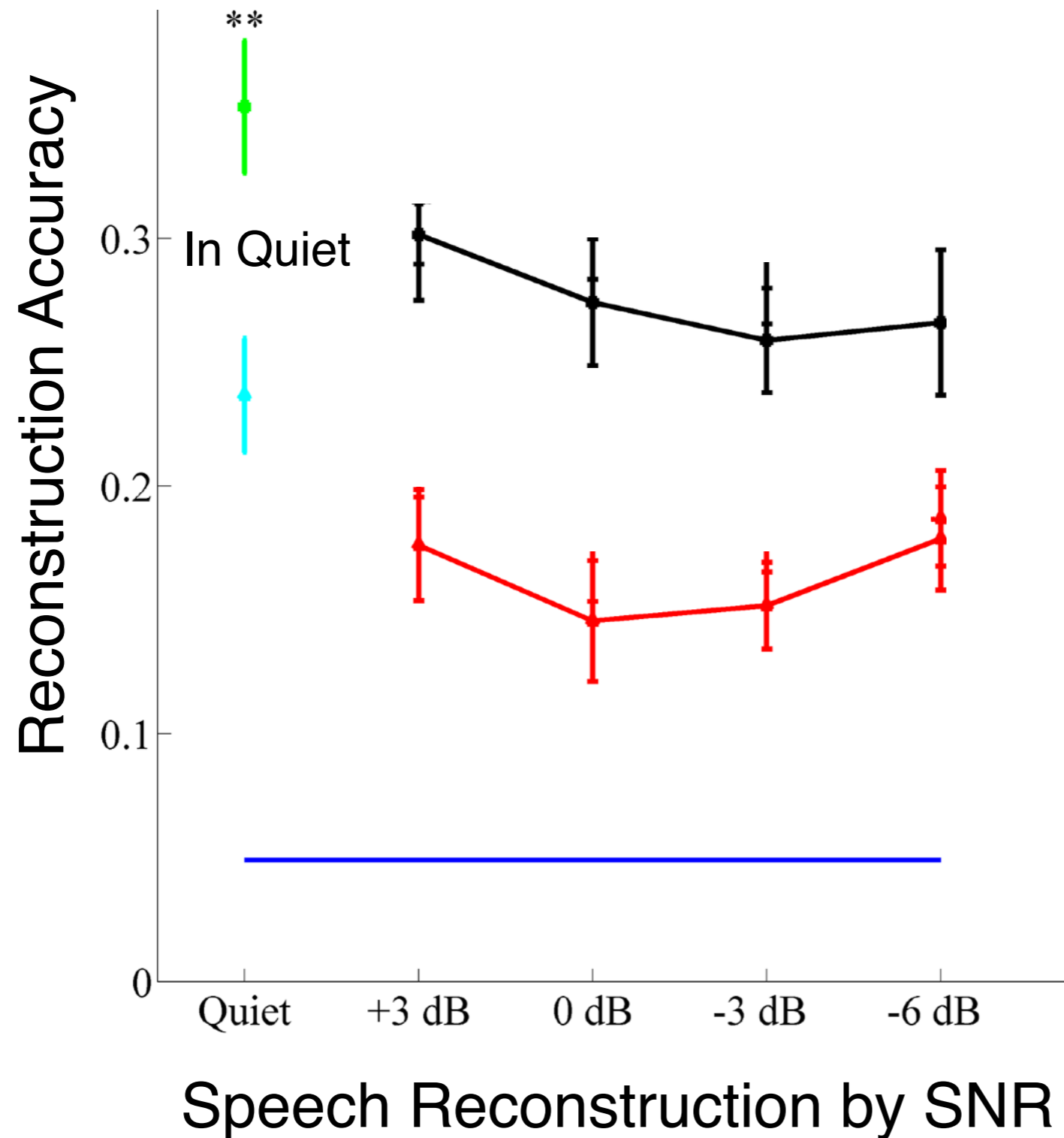
Foreground vs. Background

Early vs. Late



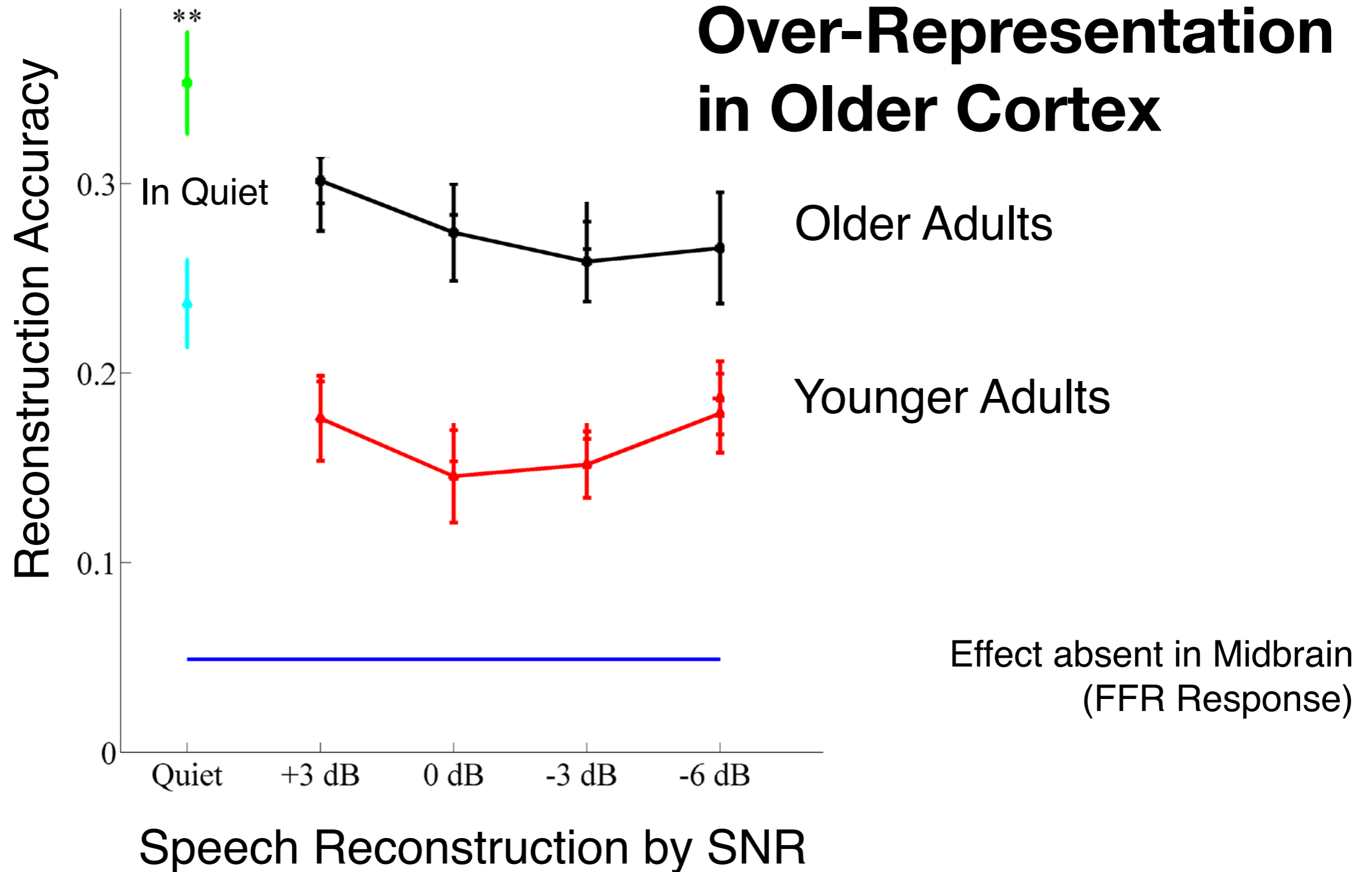
HG represents attended and unattended speech with *almost* equal fidelity

