# The Progression of Neural Speech Representations Through Auditory Cortex and Beyond, from Acoustics to Language to Semantics

## Jonathan Z. Simon

**University of Maryland**

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

Mastodon: @jzsimon@fediscience.org

Computational Sensorimotor Systems Lab

http://www.isr.umd.edu/Labs/CSSL/simonlab

University of Lübeck, 1 Aug 2023

# Acknowledgements

**Current Lab Members & Affiliates**
Morgan Belcher
*Vrishab Commuri*
Charlie Fisher
Tejas Guha
Brooke Guo
Michael Johns
Kevin Hu
*Dushyanthi Karunathilake*
Karl Lerud
*Behrad Soleimani*
Ciaran Stone
Craig Thorburn
Allie Vance

**Current & Recent Collaborators**
Samira Anderson
*Behtash Babadi*
Tom Francart
L. Elliot Hong
*Stefanie Kuchinsky*
*Ellen Lau*
Elisabeth Marsh
*Philip Resnik*

**Recent Lab Members & Affiliates**
Sahar Akram
Olivia Bermudez-Hopkins
*Shohini Bhattasali*
*Christian Brodbeck*
Regina Calloway
Francisco Cervantes Constantino
Aura Cruz Heredia
*Proloy Das*
Lien Decruy
Marisel Villafane Delgado
Nai Ding
Jason Dunlap
Sydney Hancock
Marlies Gilles
Alex Jiao
*Neha Joshi*
*Joshua Kulasingham*
Natalia Lapinskaya
Sina Miran
*Mohsen Rezaeizadeh*
Alex Presacco

Krishna Puvvada
Jonas Vanthornhout
Richard Williams
Peng Zan

**Funding & Support**

NIDCD

NIH
NIA

NSF

DARPA

UNIVERSITY OF MARYLAND
18 56

# Outline

- Introduction—Cortical representations of continuous speech

- *Early & fast* cortical representation of continuous speech

- Cortical representations of speech *meaning*

- *Progression* of representations of continuous speech through cortex (bottom-up and top-down)

- Objective measures of speech *intelligibility*

- *Directional functional connectivity* during difficult speech listening

# Outline

- **Introduction—Cortical representations of continuous speech**
- *Early & fast* cortical representation of continuous speech
- Cortical representations of speech *meaning*
- *Progression* of representations of continuous speech through cortex (bottom-up and top-down)
- Objective measures of speech *intelligibility*
- *Directional functional connectivity* during difficult speech listening

# Cortical Representations of <u>Continuous Speech</u>

## *Continuous speech*

- naturalistic

- redundant

- employs auditory cognition

- acoustically rich

- drives most auditory areas

- …

- but also complicated

If you happened to find yourself on the banks of the Ohio River on a particular afternoon in the spring of 1806—somewhere just to the north of Wheeling, West Virginia, say …
*The Botany of Desire* — Michael Pollan

Alfred the Great was a young man, three-and-twenty years of age, when he became king. Twice in his childhood, he had been taken to Rome, where the Saxon nobles were in the habit of going on journeys which they supposed to be religious; …
*A Child's History of England* — Charles Dickens

In the bosom of one of those spacious coves which indent the eastern shore of the Hudson, at that broad expansion of the river denominated by the ancient Dutch navigators …
*The Legend of Sleepy Hollow* — Washington Irving

He was an old man who fished alone in a skiff in the Gulf Stream and he had gone eighty-four days now without taking a fish. In the first forty days a boy had been with him. But after forty days without a fish …
*The Old Man and the Sea* — Ernest Hemingway

# Cortical Representations of Continuous Speech

***Temporal neural patterns $\leftrightarrows$ temporal patterns in speech***

- Generalization of "Speech Tracking"

- Need high temporal precision, for fast temporal speech features

    - EEG (electroencephalography): *whole brain*

    - MEG (magnetoencephalography): *whole brain but with strong cortical bias*

    - ECoG (electrocorticography): *placed cortical surface electrodes*

    - single- and multi-unit recording methods: *placed depth electrodes*

# Cortical Representations of Continuous Speech

## *Neural Representations of Speech*

- oscillations at pitch frequencies (primarily subcortical)    Maddox & Lee (2018) eNeuro

  - acoustic onset tracking    Daube et al. (2019) Curr Biol

    - speech envelope rhythmic following    Lalor & Foxe (2010) Eur J Neurosci

      - phoneme-based responses    Teoh et al. (2022) J Neurosci

        - phoneme-context-based responses    Brodbeck et al. (2018) Curr Biol

          - word-context-based responses    Brodbeck et al. (2022) eLife

            - semantic structure rhythm following    Ding et al. (2016) Nat Neuro

- plus connections to **intelligibility/perception/behavior**

Brodbeck & Simon (2020) *Continuous Speech Processing*, Curr Op Physiol

# Cortical Representations of Speech

- Measure *time-locked* responses to temporal pattern of speech features (in humans)

- Any speech feature of interest: acoustic envelope, lexical, pitch, semantic, etc.

- Infer spatio-temporal neural origins of neural responses



his schoolhouse    was a  low      building  of   one        large          room      rudely   constructed      of  logs

Speech envelope

"Decoder"

"Temporal response functions"
(TRFs)

# Cortical Representations: Encoding

- Predicting future neural responses from present stimulus features,
  - wide variety of stimulus features
  - via Temporal Response Function (TRF)

- Why look at encoding? It *often* tells us more about the brain
  - TRF analogous to evoked response
  - peak amplitude ≈ processing intensity
  - peak latency ≈ source location
  - multiple TRFs simultaneously

large        room        rudely        constructed        of    logs

"Temporal response functions" (TRFs)

Example: MEG Prediction of Voxel Responses

# Temporal Response Functions

# TRF Model Estimation & Fit

**Temporal Response Function (TRF) estimation:**

Stimulus and response are known; find the best TRF
 to produce the response from the stimulus:



Resp.

Stim.

Estimated TRF

Actual response

Resp.

Predicted response (Stimulus * TRF)

Lalor & Foxe (2010) *Neural Responses to Uninterrupted Natural Speech* …  Eur J Neurosci
Ding & Simon (2012) *Neural Coding of Continuous Speech in Auditory Cortex* … , J Neurophys

# Simultaneous Temporal Response Functions

- TRFs predict neural response to speech

  ‣ Analogous to evoked response

  ‣ Peak amplitude ≈ processing intensity

  ‣ Peak Latency ≈ source location

- Multiple TRFs estimated simultaneously

  ‣ compete to explain variance (advantage over evoked response)



Speech Representations        TRFs

Measured Neural signals

Predicted Neural signals

Crosse et al. (2016) *The Multivariate Temporal Response Function (mTRF) Toolbox* … , Front Hum Neurosci
Brodbeck et al. (2021) *Eelbrain: A Python Toolkit for Time-Continuous Analysis* … , bioRxiv

# Cortical Representations Across Cortex



**Post-Auditory Cortex**

**Auditory Cortex**

**Semantic Processing**
*semantic composition*

**Lexical Processing**

*word onsets*

**Higher Order Auditory**
*speech onsets*

**Primary Auditory**
*speech onsets*

attend

ignore

200 ms

attend

ignore

150 ms

*lexical cohort entropy*

attend

ignore

150 ms

*high frequency input*

fast envelope

carrier

50 ms

attend

ignore

150 ms

100 ms

*dominated by focus of selective attention*

# Outline

- Introduction—Cortical representations of continuous speech
- *Early & fast* cortical representation of continuous speech
- Cortical representations of speech *meaning*
- *Progression* of representations of continuous speech through cortex (bottom-up and top-down)
- Objective measures of speech *intelligibility*
- *Directional functional connectivity* during difficult speech listening

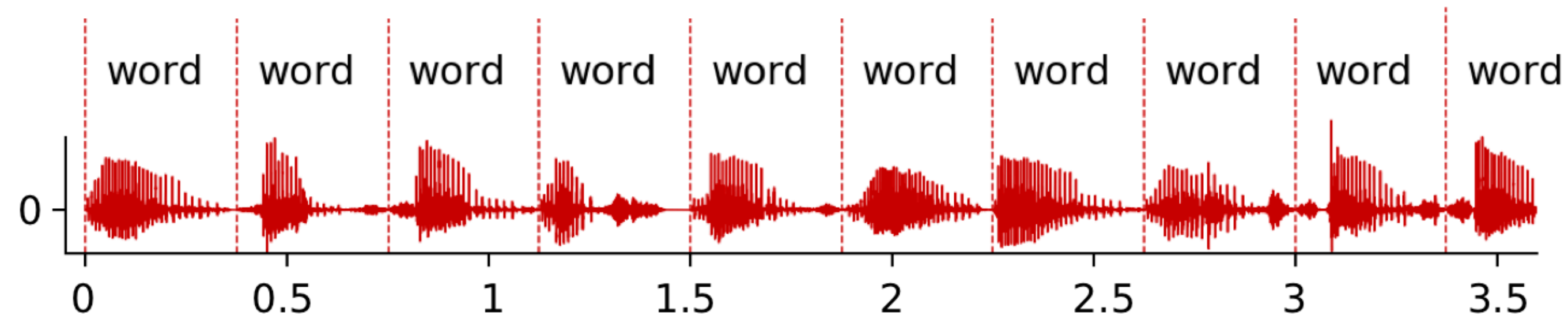# Fast & Early Cortical Representations



TRF (MEG) for
70-200 Hz
continuous speech
*envelope*

40 ms latency peak
⟹ Primary/Core auditory cortex

Kulasingham et al. (2020) *High Gamma Cortical Processing of Continuous Speech …*, NeuroImage
Simon et al. (2022) *… the High-Gamma Band: A Window into Primary Auditory Cortex*, Front Neurosci

# Fast & Early C [...] esentations



Attend Male

Ignore Male

TRF Amplitude (std units)

$2\times10^{-4}$

Time [ms]

***

Attend    Ignore

Attend > Ignore
Primary cortex response depends on selective attention

Commuri et al. (2023) *Cortical Responses ... in the High-Gamma Band Depend on Selective Attention*, bioRxiv

# Outline

- Introduction—Cortical representations of continuous speech
- *Early & fast* cortical representation of continuous speech
- Cortical representations of speech *meaning*
- *Progression* of representations of continuous speech through cortex (bottom-up and top-down)
- Objective measures of speech *intelligibility*
- *Directional functional connectivity* during difficult speech listening

# Cortical Representations Across Cortex



**Auditory Cortex**

**Post-Auditory Cortex**

*Semantic Processing*

*Higher Order Auditory*

*Lexical Processing*

semantic composition

*Primary Auditory*

word onsets

200 ms

speech onsets

attend

speech onsets

attend

ignore

high frequency input

150 ms

ignore

lexical cohort entropy

attend

ignore

fast envelope

150 ms

ignore

150 ms

carrier

100 ms

50 ms

*dominated by focus of selective attention*

# Speech Understanding/Meaning

- Behavioral correlates of speech understanding

  - implies language comprehension

  - structural comprehension

    - sentence structure

    - other structures, e.g. poetic, logical

- Neural correlates of speech understanding

  - rhythms of structural comprehension/meaning,
    even if *fully absent in the acoustics*

    - sentence structures

    - poetic structures

    - mathematical structures

Ding et al., Nat Neurosci 2016
Teng et al., Curr Biol 2020

# Isochronous Speech

Acoustics



Acoustical Spectrum (envelope)

# Isochronous Arithmetic



Acoustics

Acoustical Spectrum

Neural Spectrum

Kulasingham et al. (2021) *Cortical Processing of Arithmetic and Simple Sentences …*, J Neurosci

# Isochronous Cocktail Party

# Isochronous Cocktail Party

**Neural Spectrum**

**Attend to Sentences**

**Attend to Equations**

**TRFs**

left hemisphere sources

right hemisphere sources

left hemisphere sources

right hemisphere sources

510 ms

810 ms

1640 ms

1850 ms

1940 ms

1200 ms

1920 ms

2040 ms

Frequency [Hz]

Time [ms]

Amplitude [a.u]

[fT/√Hz]