Younger vs. Older Adults

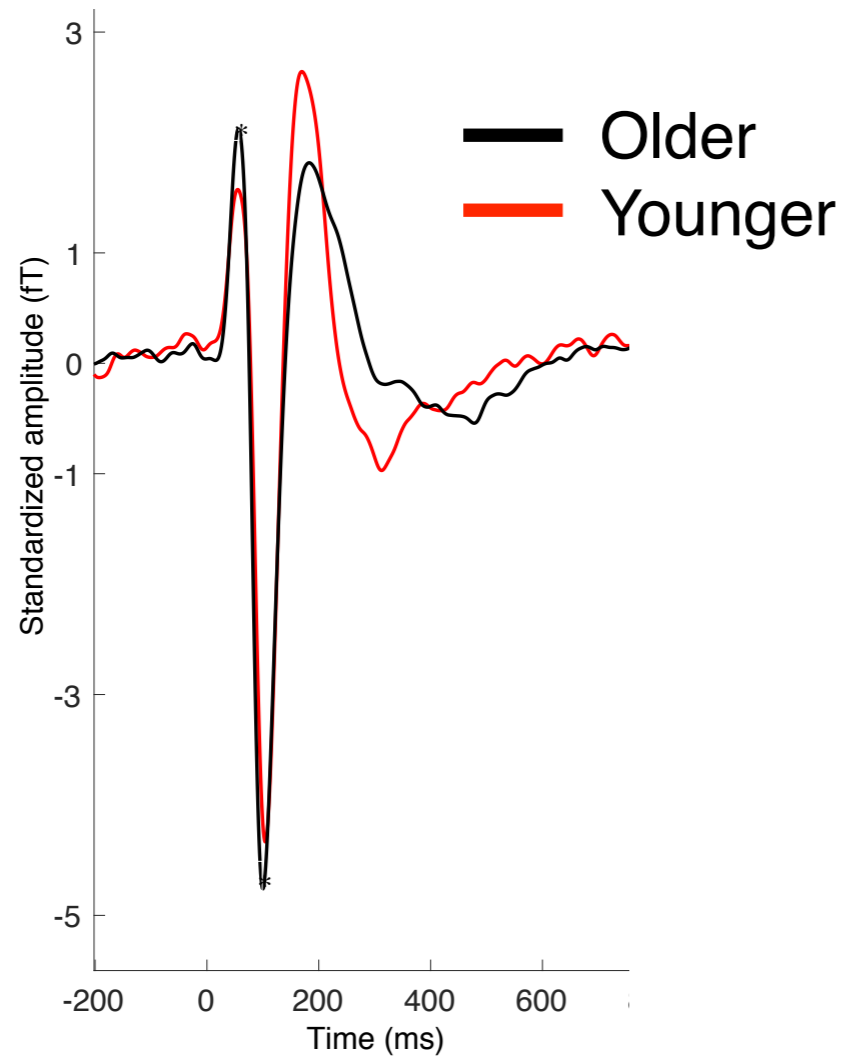


Younger vs. Older Adults

Over-Representation in Older Cortex

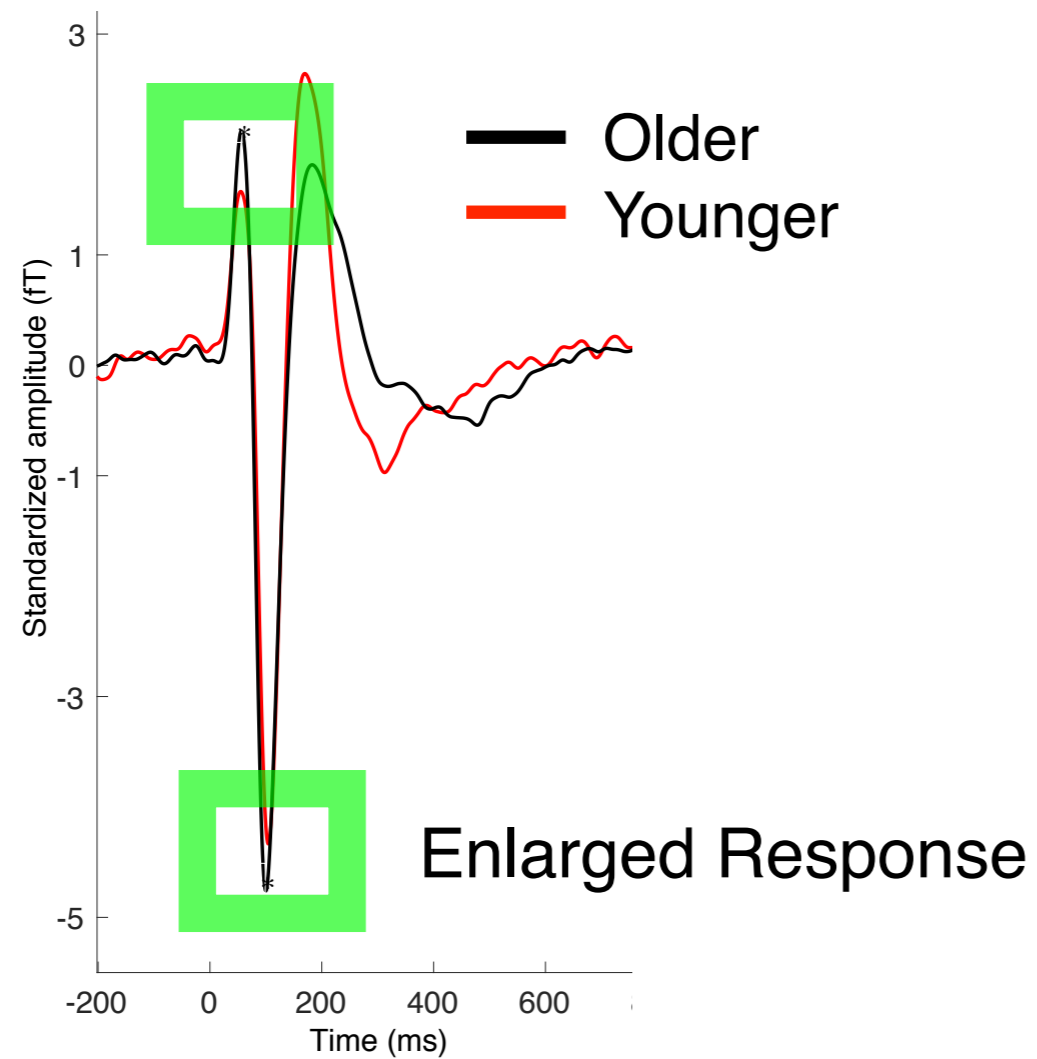


Older Enlarged Response



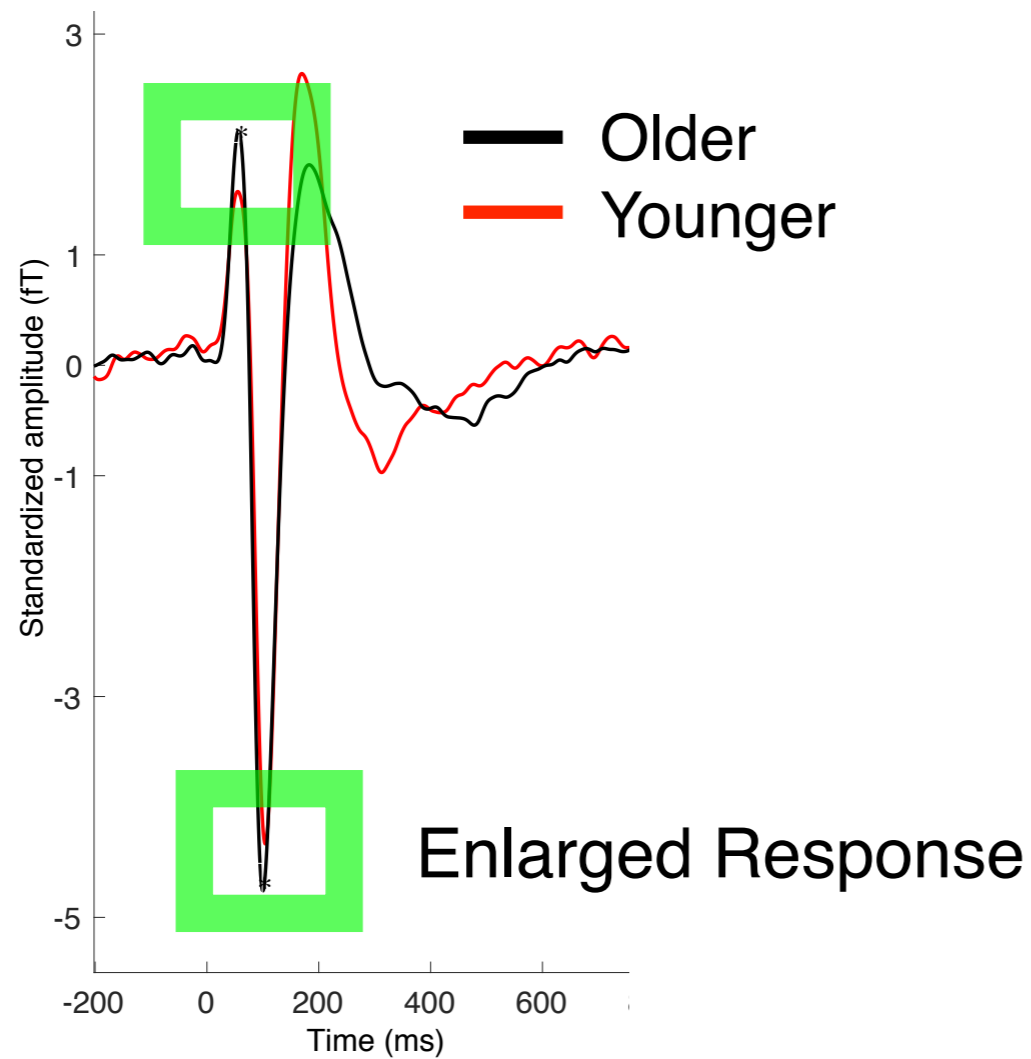
Average Responses to Pure Tone

Older Enlarged Response

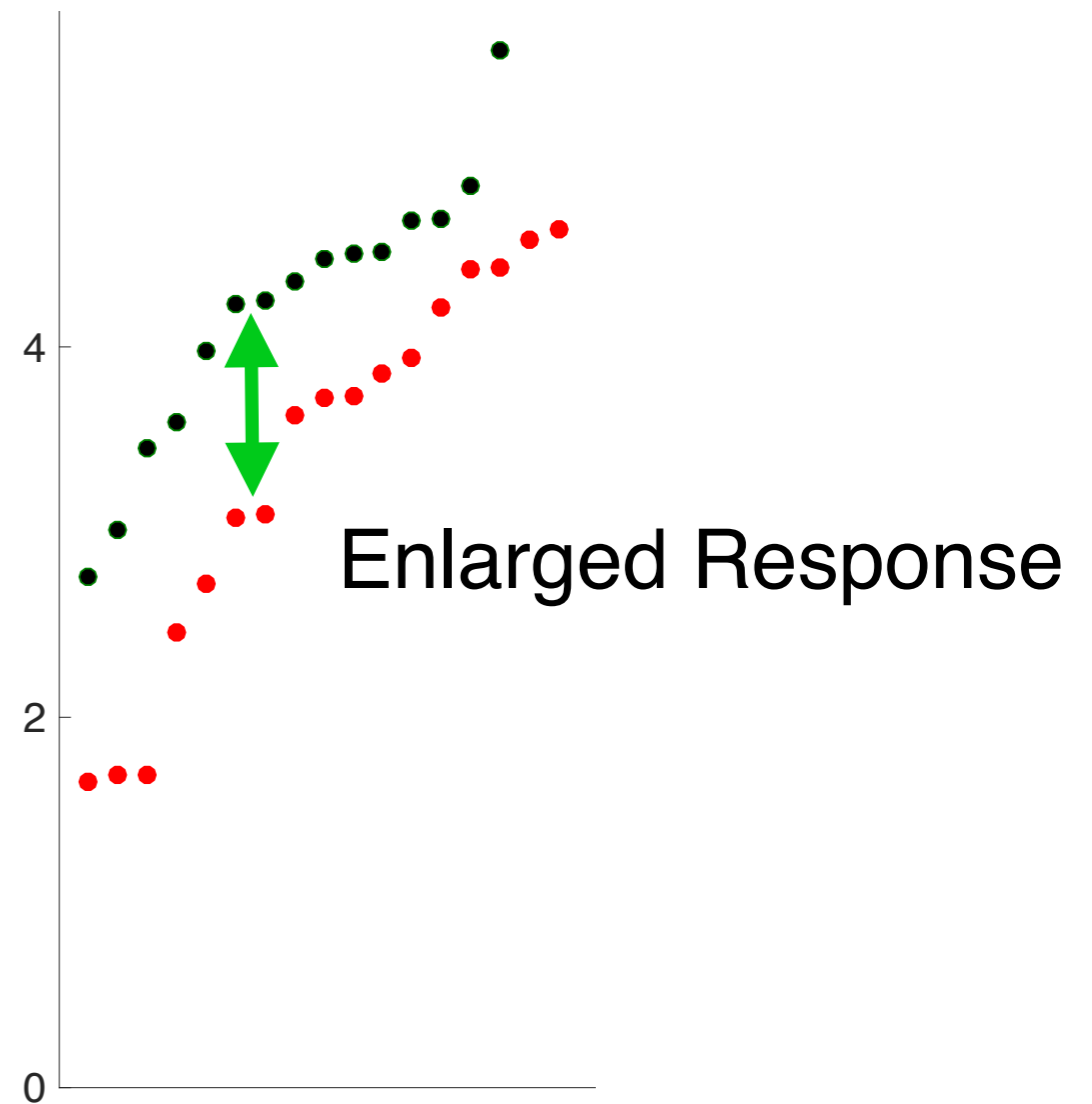


Average Responses to Pure Tone

Older Enlarged Response



Average Responses to Pure Tone

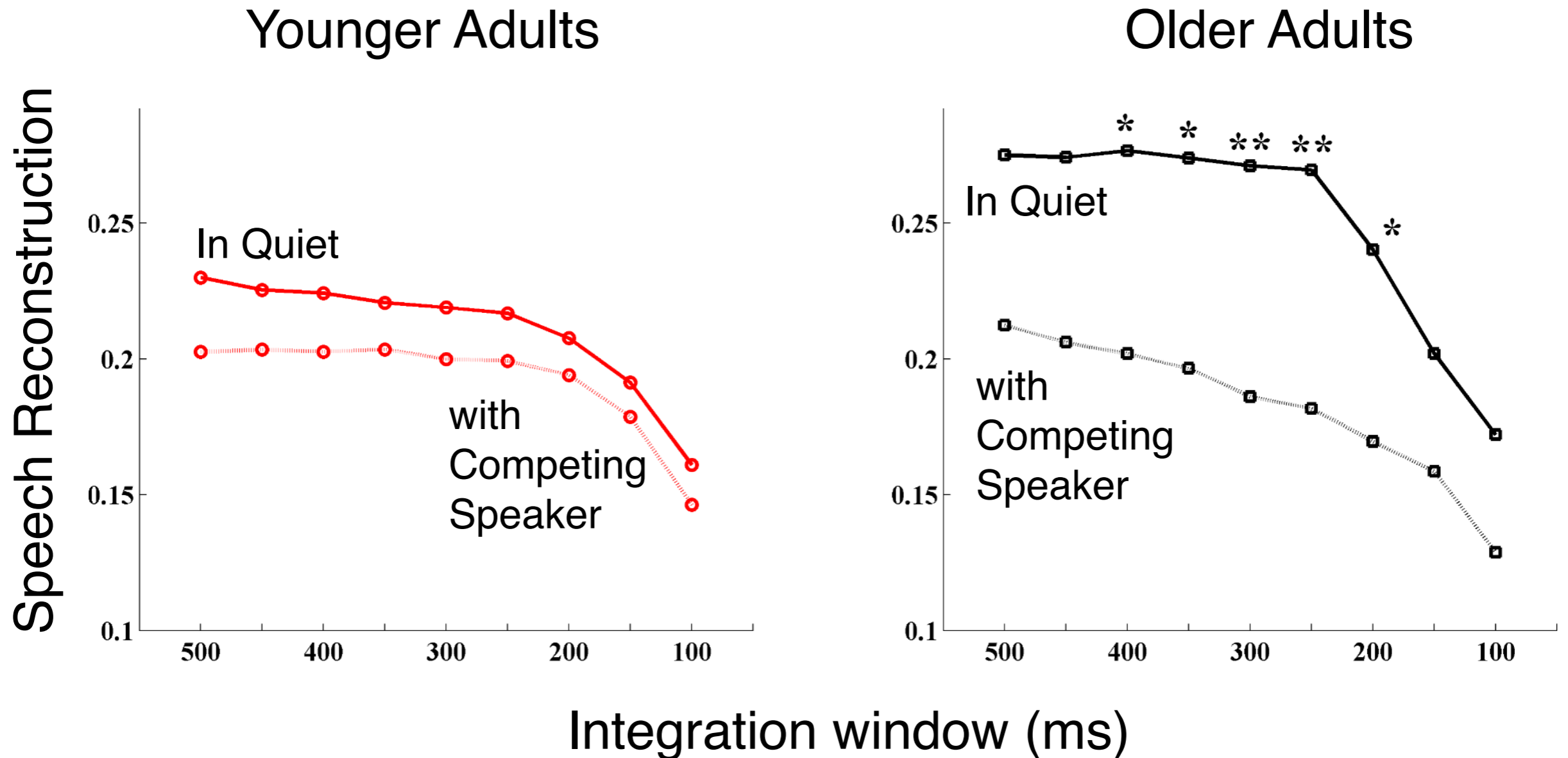


M100 Power by Subject

Also true for pure tones

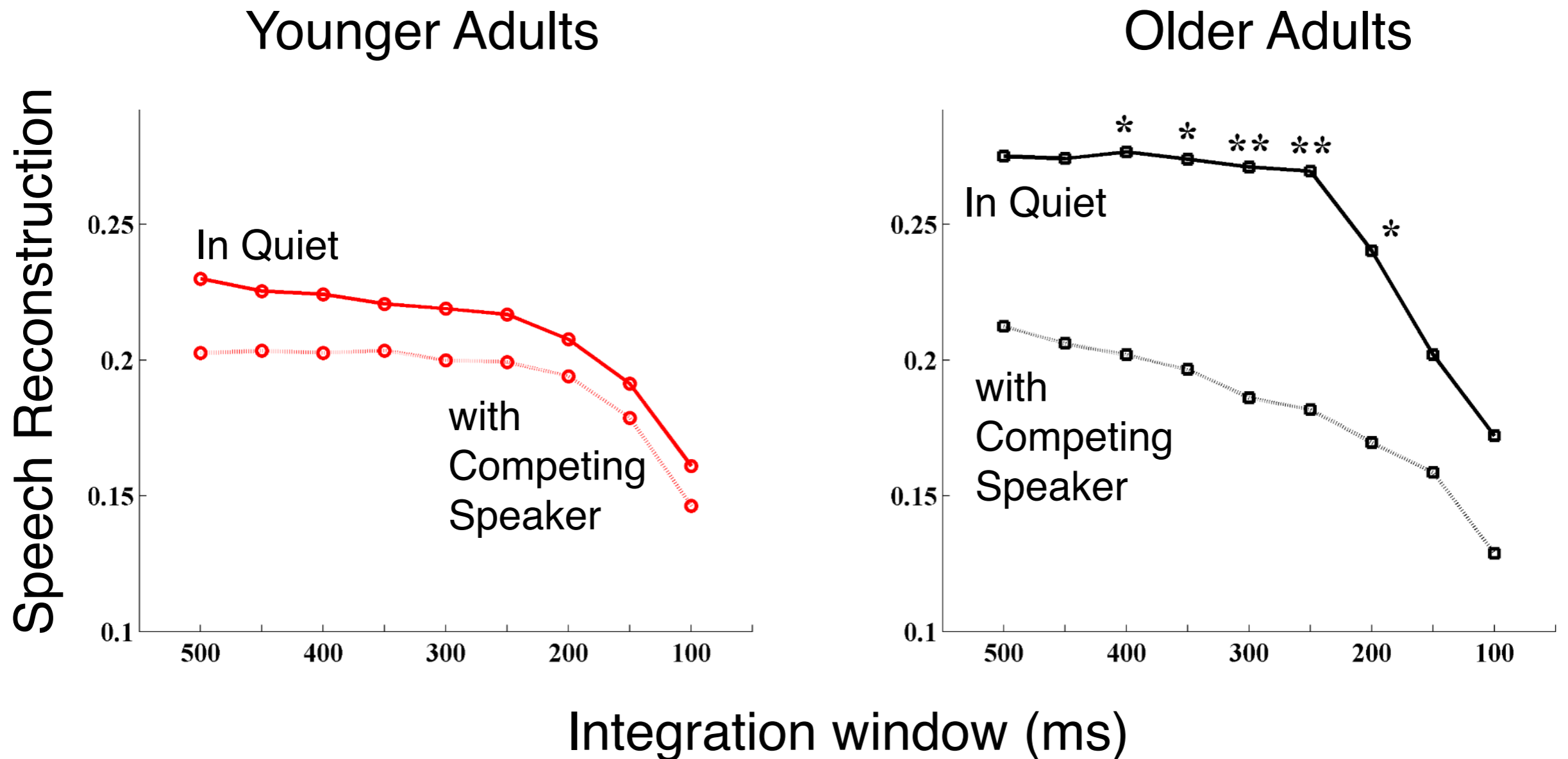
Younger vs. Older Adults

Temporal Integration

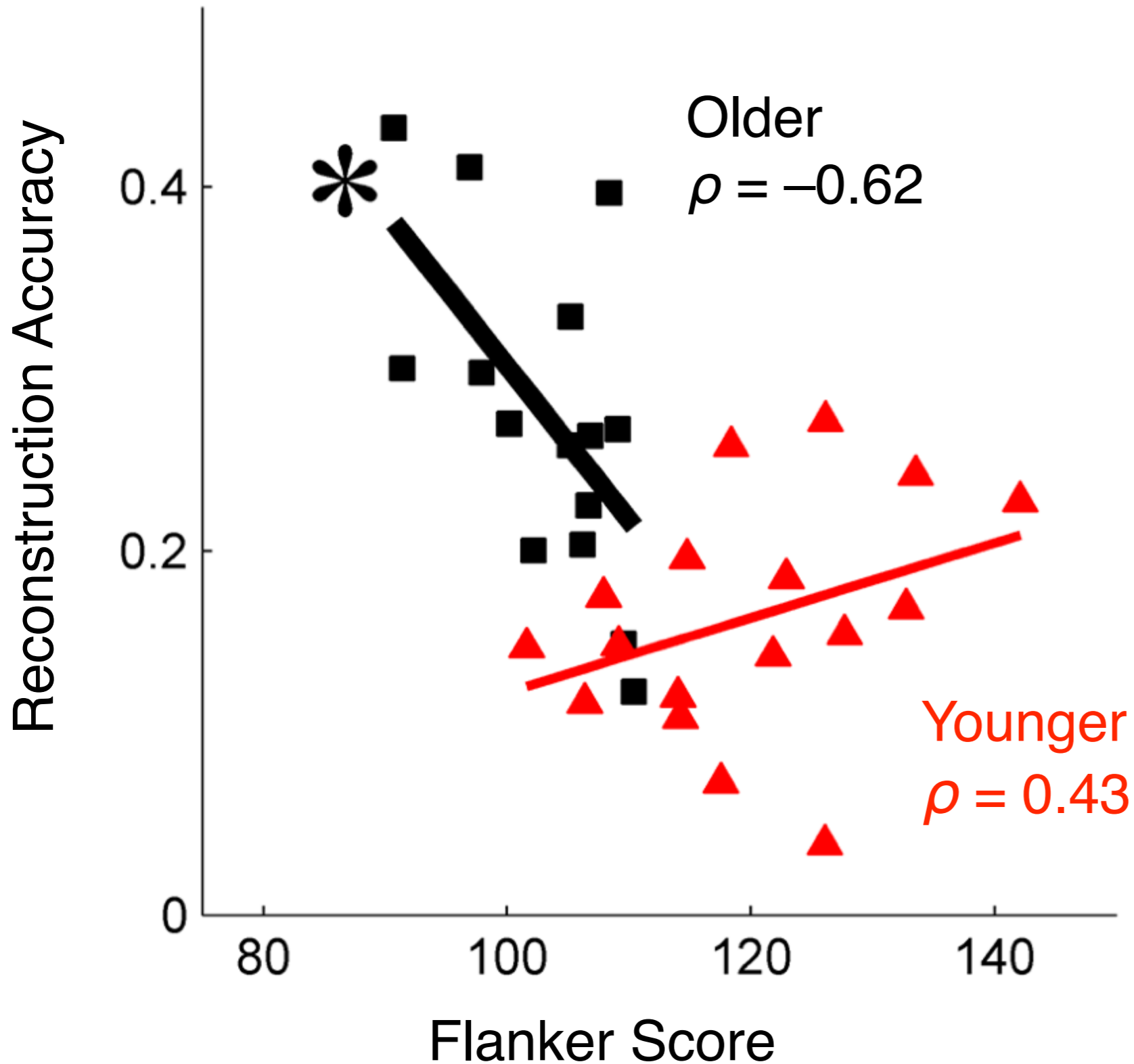