# Representations of Understanding



Attend to Sentences

Attend to Equations
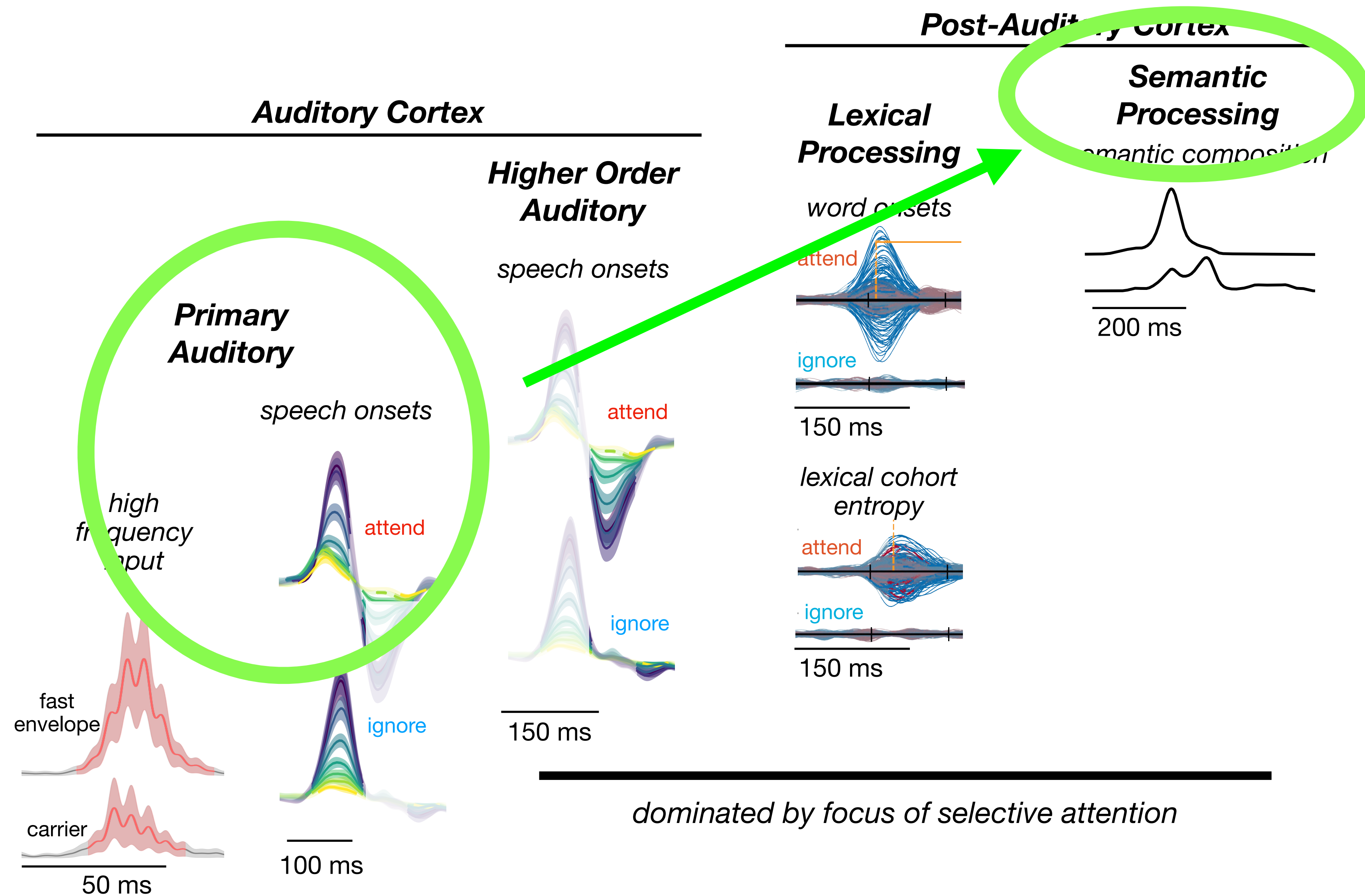
**Neural Correlation with Behavior**

# Neural Markers of Comprehension

- Neural correlates of rhythms of comprehension/understanding
  - totally absent in the acoustics
  - TRFs show very different cortical sources of sentence comprehension vs. mathematical equation comprehension
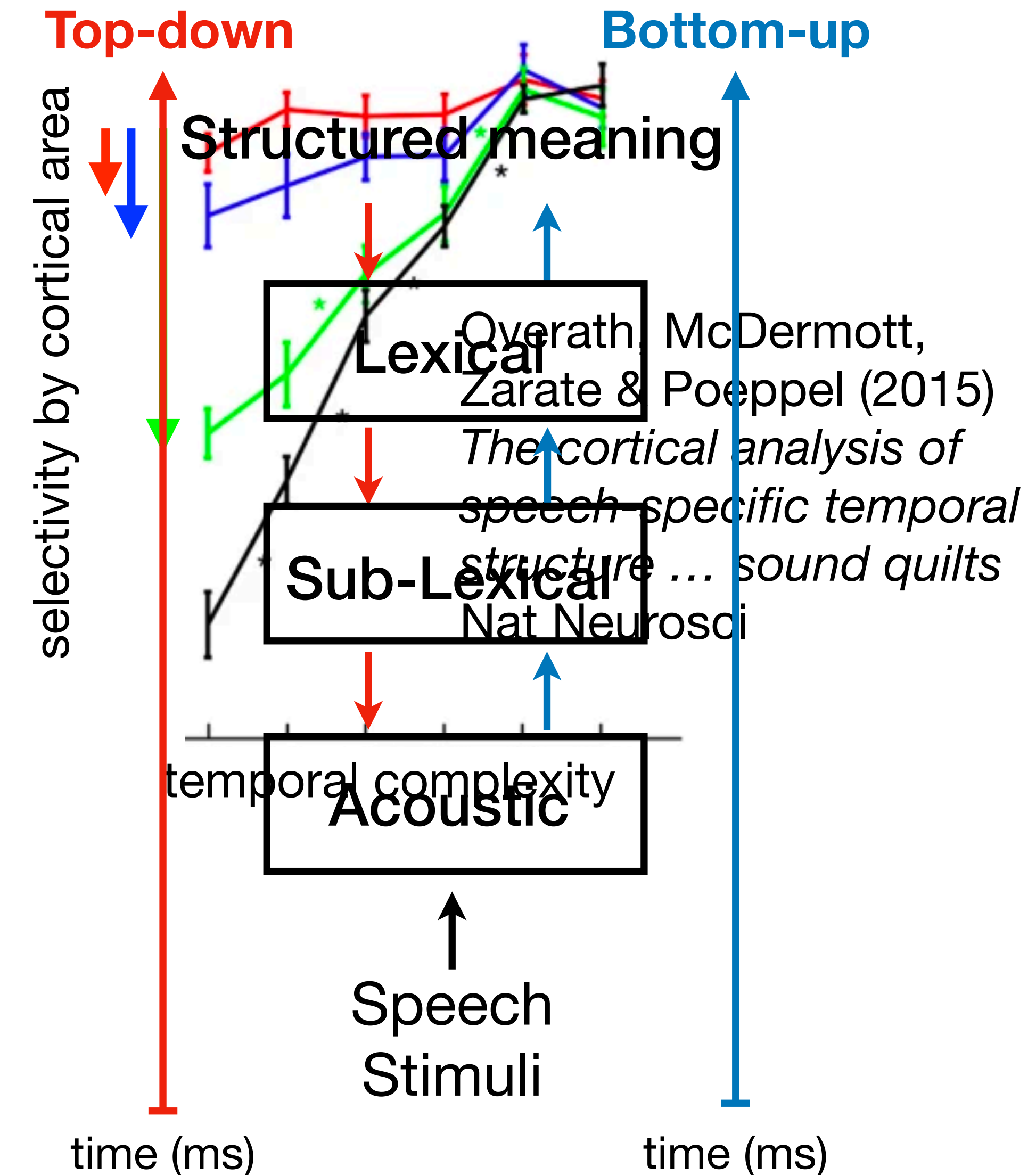  - neural responses correlated with behavior

# Outline

- Introduction—Cortical representations of continuous speech
- *Early & fast* cortical representation of continuous speech
- Cortical representations of speech *meaning*
- *Progression* of representations of continuous speech through cortex (bottom-up and top-down)
- Objective measures of speech *intelligibility*
- *Directional functional connectivity* during difficult speech listening

# Cortical Representations Across Cortex



**Post-Auditory Cortex**

**Auditory Cortex**

**Semantic Processing**

*semantic composition*

**Lexical Processing**

**Higher Order Auditory**

*speech onsets*

*word onsets*

attend

ignore

200 ms

**Primary Auditory**

*speech onsets*

attend

*high frequency input*

*lexical cohort entropy*

attend

ignore

150 ms

ignore

ignore

150 ms

*fast envelope*

ignore

150 ms

*dominated by focus of selective attention*

*carrier*

100 ms

50 ms

# Progression of Speech Representations

- Previous fMRI research on which brain regions process which speech and language features

- Progression of feature-based (bottom-up) levels
  - complex auditory stimulus, to
  - speech sounds, to
  - linguistic information via speech sounds

- Not all processing is straight bottom up
  - selective attention
  - secondary processing upon "error" detection

- MEG & EEG excel at showing temporal (i.e., latency) progression of processing



**Top-down**    **Bottom-up**

selectivity by cortical area

Structured meaning

Lexical

Overath, McDermott, Zarate & Poeppel (2015) *The cortical analysis of speech-specific temporal structure … sound quilts* Nat Neurosci

Sub-Lexical

temporal complexity

Acoustic

Speech Stimuli

time (ms)    time (ms)

Brodbeck et al. (2022) *Parallel Processing in Speech Perception: Local and Global Representations…*, eLife

# Experimental Design

**Task**            Listening to 1-minute long passages
                    The Botany of Desire (Michael Pollan)

**Stimuli**         4 passage types

                        – Speech modulated noise

                        – Non-words
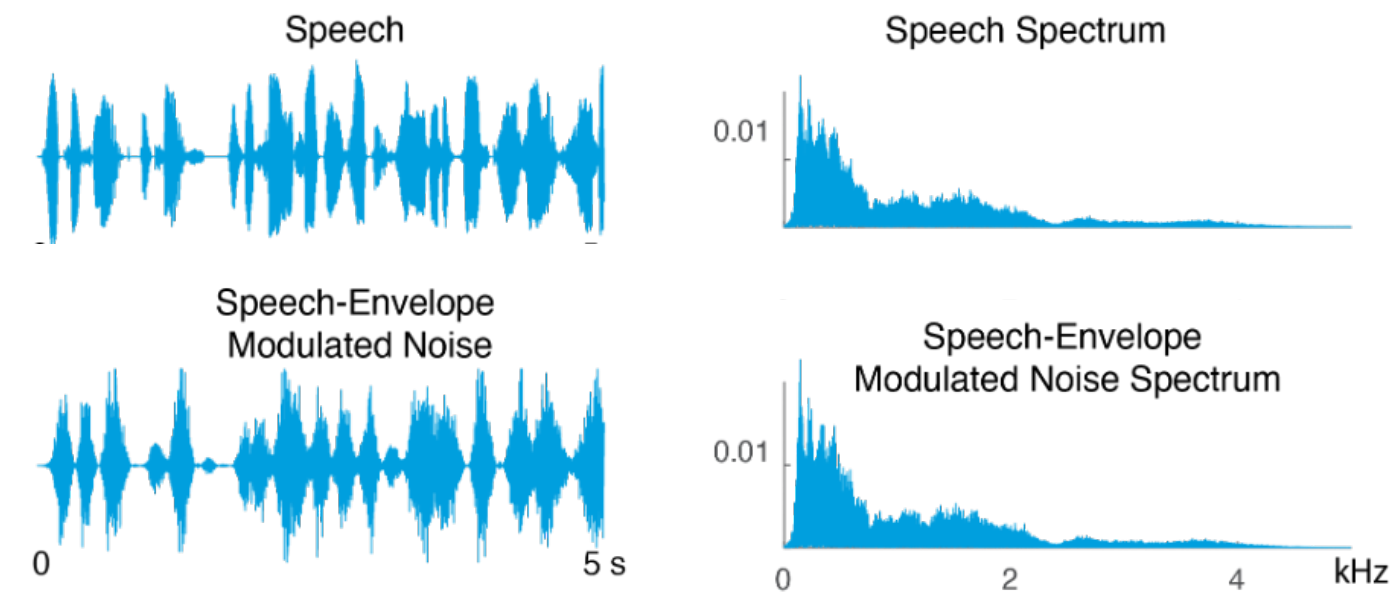
                        – Scrambled words

                        – Narrative

                    Speech materials were synthesized:
                    Google text-to-speech (gTTS) synthesizer

# Experimental Design

**Speech-envelope Modulated Noise**

Speech

Speech Spectrum
0.01

Speech-Envelope Modulated Noise

Speech-Envelope Modulated Noise Spectrum
0.01

0          5 s

0     2     4     kHz

**Non-words**

Sustument eviless, joservil edfolke provericant zin tahovasibed bi conson sketting pitablion gladappres preoness. Feno unknoways, chasizer, giiz, warrowied tanatum impinges. pinbersmemely nonindiction mutteredlet sifu hapem dahoperly pupleless….

**Scrambled words**

A liquid is only speak, second even for good reach the attack us. Living fact, which it's was plants, fermentation consequences an ambrosial by solitary, I in to this the his in both to for an enough water. Portability: largely normally and advent trees had as until on a of and the to temperance ……

**Narrative**

If you happened to find yourself on the banks of the Ohio River on a particular afternoon in the spring of 1806-somewhere just to the north of Wheeling, West Virginia, say, you would probably have noticed a strange makeshift craft drifting lazily down the river. At the time, this particular …..
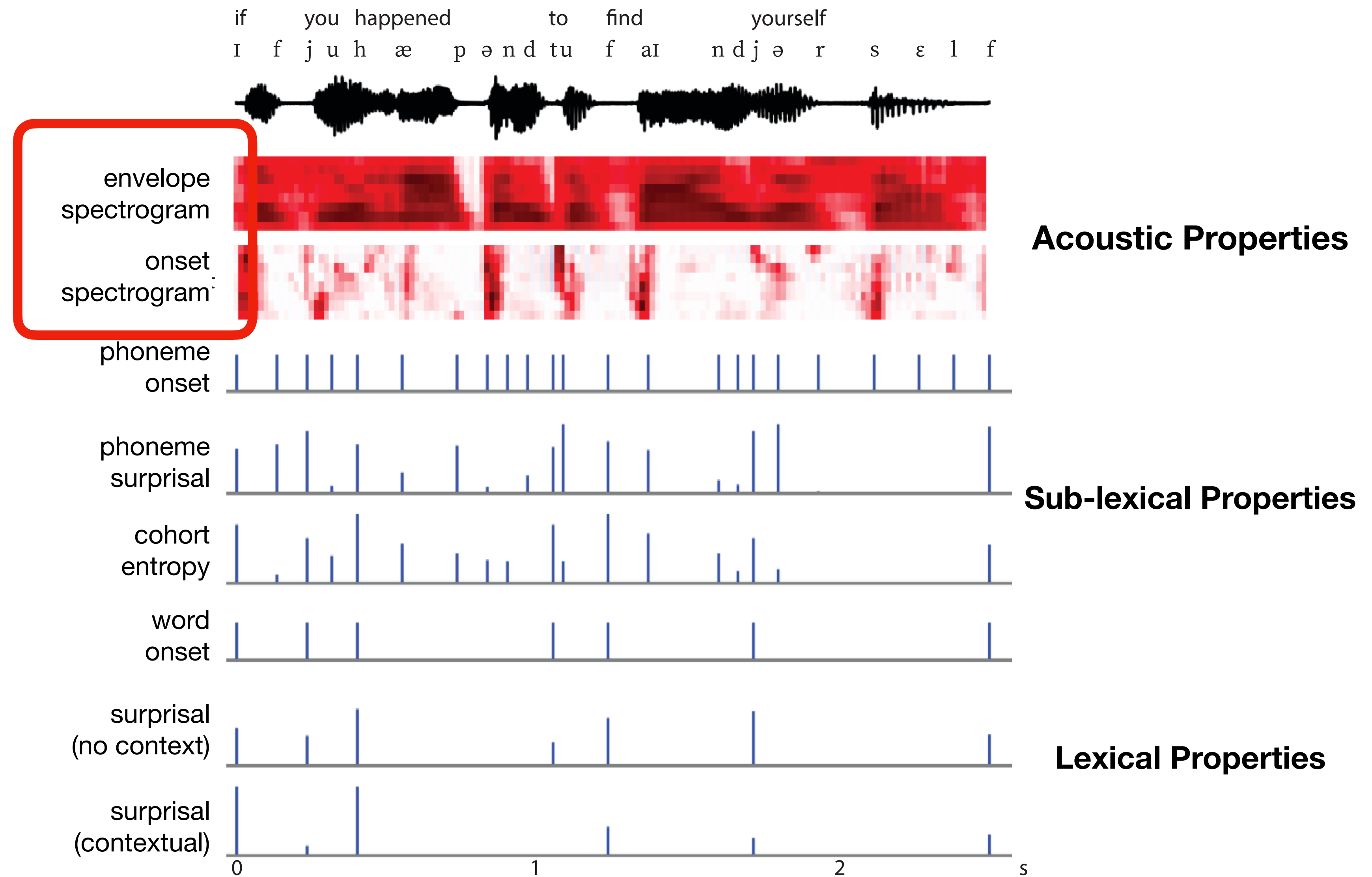
continuous-speech-like prosody and rhythm

Karunathilake et al. *in preparation*

# Simultaneous Temporal Response Functions

- TRFs predict neural response to speech

  ‣ Analogous to evoked response

  ‣ Peak amplitude ≈ processing intensity

  ‣ Peak Latency ≈ source location

- Multiple TRFs estimated simultaneously

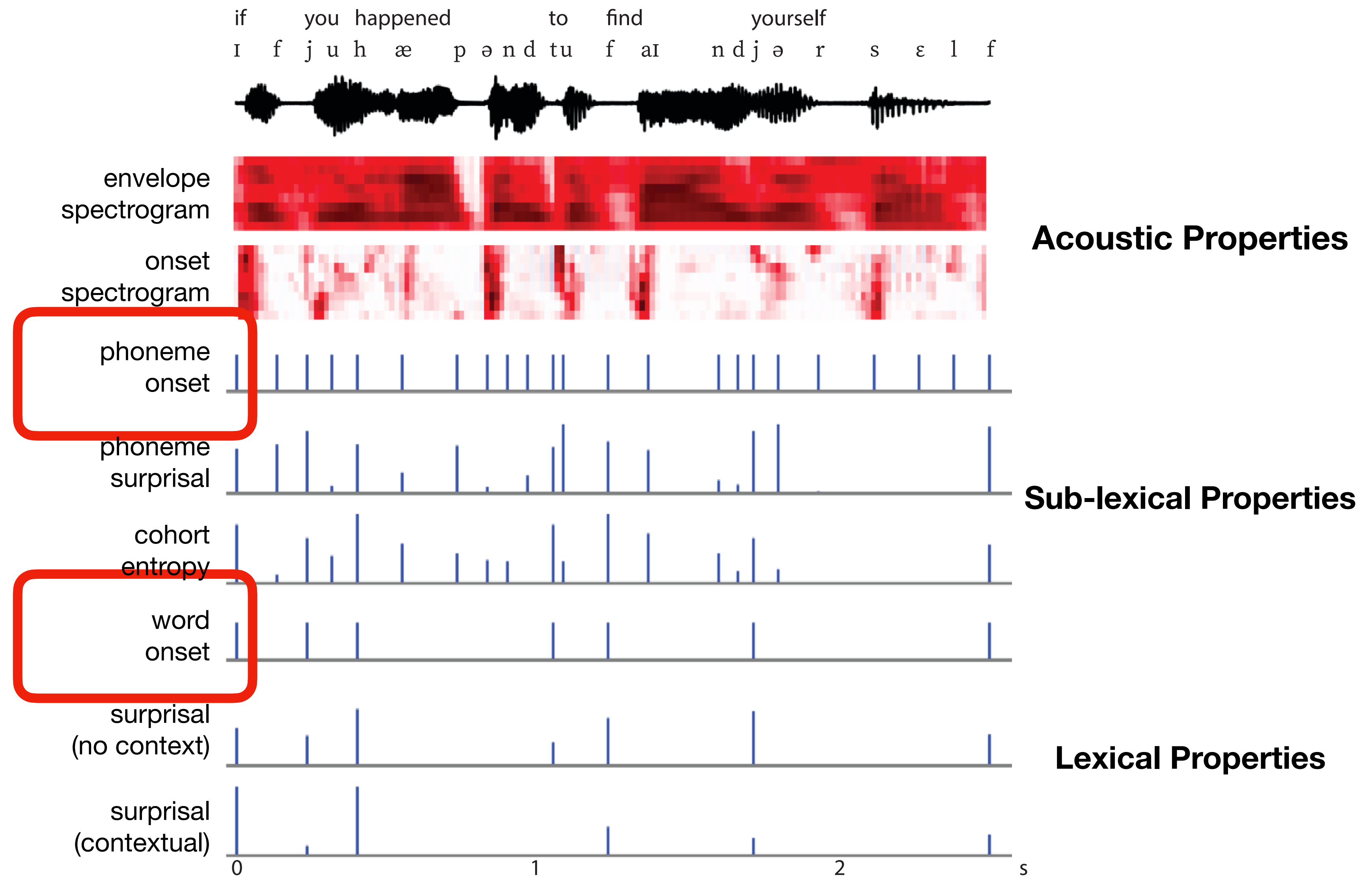  ‣ compete to explain variance (advantage over evoked response)



Speech Representations

TRFs

Measured Neural signals

Predicted Neural signals

# Speech Representations



if    you happened      to  find     yourself
ɪ   f  j u h   æ   p ə n d  t u   f  aɪ   n d j ə   r   s   ɛ   l   f

envelope spectrogram

onset spectrogram

**Acoustic Properties**

phoneme onset

phoneme surprisal

cohort entropy

word onset

**Sub-lexical Properties**

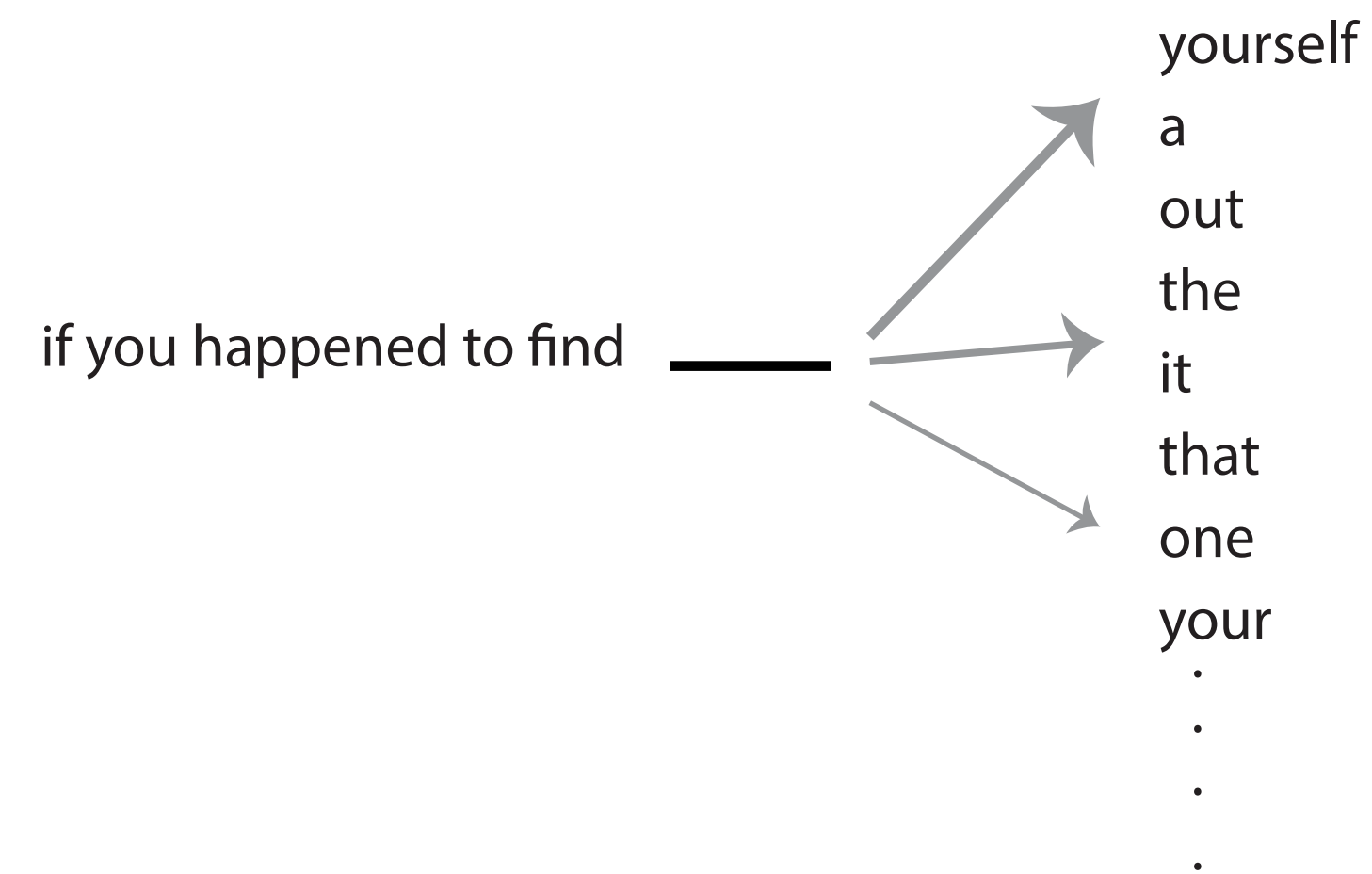surprisal (no context)

surprisal (contextual)

**Lexical Properties**
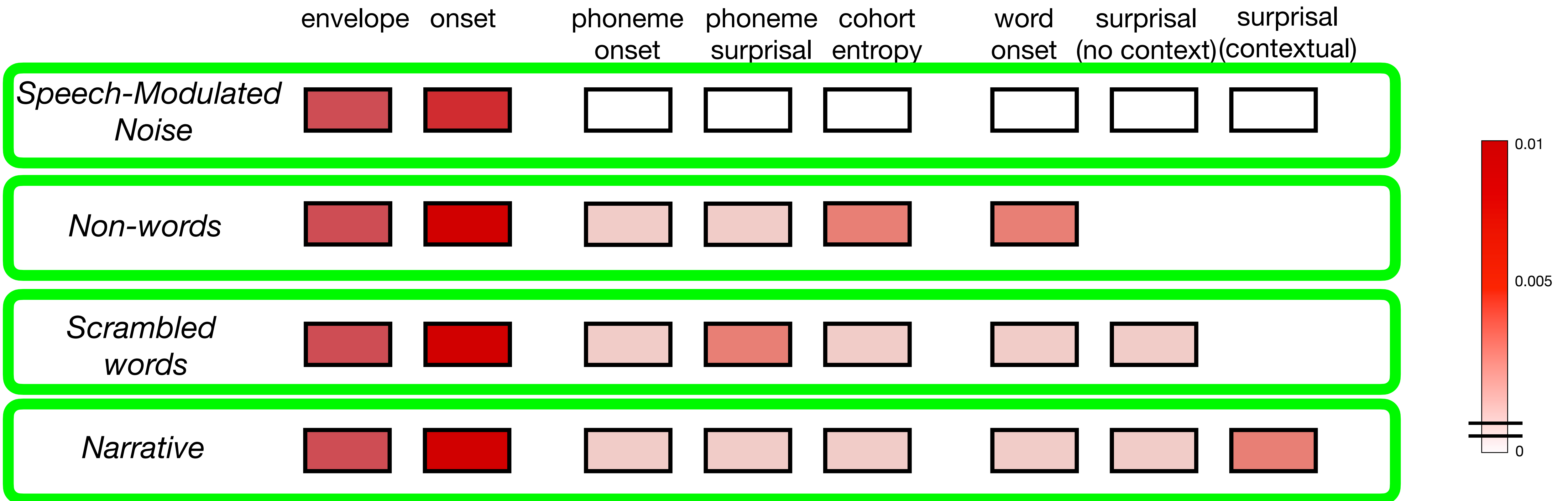
0                    1                    2        s

# Speech Representations

# Speech Representations

# Speech Representations

# Speech Representations

# Speech Representations

# Neural Prediction Results

**Emergence of neural features as the incremental processing occur**



| | envelope | onset | phoneme onset | phoneme surprisal | cohort entropy | word onset | surprisal (no context) | surprisal (contextual) |
|---|---|---|---|---|---|---|---|---|
| *Speech-Modulated Noise* | ■ | ■ | □ | □ | □ | □ | □ | □ |
| *Non-words* | ■ | ■ | ■ | ■ | ■ | ■ | | |
| *Scrambled words* | ■ | ■ | ■ | ■ | ■ | ■ | ■ | |
| *Narrative* | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ |

- Acoustic features are encoded for both non-speech and speech stimuli

- (Sub)-lexical features are encoded only when (sub)-lexical boundaries are intelligible

- Context based word surprisal emerges for narrative passage

- When context supports, context based surprisal is better tracked compared to naive surprisal

# Hemispheric Lateralization Results



_**Speech feature**_

**Envelope Onset**

Envelope

Phoneme Onset
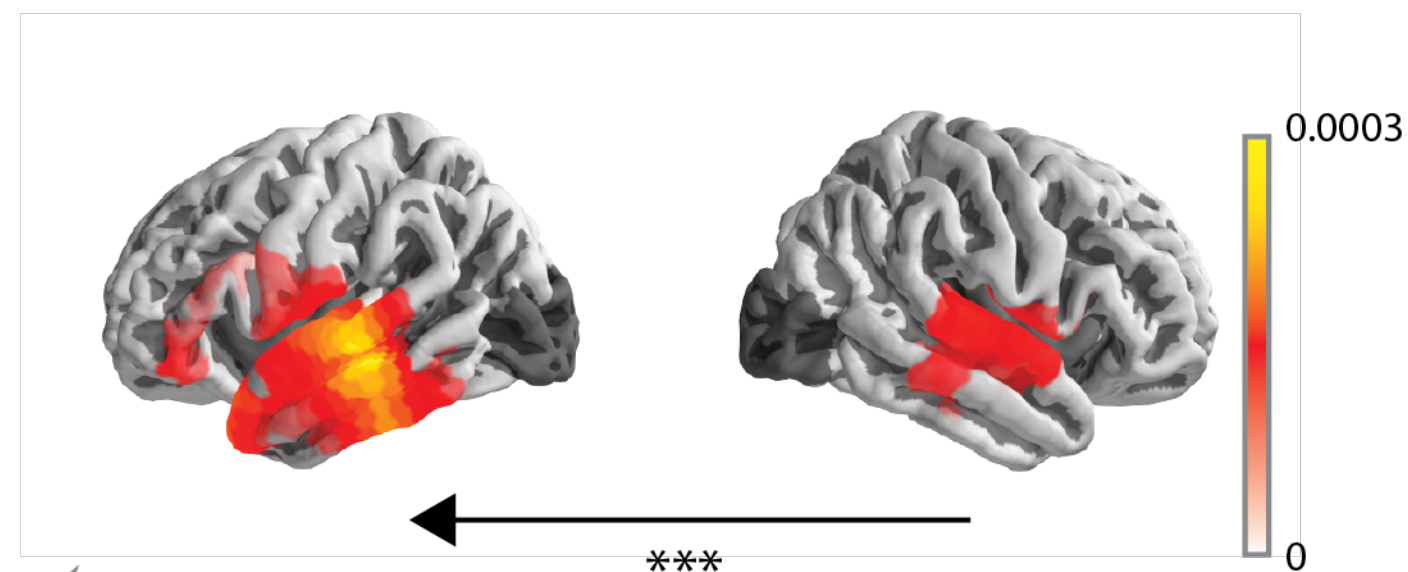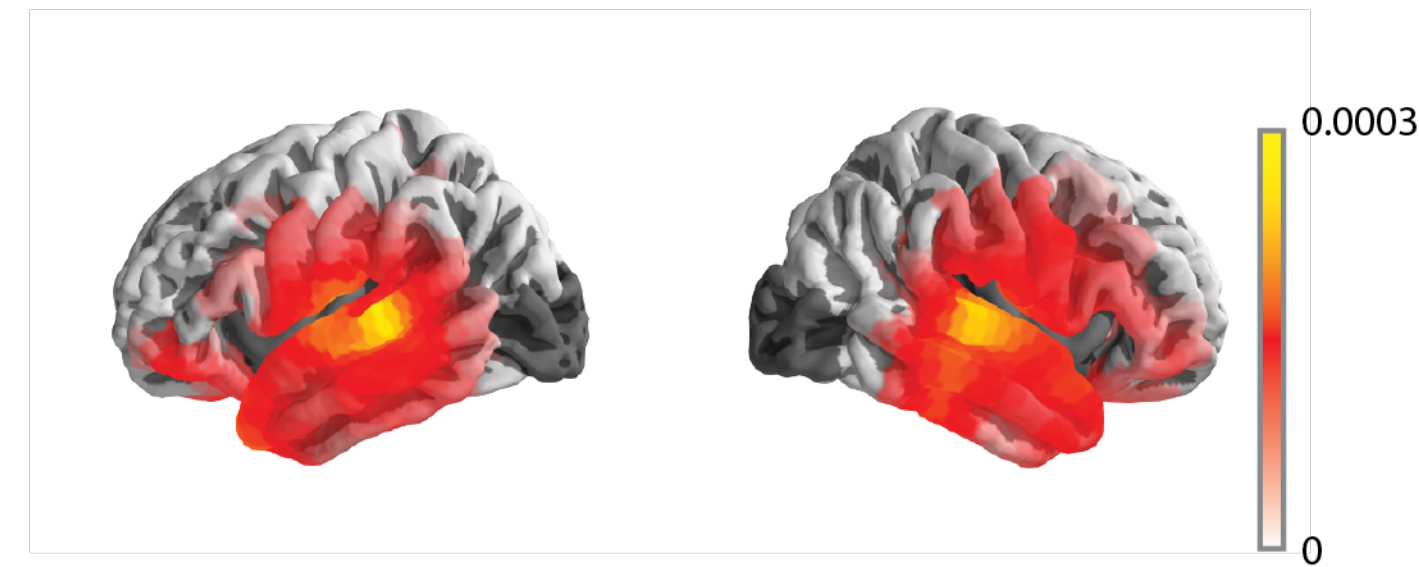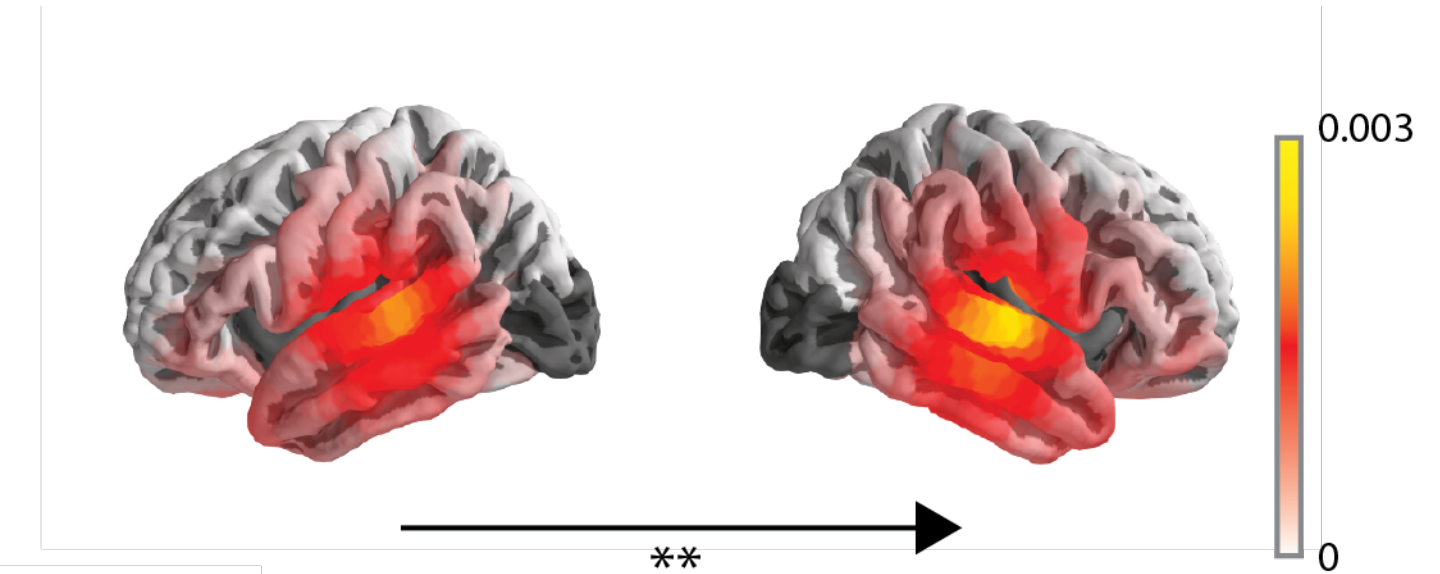
**Phoneme Surprisal**

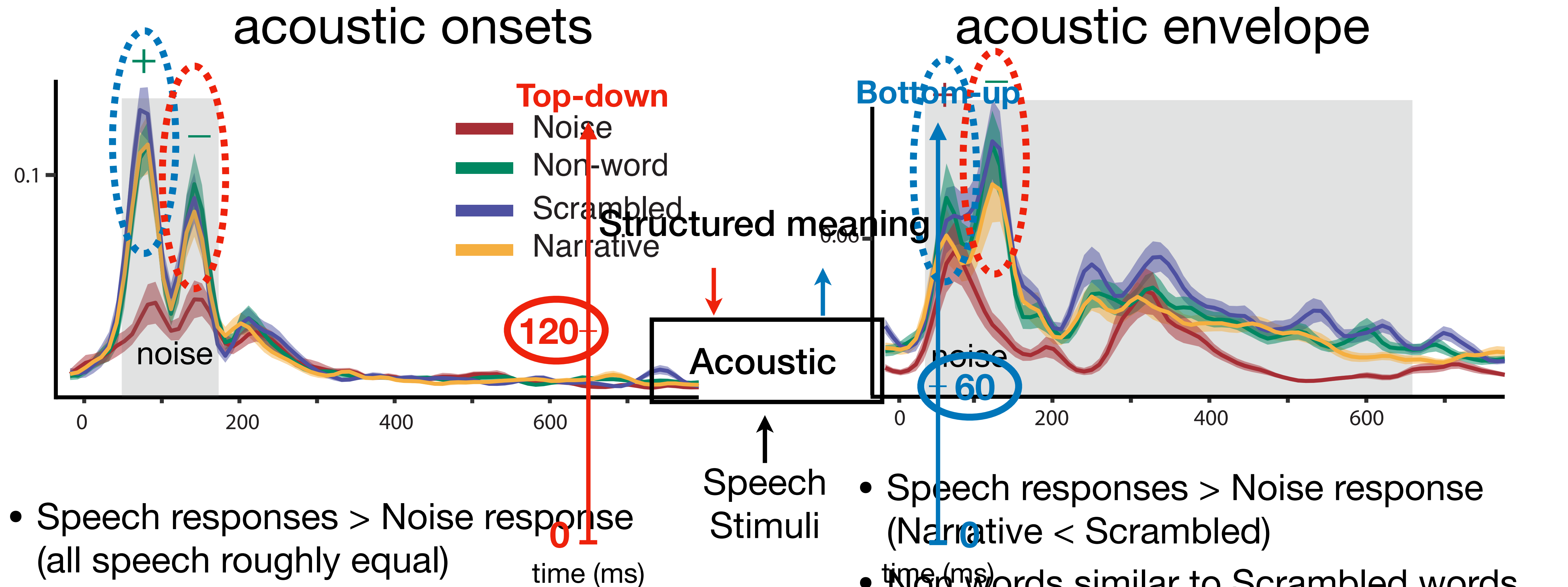Cohort Entropy

Word Onset

Unigram Surprisal

**GPT2 Surprisal**

**Left Lateralized**      **Bilateral**      **Right Lateralized**

Note: lateralization results can be task dependent

# Acoustic TRF Results



**acoustic onsets**

**acoustic envelope**

**Top-down**

Noise
Non-word
Scrambled
Narrative

**Structured meaning**

**Bottom-up**

120

0

Acoustic

Speech
Stimuli

time (ms)

- Speech responses > Noise response
  (all speech roughly equal)