Younger vs. Older Adults

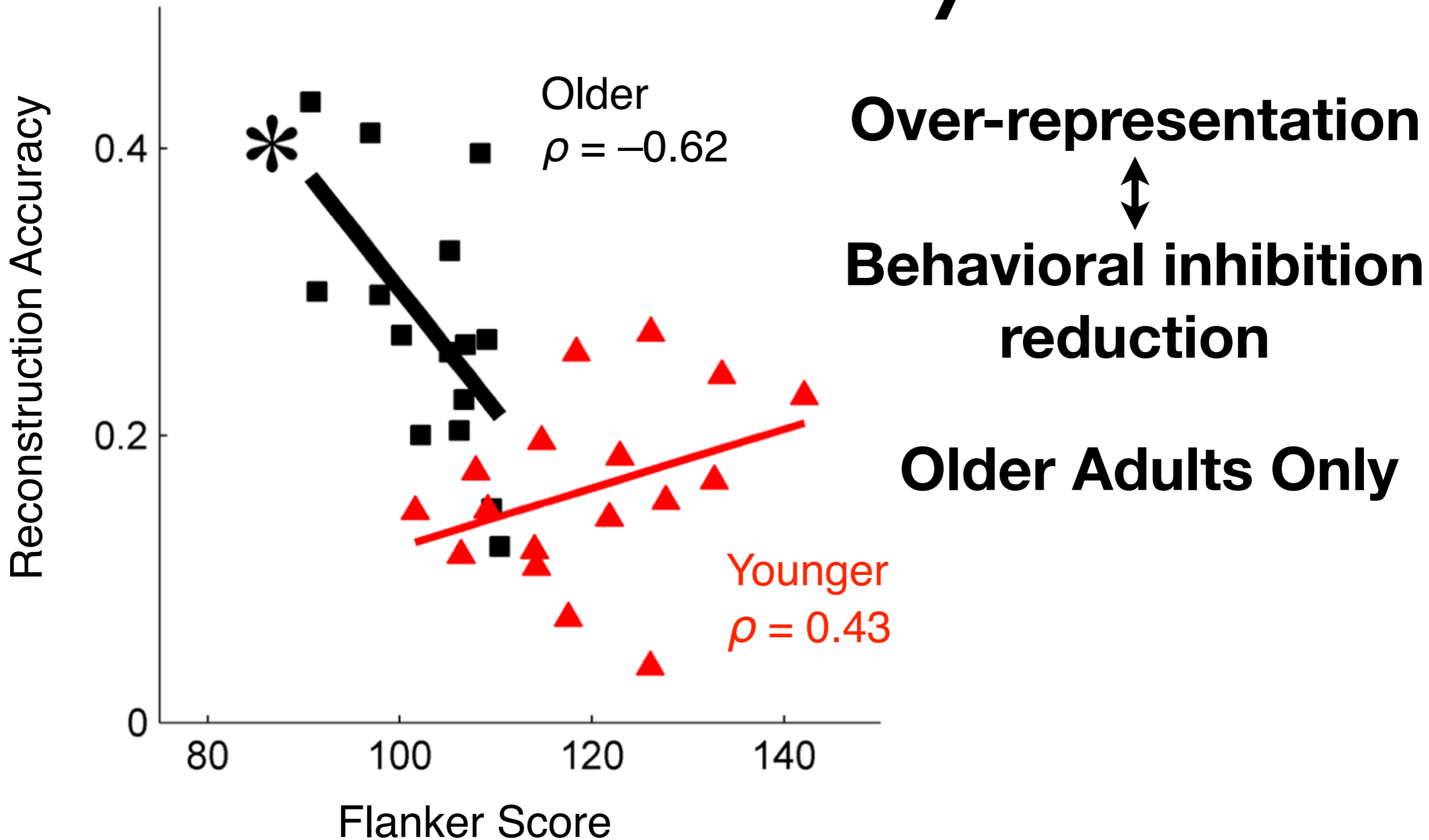
Temporal Integration



Neural vs Inhibitory Control



Neural vs Inhibitory Control

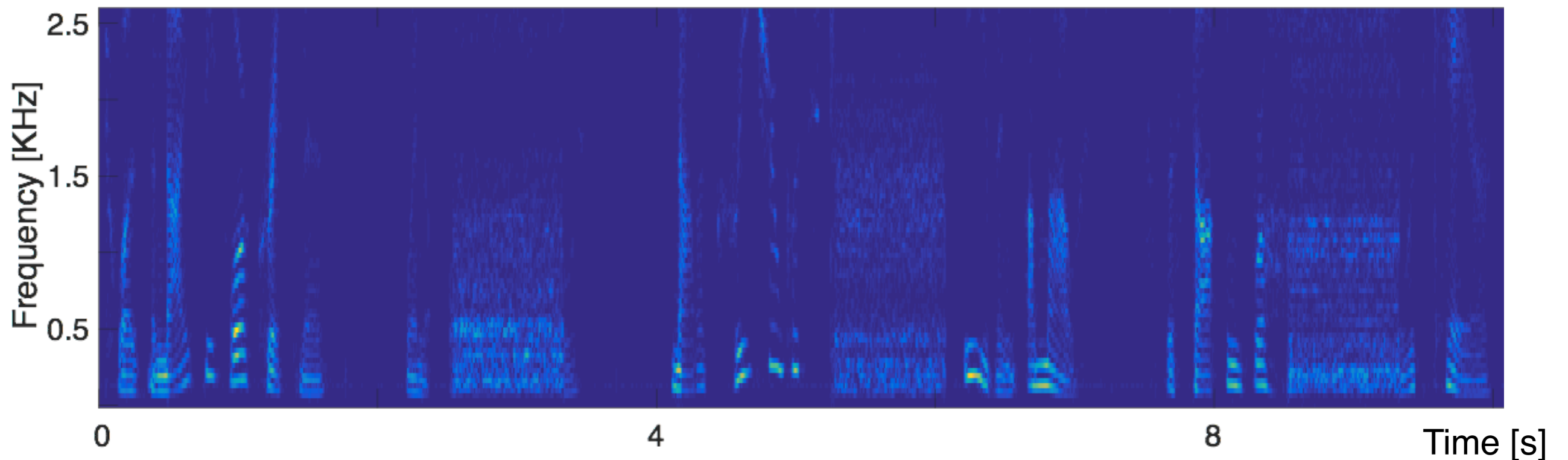
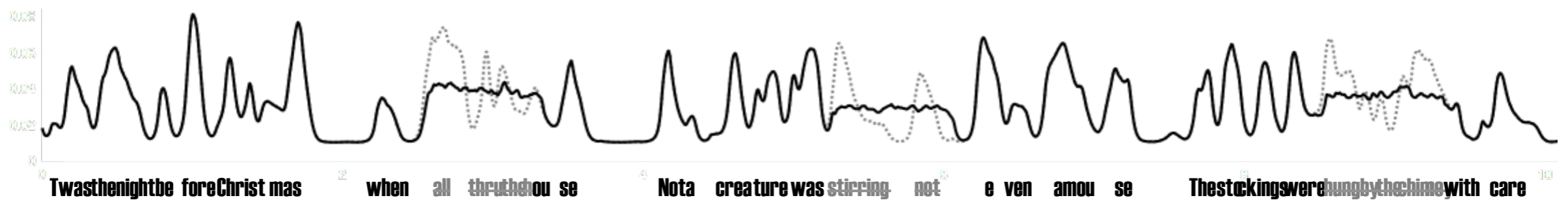


Missing Speech Restoration

- Can sustained, strongly non-stationary, speech be “restored”?
 - ▶ Might be aided by contextual knowledge/familiarity
 - ▶ Might be aided by strong rhythmicity

Missing Speech Restoration

- Can sustained, strongly non-stationary, speech be “restored”?
 - ▶ Might be aided by contextual knowledge/familiarity
 - ▶ Might be aided by strong rhythmicity



Missing Speech: Context