- Speech responses > Noise response
  (Narrative < Scrambled)

- Non words similar to Scrambled words

- Noise response lacks 2nd peak ~120 ms
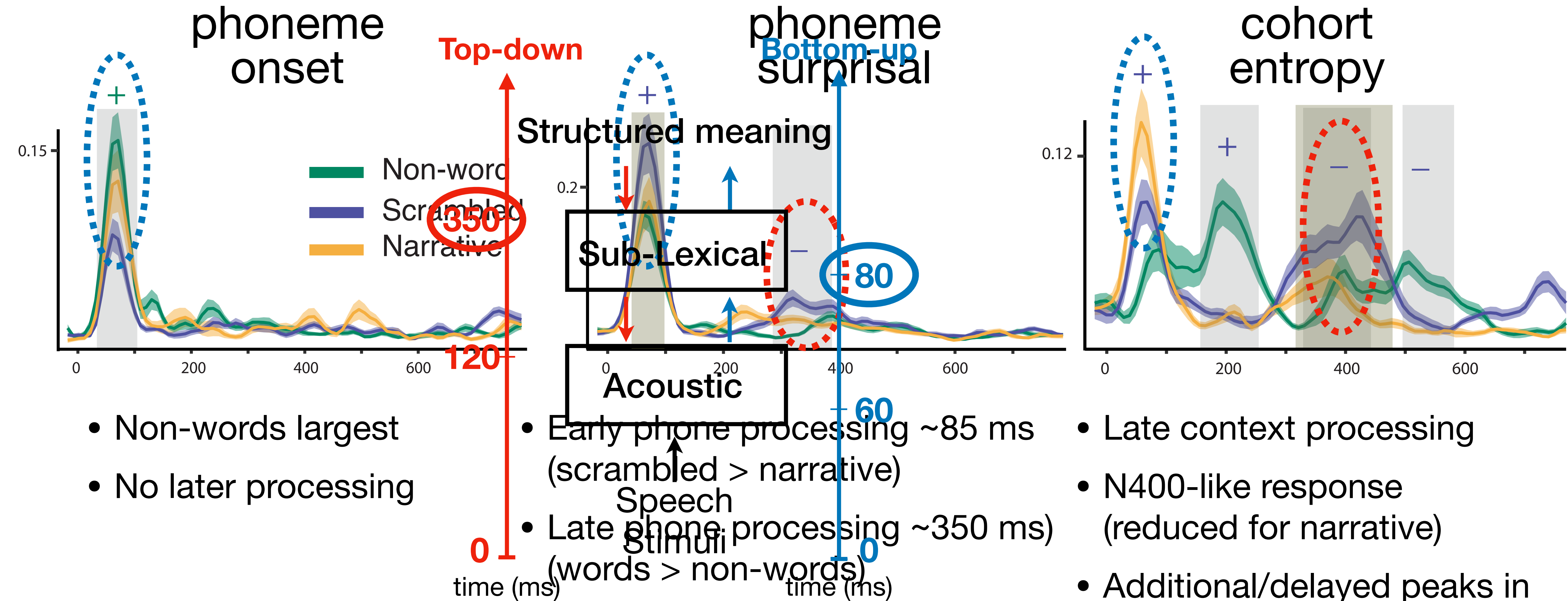
60 ms: acoustic bottom-up processing

120 ms: acoustic but attention-dependent

right hemisphere shown
condition based differences similar in left

# Phonemic TRF Results



## phoneme onset

**Top-down**

- Non-words largest
- No later processing

## phoneme surprisal

**Bottom-up**

Structured meaning

Sub-Lexical

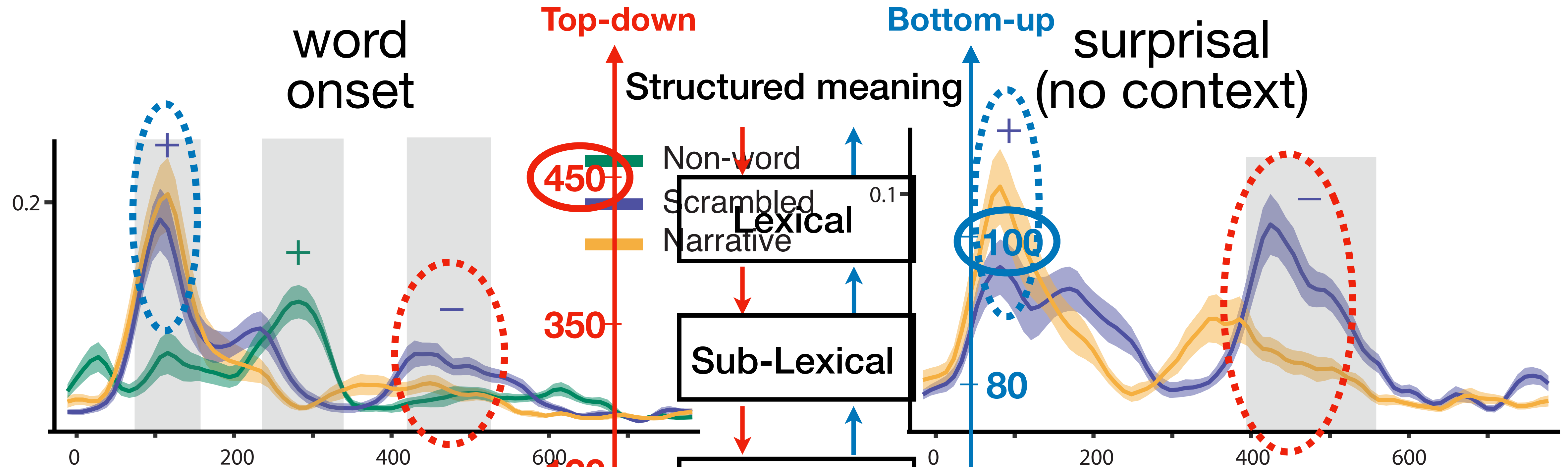Acoustic

Speech Stimuli

time (ms)

- Early phone processing ~85 ms (scrambled > narrative)
- Late phone processing ~350 ms) (words > non-words)

## cohort entropy

- Late context processing
- N400-like response (reduced for narrative)
- Additional/delayed peaks in non-words (difference in stimulus distributions)

left hemisphere shown (right similar)

**85 ms: simple phoneme processing**
**350 ms: additional further processing**

Legend: Non-word, Scrambled, Narrative

# Word-based TRF Results



**word onset**

**Top-down**

**Bottom-up**

**surprisal (no context)**

**Structured meaning**
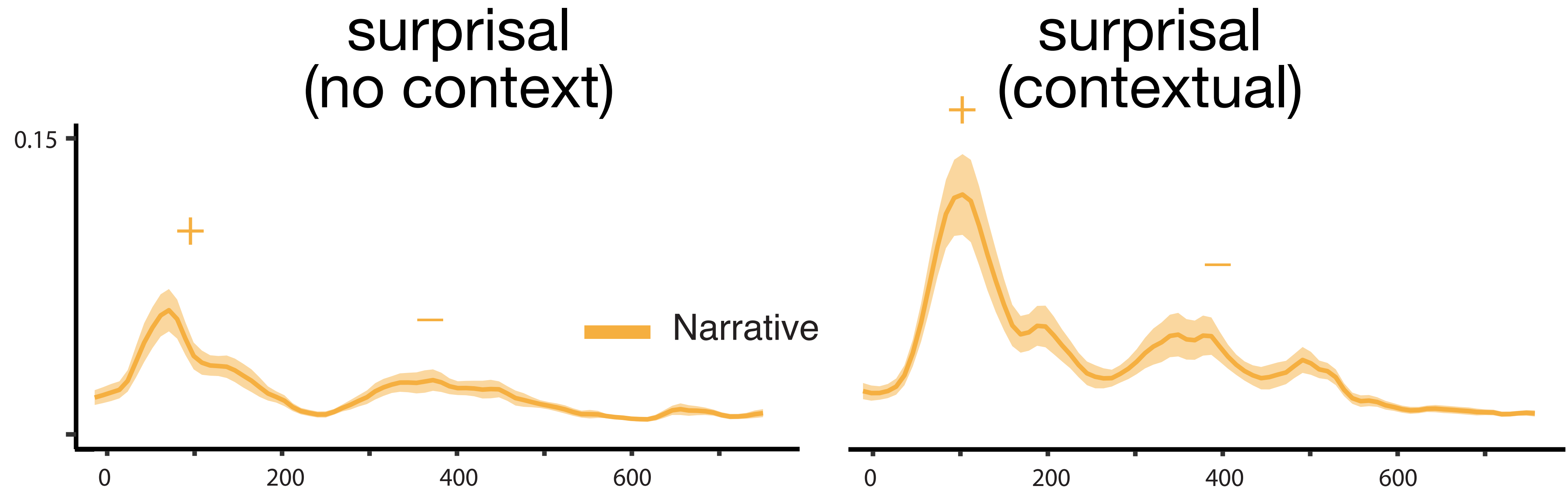
- 450+
- Non-word — Scrambled — Narrative
- Lexical
- Sub-Lexical
- Acoustic

- Scrambled ≈ narrative for rapid processing
- Scrambled words > narrative at ~450 ms
- words: Left hemi > Right (non-words: L ≈ R)

- N400 like response
- Reduction in surprisal when context
- Left hemi > Right hemi
- Right hemisphere: Scrambled ≈ Narrative

100 ms: simple word processing
450 ms: "error" correction processing

Speech Stimuli

left hemisphere shown
(right much weaker except for non-word onset)

# Contextual Word Surprisal Results



surprisal
(no context)

surprisal
(contextual)

0.15

+

−

Narrative

+

−

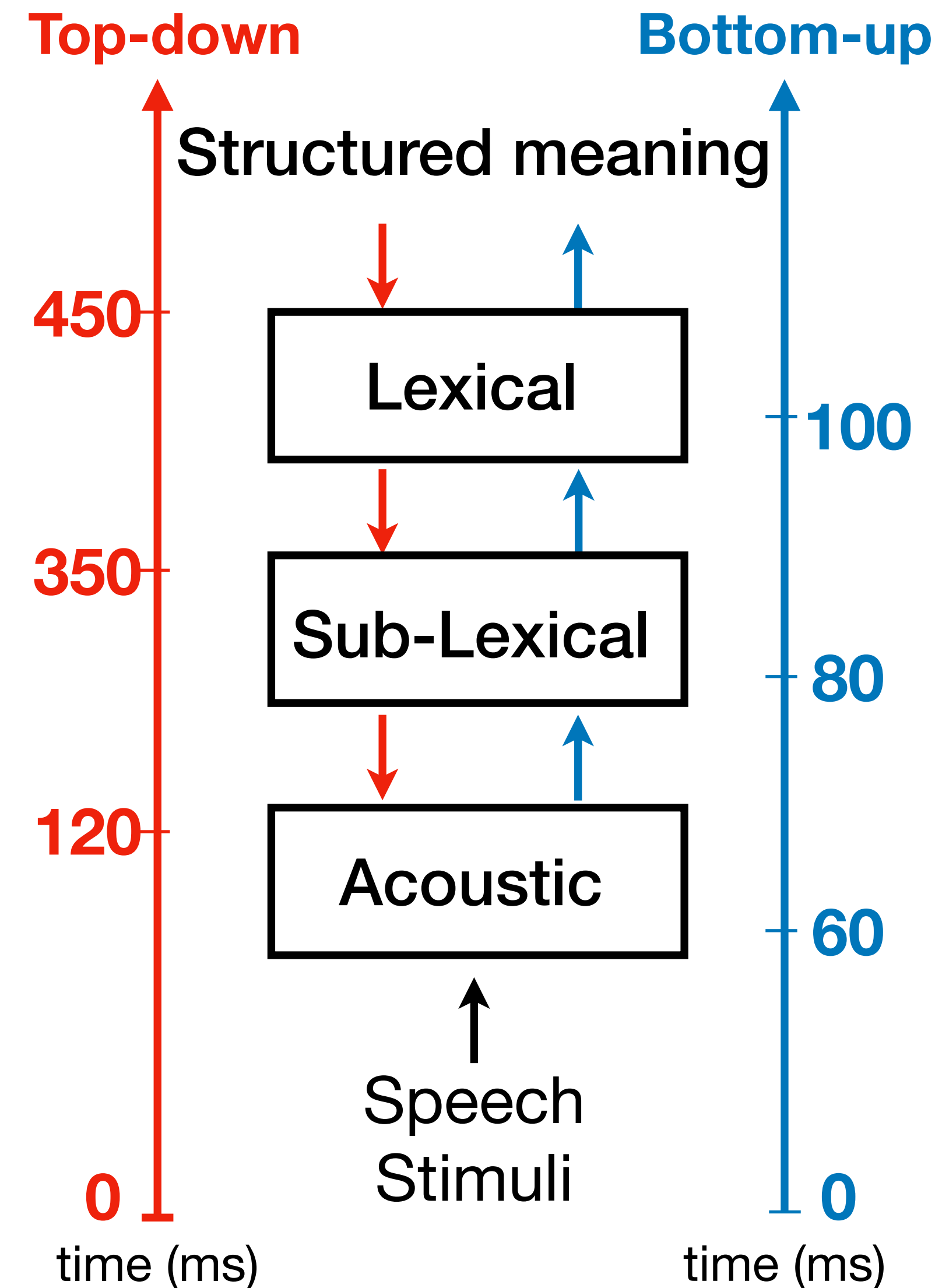0    200    400    600

0    200    400    600

- When context helps, context-based surprisal is better tracked than raw surprisal

- N400 like response in both predictors

left hemisphere shown
(right much weaker)
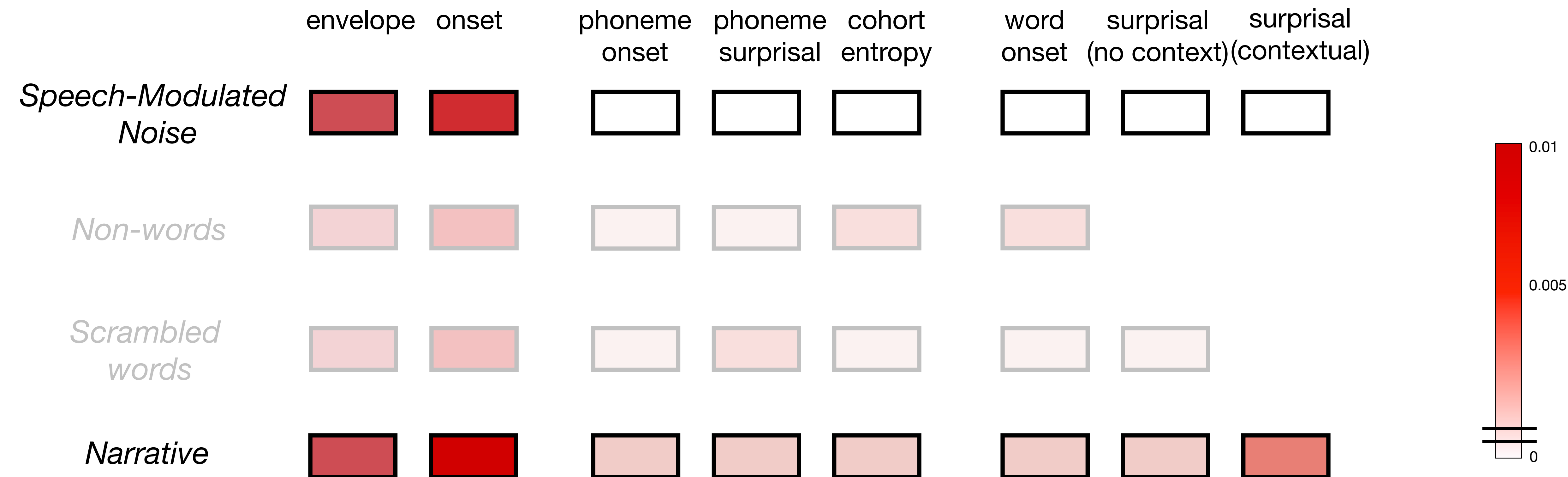
# Neural Speech Processing Progression

- Cortical response time-locks to emergent features from acoustics to context as incremental steps in the processing of speech input occur

- Higher level processing / top-down mechanisms may affect lower level speech processing

- Linguistic features are processed when the linguistic boundaries are intelligible

- Lower-level acoustic feature responses are bilateral but right lateralized whereas, context based responses are strongly left lateralized

**Top-down**   **Bottom-up**

Structured meaning

450 — Lexical

350 — Sub-Lexical

120 — Acoustic

Speech Stimuli

100

80

60

0

time (ms)   time (ms)

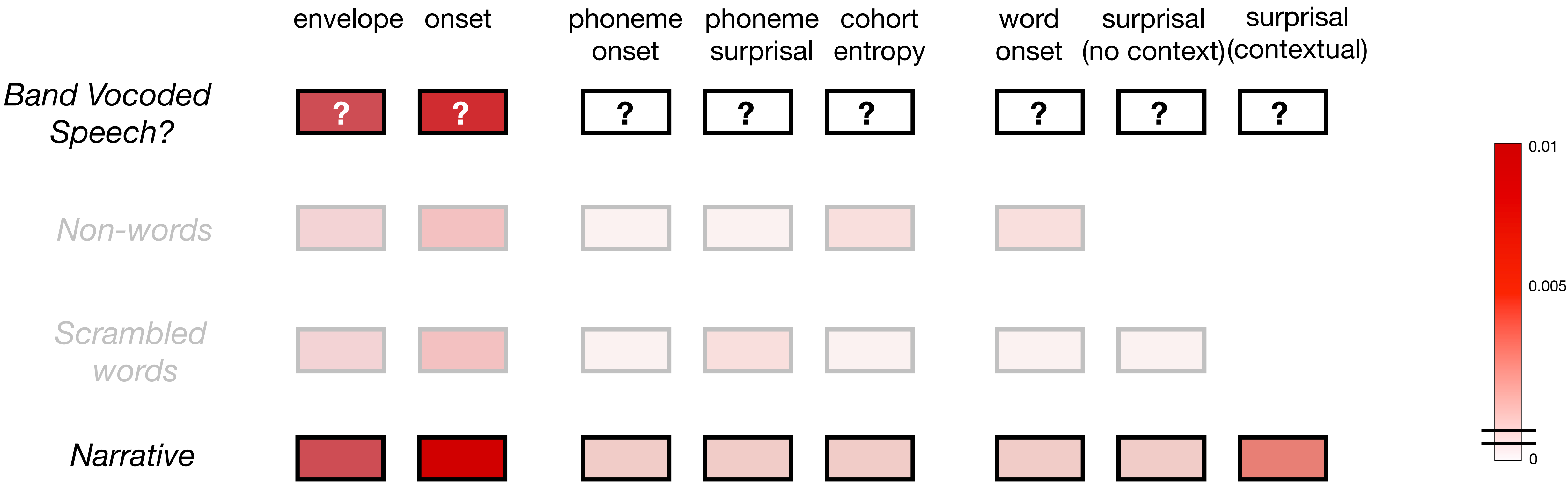Karunathilake et al. *in preparation*

# Outline

- Introduction—Cortical representations of continuous speech
- *Early & fast* cortical representation of continuous speech
- Cortical representations of speech *meaning*
- *Progression* of representations of continuous speech through cortex (bottom-up and top-down)
- Objective measures of speech *intelligibility*
- *Directional functional connectivity* during difficult speech listening

# Previous Neural Prediction Results

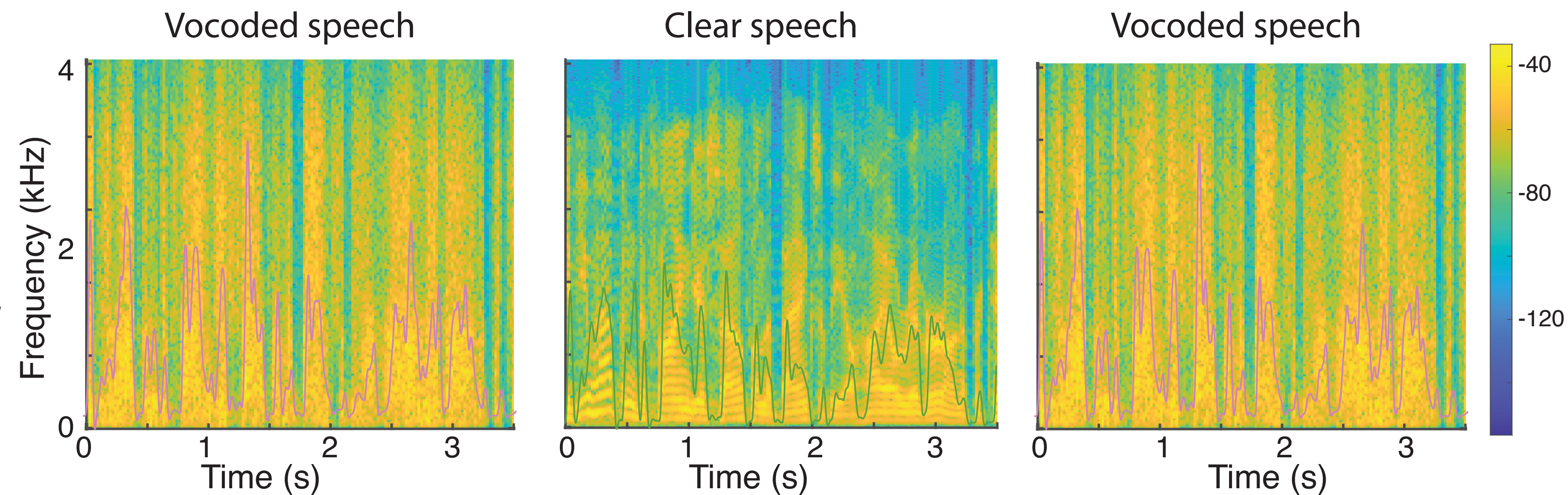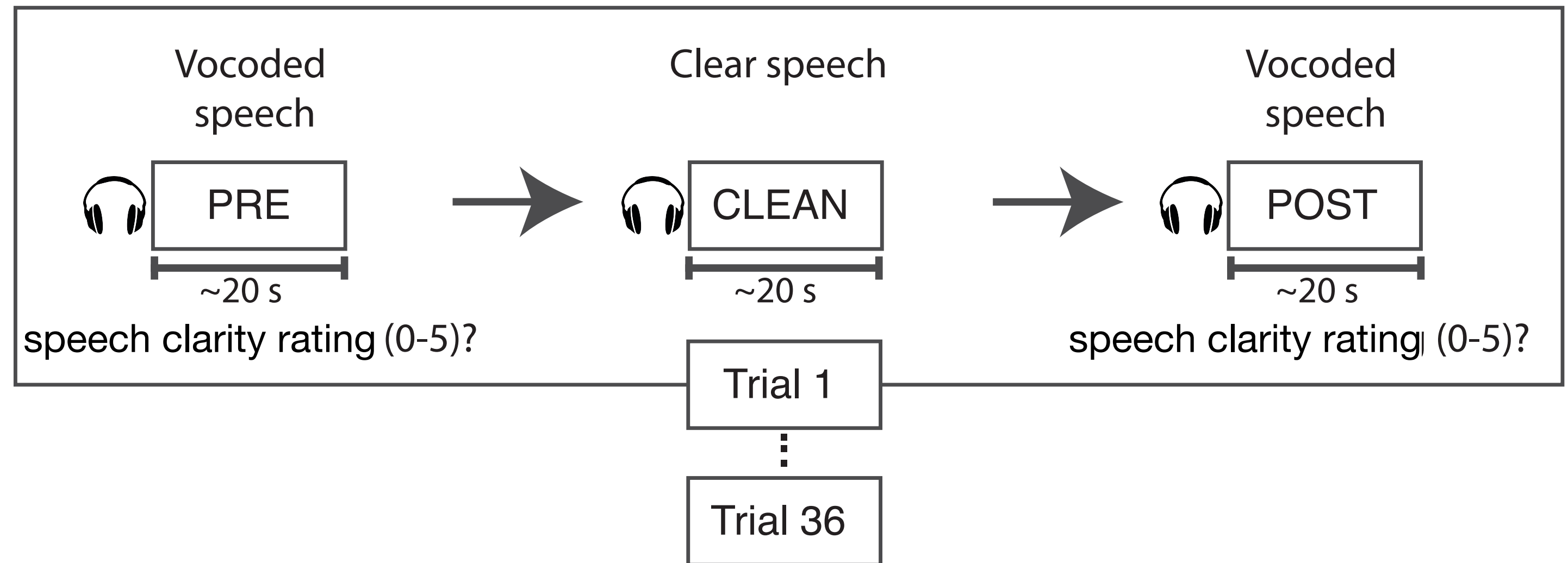# Possible Neural Prediction Results
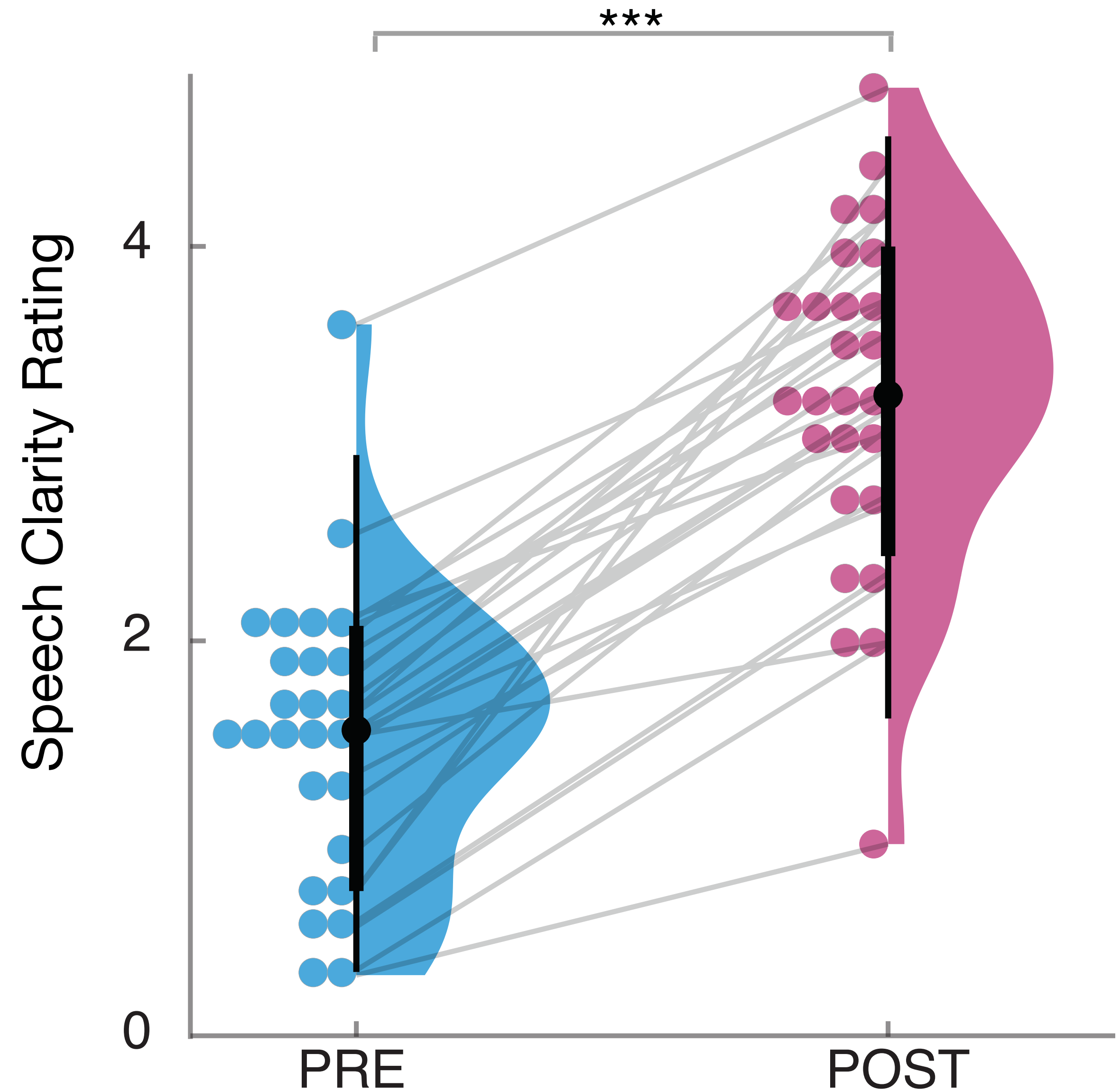
# Intelligibility Experimental Design