Twas the night before Christmas, when all through the house
not a creature was stirring, not even a mouse.
The stockings were hung by the chimney with care,
in hopes that St. Nicholas soon would be there.

The children were nestled all snug in their beds,
while visions of sugar plums danced in their heads.
And Mama in her 'kerchief, and I in my cap,
had just settled our brains for a long winter's nap.

When out on the lawn there arose such a clatter,
I sprang from my bed to see what was the matter.
Away to the window I flew like a flash,
tore open the shutter, and threw up the sash.

The moon on the breast of the new-fallen snow
gave the lustre of midday to objects below,
when, what to my wondering eyes should appear,
but a miniature sleigh and eight tiny reindeer.

With a little old driver, so lively and quick,
I knew in a moment it must be St. Nick.
More rapid than eagles, his coursers they came,
and he whistled and shouted and called them by name.

"Now Dasher! Now Dancer! Now, Prancer and Vixen!
On, Comet! On, Cupid! On, Donner and Blitzen!
To the top of the porch! To the top of the wall!
Now dash away! Dash away! Dash away all!"

As dry leaves that before the wild hurricane fly,
when they meet with an obstacle, mount to the sky
so up to the house-top the coursers they flew,
with the sleigh full of toys, and St. Nicholas too.

And then, in a twinkling, I heard on the roof
the prancing and pawing of each little hoof.
As I drew in my head and was turning around,
down the chimney St. Nicholas came with a bound.

He was dressed all in fur, from his head to his foot,
and his clothes were all tarnished with ashes and soot.
A bundle of toys he had flung on his back,
and he looked like a peddler just opening his pack.

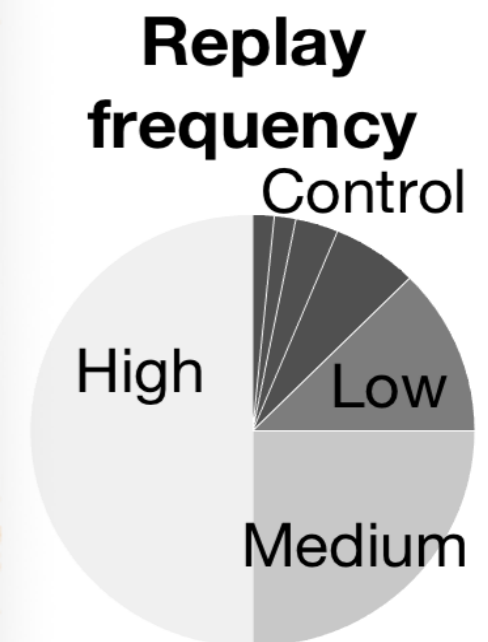
His eyes--how they twinkled! His dimples, how merry!
His cheeks were like roses, his nose like a cherry!
His droll little mouth was drawn up like a bow,
and the beard on his chin was as white as the snow.