- Manipulate intelligibility but keep acoustics unchanged

  - Speech acoustics: three-band noise-vocoded speech

  - Intelligibility manipulated via priming

- Hypothesized intelligibility measure(s)

  - word boundaries

*"Slice an apple through at its equator, and you will find five small chambers arrayed in a perfectly symmetrical starburst—a pentagram."*



Karunathilake et al. (2023) *Neural Tracking Measures of Speech Intelligibility…*, bioRxiv

# Intelligibility Behavioral Results

Speech Clarity **increases** from PRE condition to POST condition



Karunathilake et al. (2023) *Neural Tracking Measures of Speech Intelligibility*…, bioRxiv
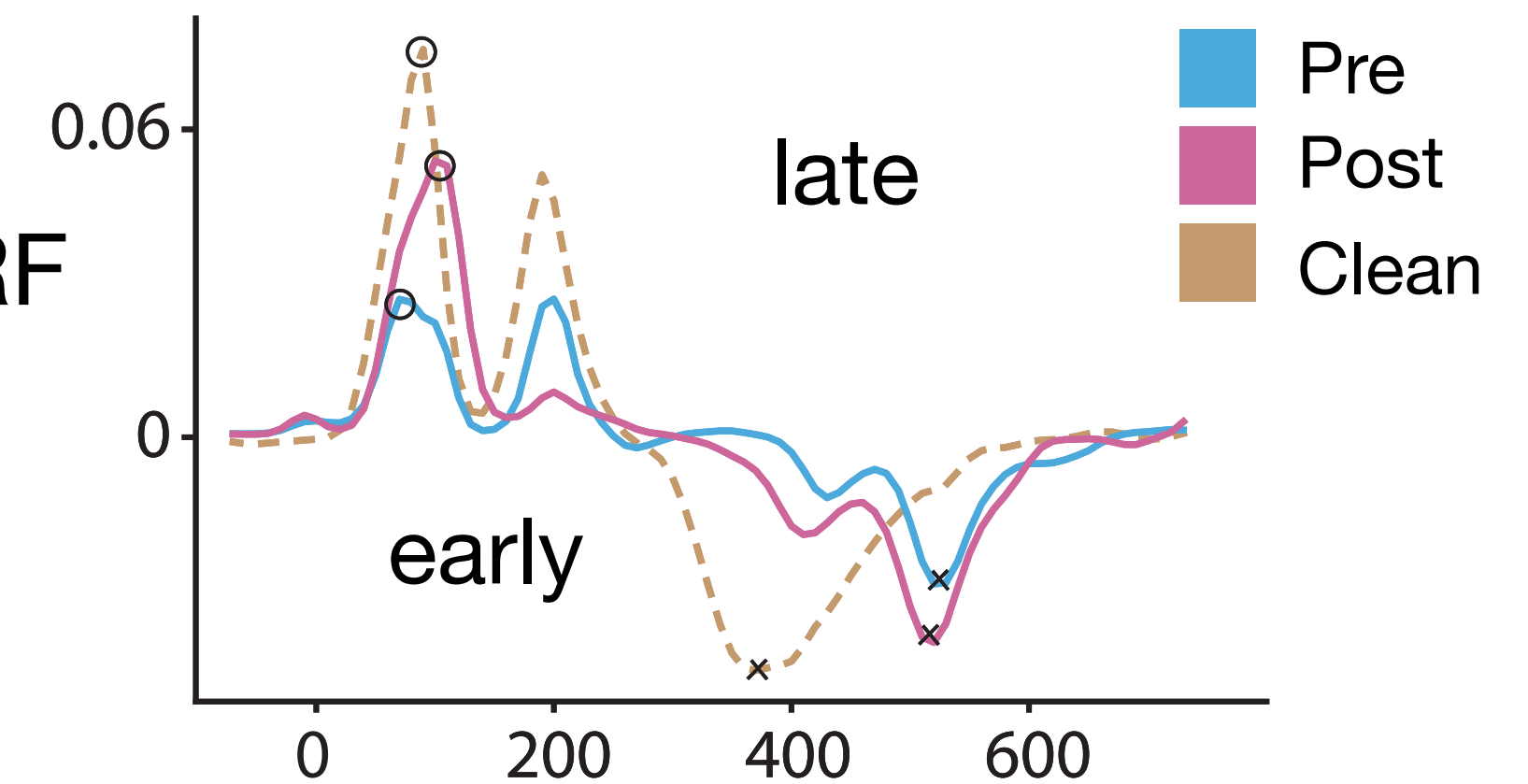
# Intelligibility Neural Results

- Word onset TRF shows both early (+) and late (-) processing stages

- Physiological response increases Pre→Post
  - Only in left hemisphere
  - Late processing stage shows larger change than early

- Physiological Word Onset response
  *Objective measure of intelligibility*
  - Acoustic responses: no change
  - Response to Word Surprisal: *Additional intelligibility measure*
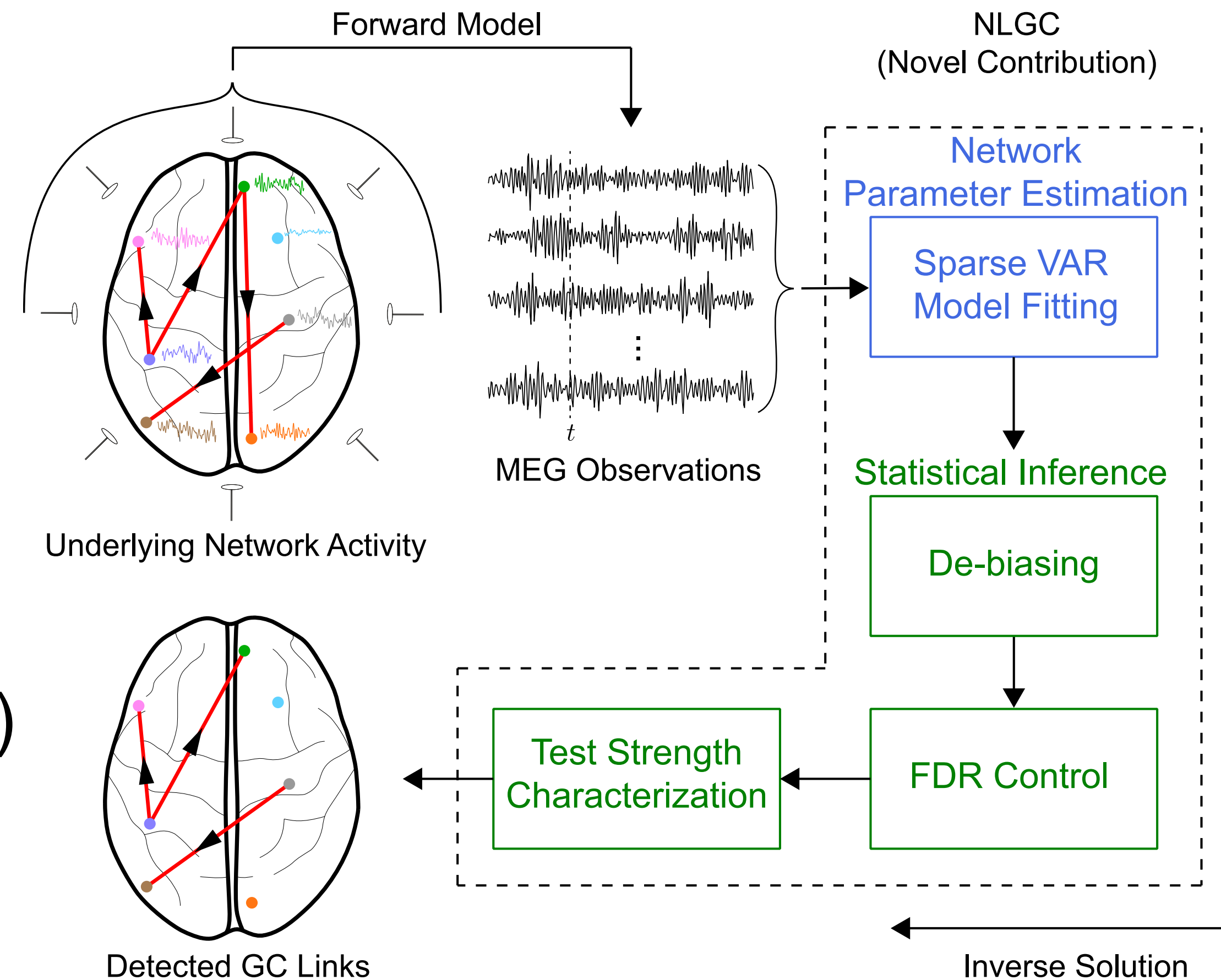
Word onset TRF

*single subject*

late

early

0.06

0

| Pre |
| Post |
| Clean |

0   200   400   600

Karunathilake et al. (2023) *Neural Tracking Measures of Speech Intelligibility...*, bioRxiv

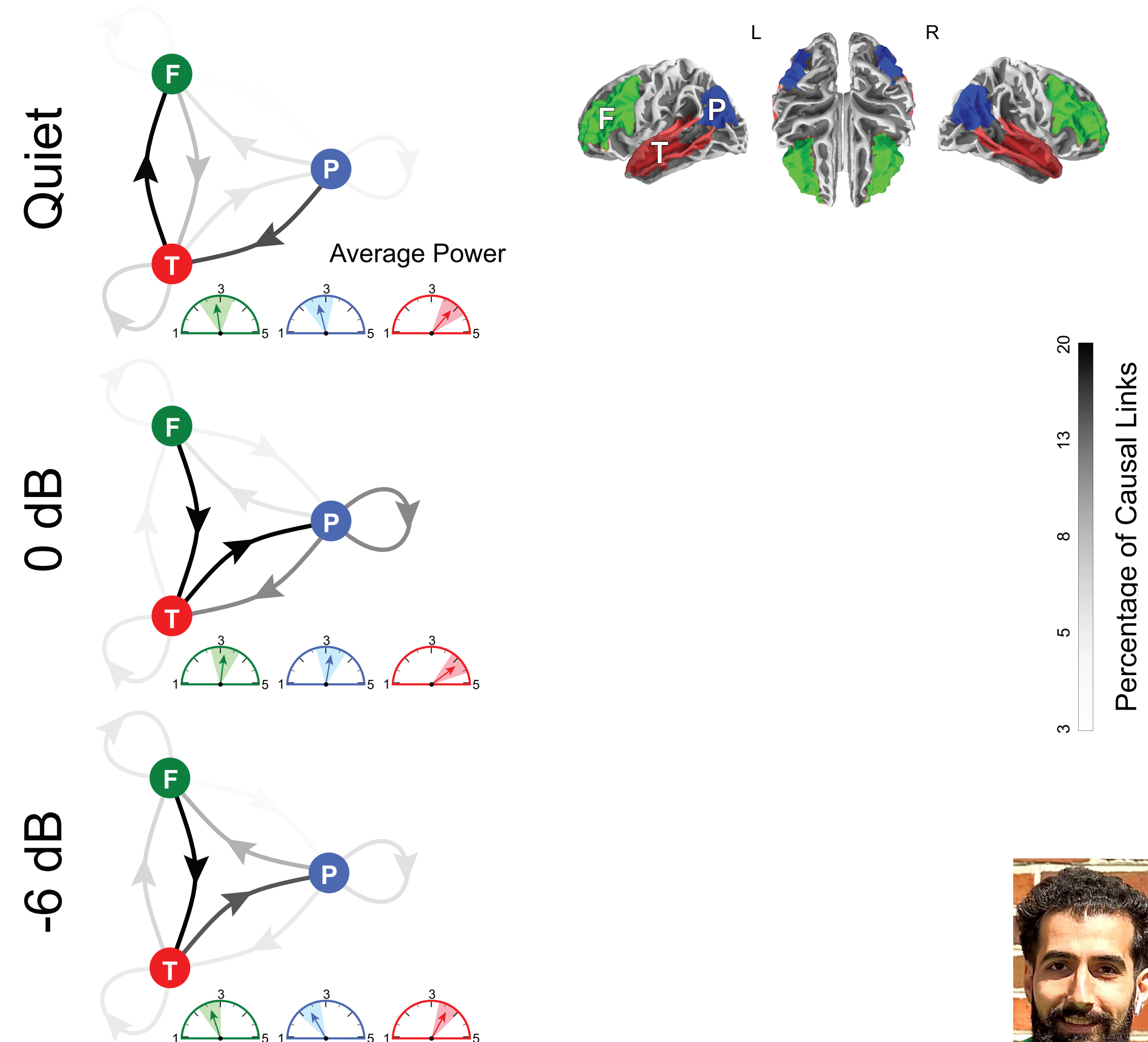# Outline

# Directional Functional Connectivity

- Novel method, based on Granger Causality (if source A can predict source B)