The stump of a pipe he held tight in his teeth,
and the smoke it encircled his head like a wreath.
He had a broad face and a little round belly,
that shook when he laughed, like a bowl full of jelly.

He was chubby and plump, a right jolly old elf,
and I laughed when I saw him, in spite of myself.
A wink of his eye and a twist of his head
soon gave me to know I had nothing to dread.

He spoke not a word, but went straight to his work,
and filled all the stockings, then turned with a jerk.
And laying his finger aside of his nose,
and giving a nod, up the chimney he rose.

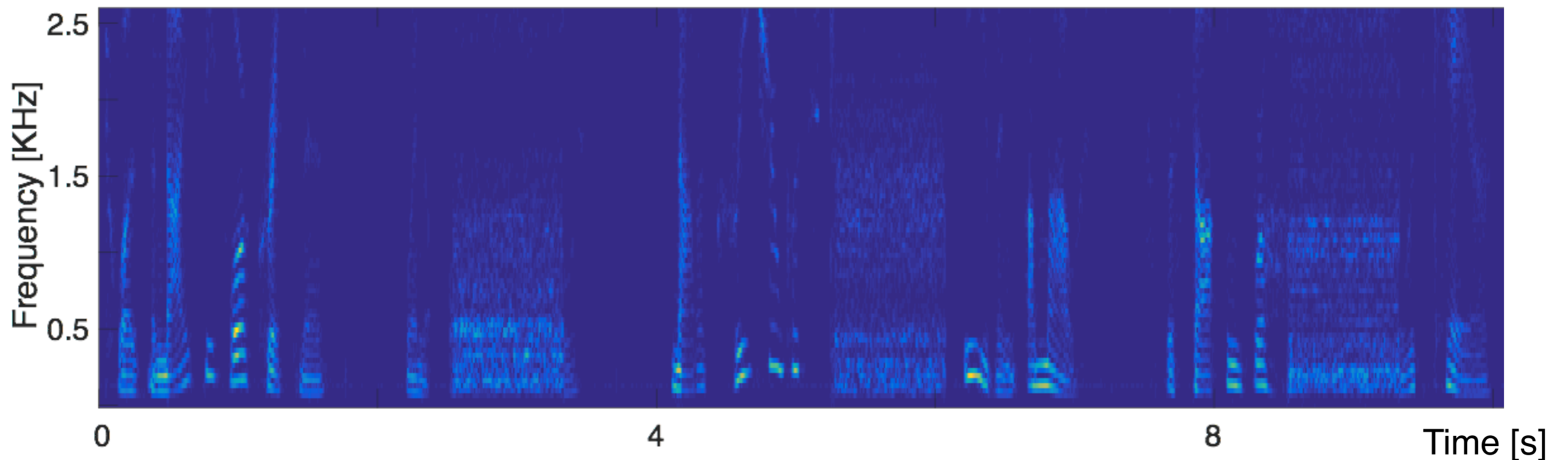
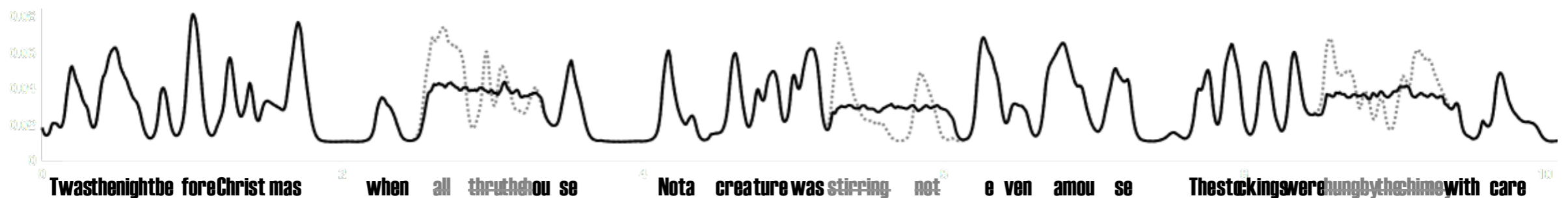
He sprang to his sleigh, to his team gave a whistle,
And away they all flew like the down of a thistle.
But I heard him exclaim, 'ere he drove out of sight,
"Happy Christmas to all, and to all a good night!"



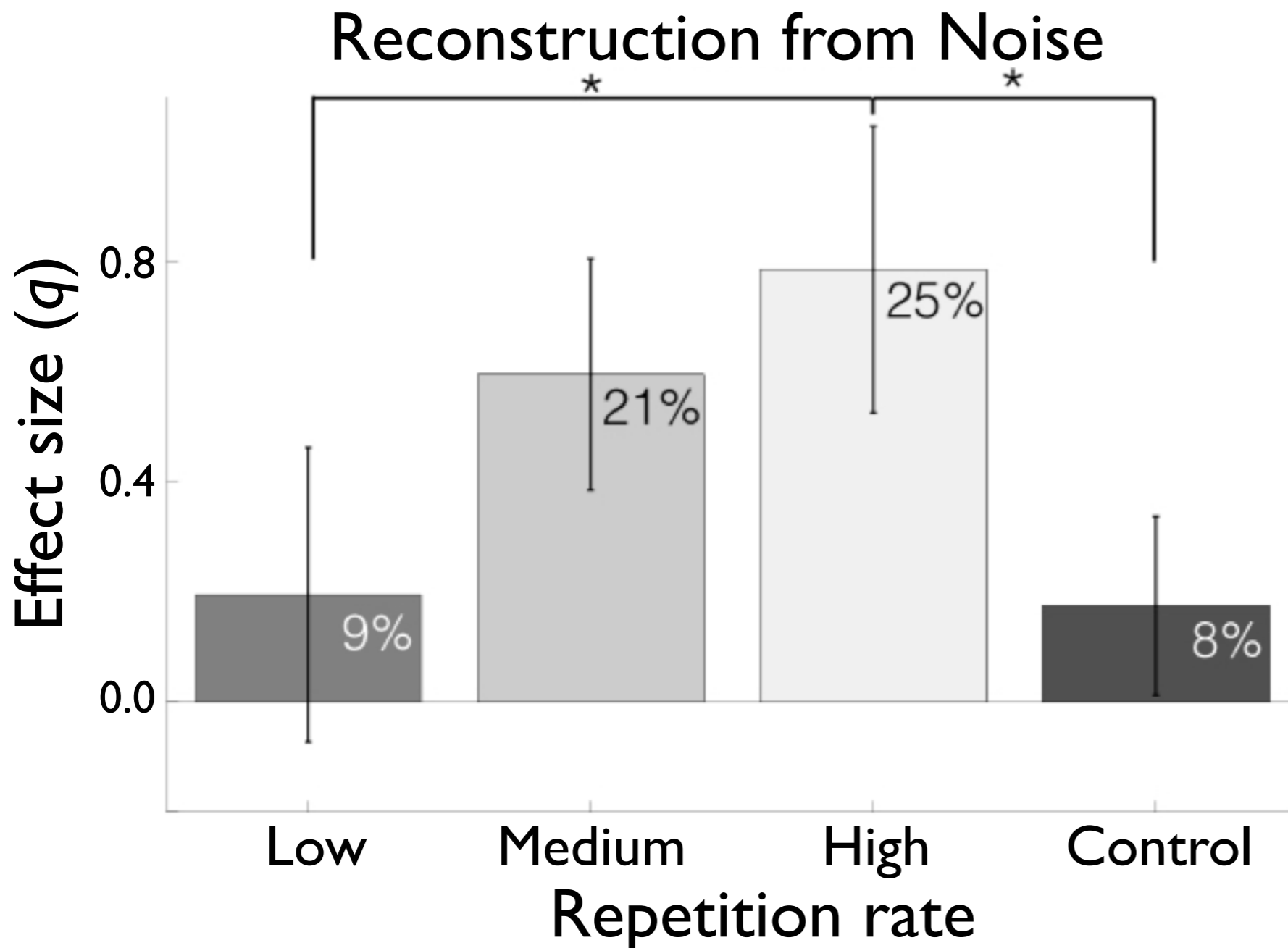
- Hypothesis: contextual knowledge of missing speech can be controlled by exposure to the speech

Missing Speech Restoration

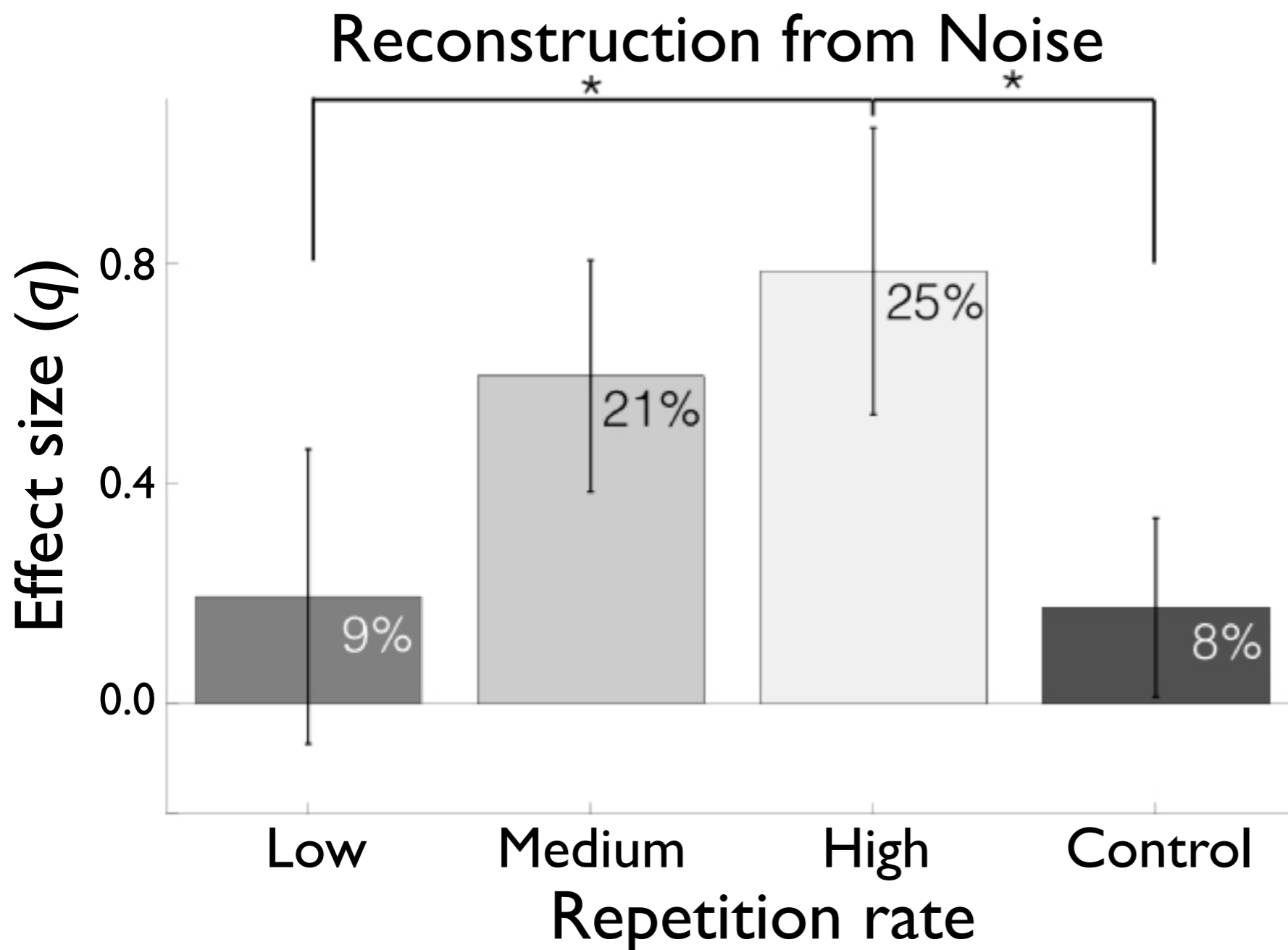
- Can sustained, strongly non-stationary, speech be “restored”?
 - ▶ Might be aided by contextual knowledge/familiarity
 - ▶ Might be aided by strong rhythmicity



Missing Speech “Reconstruction”



Missing Speech “Reconstruction”



- Decoding of the **missing** speech improves with prior experience
- Performance is a considerable fraction of that for clean speech

Summary

- Cortical representations of speech
 - representation of envelope (up to ~ 10 Hz)
 - robust against a variety of noise types
 - robust against competing speech!
- Object-based representation at 100 ms latency (PT), but not by 50 ms (HG)
- Aging shows over-representation (and time integration deficits)
- Applies to acoustically missing internal speech

Thank You

Acknowledgements

Current Lab Members & Affiliates

Christian Brodbeck

Alex Presacco

Proloy Das

Alex Jiao

Dushyanthi Karunathilake

Joshua Kulasingham

Natalia Lapinskaya

Sina Miran

David Nahmias

Peng Zan

Pirazh Khorramshahi

Huan Luo

Mahshid Najafi

Krishna Puvvada

Jonas Vanthornhout

Ben Walsh

Yadong Wang

Juanjuan Xiang

Jiachen Zhuo

Tom Francart

Jonathan Fritz

Michael Fu

Stefanie Kuchinsky

Steven Marcus

Cindy Moss

David Poeppel

Shihab Shamma

Past Lab Members & Affiliates

Nayef Ahmar

Sahar Akram

Murat Aytakin

Francisco Cervantes Constantino

Maria Chait

Marisel Villafane Delgado

Kim Drnec

Nai Ding

Victor Grau-Serrat

Julian Jenkins

Collaborators

Pamela Abshire

Samira Anderson

Behtash Babadi

Catherine Carr

Monita Chatterjee

Alain de Cheveigné

Stephen David

Didier Depireux

Mounya Elhilali

Past Undergraduate Students

Nicholas Asendorf

Ross Baehr

Anurupa Bhonsale

Sonja Bohr

Elizabeth Camenga

Julien Dagenais

Katya Dombrowski

Kevin Hogan

Andrea Shome

James Williams

Funding NIH (*NIDCD, NIA, NIBIB*); NSF; *DARPA*; USDA; UMD