  - Directional (bi-directional allowed)

  - Localizes neural sources & GC link strengths simultaneously

  - source currents: latent sparse vector autoregressive (VAR) processes

- Network Localized Granger Causality (NLGC)

  - source spread & other biases minimized

  - robust against source model mismatch

  - parametrized by false discovery rate

  - intrinsically statistically robust



Forward Model

NLGC (Novel Contribution)

Network Parameter Estimation

Sparse VAR Model Fitting

MEG Observations

Underlying Network Activity

Statistical Inference

De-biasing

FDR Control

Test Strength Characterization

Detected GC Links

Inverse Solution

Soleimani et al. (2022) *NLGC: Network Localized Granger Causality with Application to …*, NeuroImage

# Cocktail Party Speech Results

Theta band example

- Speech in quiet connectivity: dominantly Temporal→Frontal and Parietal→Temporal

- Cocktail Party listening (moderate SNR): Temporal-Frontal switches direction; Parietal-Temporal now bi-directional

- Cocktail Party listening (poor SNR): Temporal←Frontal remains; Parietal→Temporal dominant



Soleimani et al. (2023) … *Cortical Directional Connectivity during Difficult Listening*… bioRxiv
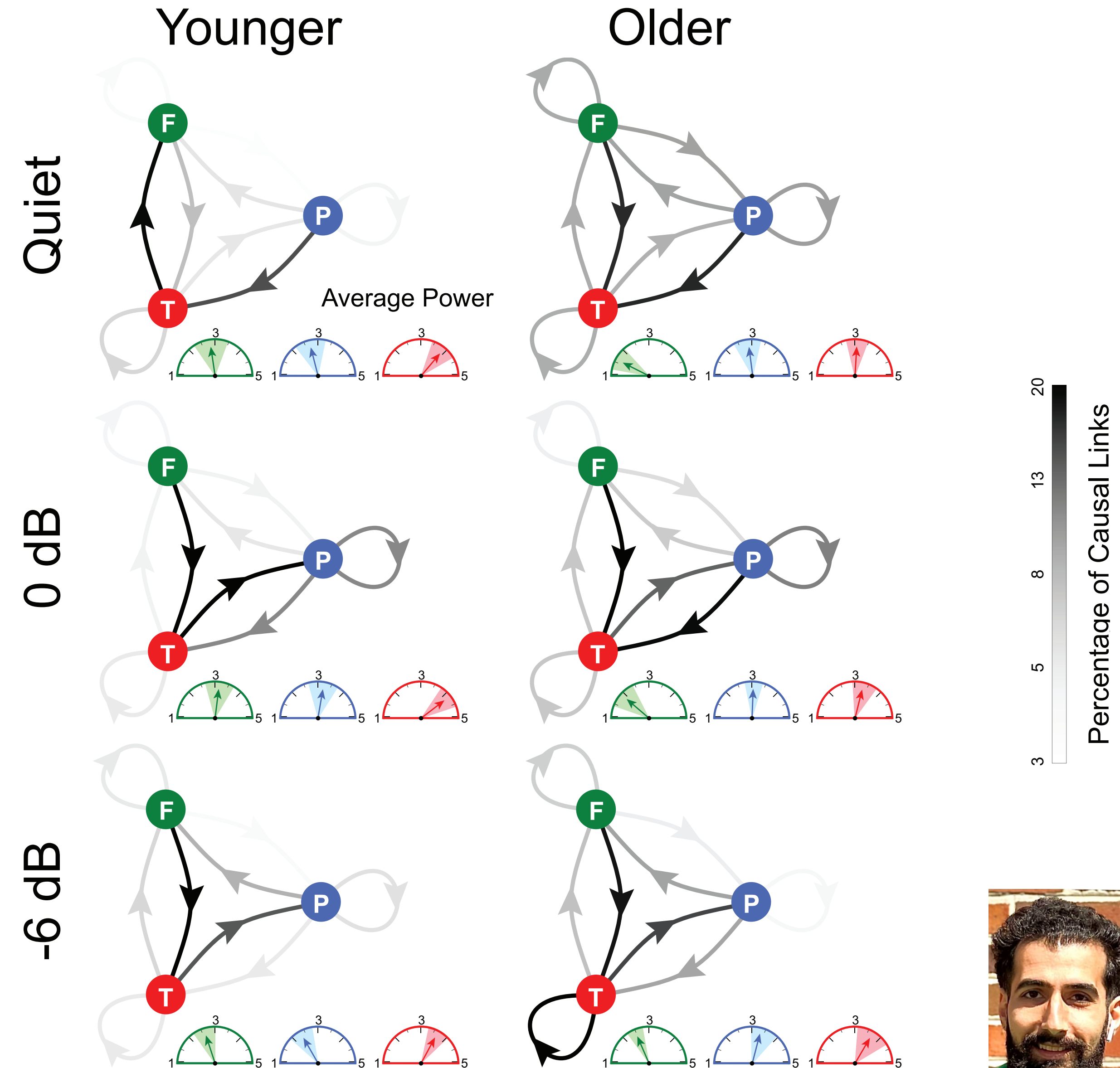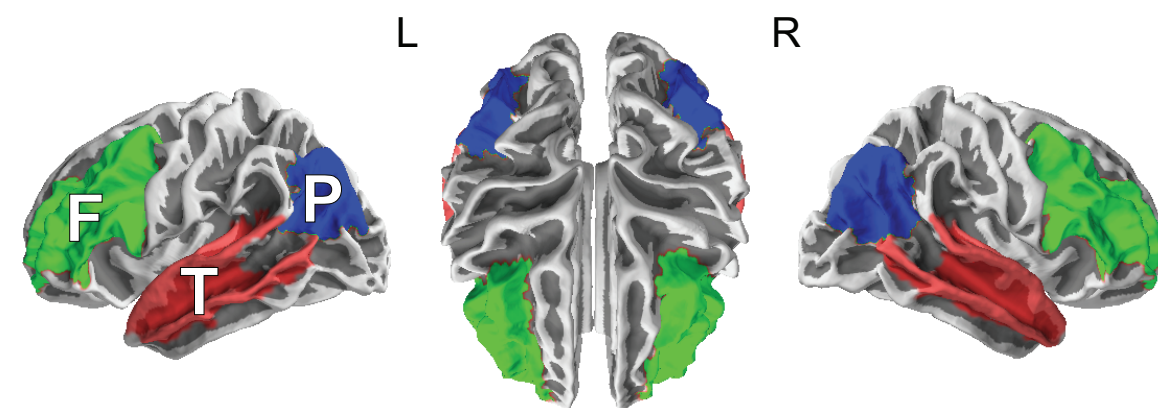
# Cocktail Party Speech Results

Older Listeners exhibit strongly different connectivity

- Older speech in quiet connectivity: similar to Younger cocktail party listening connectivity



Soleimani et al. (2023) … *Cortical Directional Connectivity during Difficult Listening…* bioRxiv

# Cocktail Party Speech Results

"Excitatory/Inhibitory" balance changes with task difficulty for Older Listeners only

- VAR (IIR filter) coefficients reveal neural signal transformation between sources
- coefficients > 0: "Excitatory"/facilitative
- coefficients < 0: "Inhibitory"/suppressive
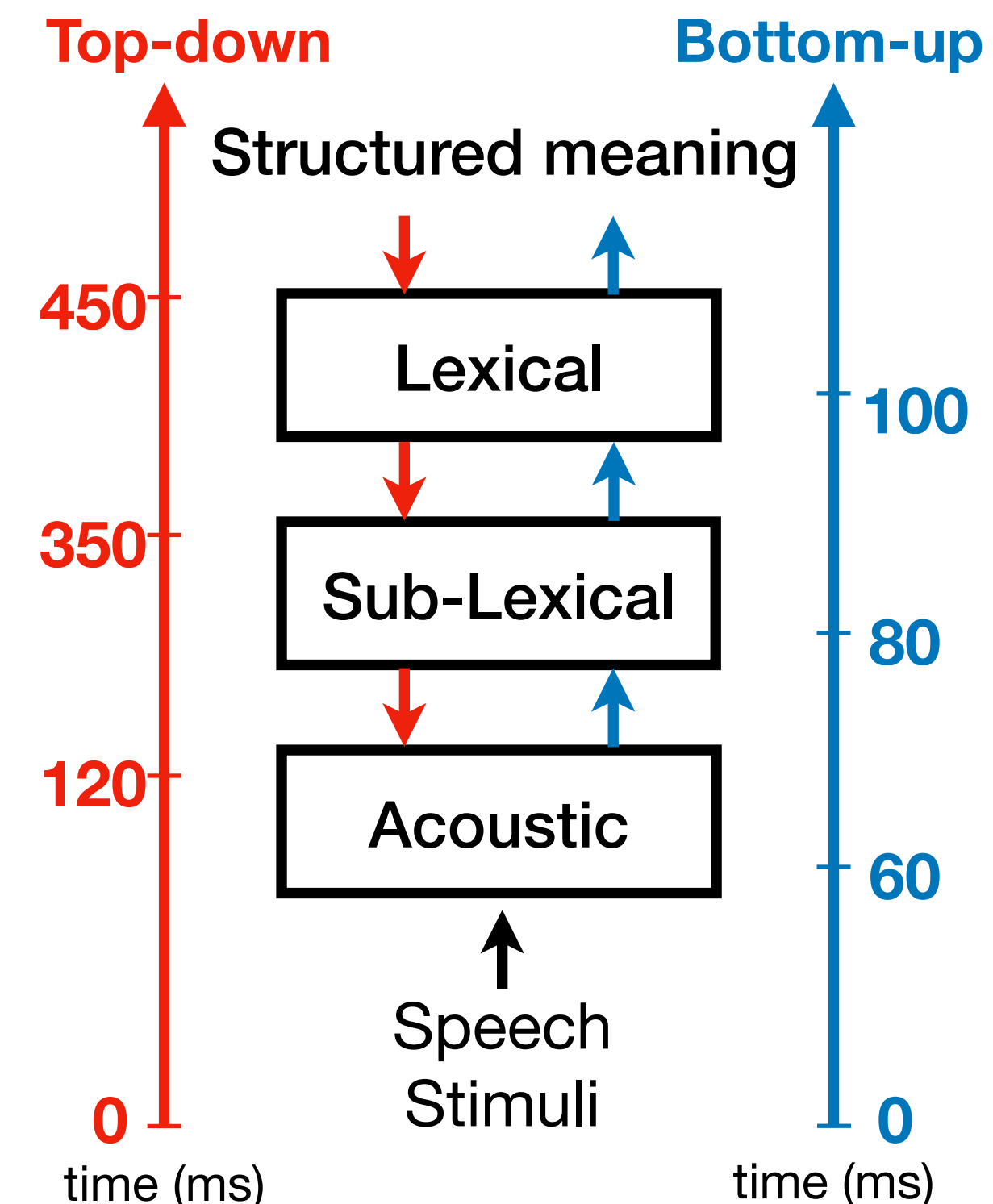- mixed coefficients: sharpening filter



Nature of the Links in Theta Band

Soleimani et al. (2023) … *Cortical Directional Connectivity during Difficult Listening*… bioRxiv

# Final Summary

*temporal patterns in **speech acoustics***
*temporal **neural** patterns* ⇄ *temporal patterns in **speech perception***
*temporal patterns in **language perception***
*temporal patterns in **understanding***

- Cortical responses time-lock to emergent features

- Higher level processing / top-down mechanisms may affect lower level

- Linguistic features processed only when linguistic boundaries intelligible

- Acoustic responses: bilateral but right lateralized; context-based responses strongly left lateralized

# thank you

These slides
available at:
ter.ps/simonpubs

Mastodon: @jzsimon@fediscience.org

http://www.isr.umd.edu/Labs/CSSL/simonlab