



**COMPUTATIONAL
AUDITORY NEURAL
SYSTEMS LAB**

Neural Representations of Same-Species Vocalizations in a Human Primate Model

Jonathan Z. Simon

University of Maryland

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

ARO Seminar Series, 26 March 2025



**COMPUTATIONAL
AUDITORY NEURAL
SYSTEMS LAB**

Neural Representations of Same-Species Vocalizations in a Human Primate Model

Jonathan Z. Simon

University of Maryland

humans

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

<https://cansl.umd.edu>
ARO Seminar Series, 26 March 2025



Speech
COMPUTATIO

NAL
AUDITORY
NEURAL

SYSTEMS
LAB

Neural Representations of Same-Species Vocalizations

in a Human
Primate Model

University of Maryland

humans

Jonathan Z. Simon

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

Charlie Fisher

Brooke Guo

Kevin Eguida

Ruwanthi Abeysekara

Michael Johns

Dushyanthi Karunathilake Craig Thorburn

London Dixon

Joshua Kulasingham

Shohini Bhattasali

Christian Brodbeck

<https://cansl.umd.edu>

ARO Seminar Series, 26 March 2025

Lab Members & Affiliates Karl Lerud

Vrishab Commuri

Thanks to

Faculty Collaborators Funding &

Support Samira Anderson (UMD)

Behtash Babadi (UMD)

Ellen Lau (UMD)

Philip Resnik (UMD)

Shihab Shamma (UMD)

Stefanie Kuchinsky (Walter Reed)

Elisabeth Marsh (Johns Hopkins)

Tom Francart (KU Leuven)

John Mosher (UTHealth)

L. Elliot Hong (UTHealth)



Wednesday, May 28, 2025

an e nsue or ysems esearc n . mon s co-recor o e

uory eura ysems aoraory . e s curreny e o e Maryland Magnetoencephalography Center, and director of the Computati
R01 grant "Multilevel Auditory Processing of Continuous Speech, fr

Auditory Cortex and Thalamus Seminar Series

From
12:00-1:00 PM EST

Auditory Neural Systems Laboratory
(CANSL). He is currently the PI of the
Acoustics to Language."
R01 grant "Multilevel Auditory

Processing of Continuous Speech,
frAcoustics to Language."

: Association for Research in Otolaryngology headquarters@aro.org

Subject: Join us for a NEW Seminar Series on the Auditory Cortex and Thalamus
Date: March 18, 2025 at 3:50PM

UPCOMING PRESENTERS

To: jzsimon@umd.edu

PAST

RECORDINGS

If you have missed any of the past Seminar
AUDITORY CORTEX AND THALAMUS

Series you can watch the full recordings
NOW on ARO's Official YouTube Channel!



These sessions are FREE to all. However, you must register to attend the webinar. Don't miss this

opportunity to engage and ask questions!

These sessions will be recorded for later viewing.

Join ARO and the Education Committee for an
*A special thank you to the ARO Education Committee for organizing
these enlightening Seminar Series on the
exciting sessions!*

Auditory Cortex and Thalamus .



Should you have any questions, please don't hesitate to
contact the ARO
Executive Office, at headquarters@aro.org or 615.432.0100. Thank you!

We are excited to invite you to the
first talk in the new series,

Association for Research in Otolaryngology | 5034A Thoroughbred Lane | Brentwood, TN

37027 **Corrected title**

US

“ Neural Representations of Same-Species

[Unsubscribe](#) | [Update Profile](#) | [Constant Contact Data Notice](#)

Vocalizations in a Human Primate Model”

presented by Dr. Jonathan Z. Simon.

UPCOMING PRESENTERS

Dr. Lori L. Holt

Univerity of Texas at

Austin Dr. Lori L. Holt

Wednesday, April 30,

2025 University

~~Univerity of Texas at~~

~~Austin~~ 12:00-1:00 PM EST

~~Wednesday, April 30,~~

~~2025~~ 12:00-1:00 PM EST

[REGISTER HERE](#)

[REGISTER HERE](#)

Dr. Ross Williamson

University

Univerisry of Pittsburgh

Dr. Ross Williamson

Wednesday, May 28,

2025 Univerisry of

~~Pittsburgh~~
~~12:00-1:00 PM EST~~
~~Wednesday, May 28,~~

~~2025 12:00-1:00 PM EST~~
~~REGISTER HERE~~

REGISTER HERE

Outline

- Auditory neurophysiology in animals *vs.* non-invasive neural recordings in humans — where is there common ground? ➔ here, human recordings = electroencephalography (EEG) & magnetoencephalography (MEG) ●
Neural processing of same-species-vocalizations *and* neural processing of speech

- ➔ speech as vocalization that is also a carrier for language
- Categorical perception & neural processing of elements of vocalization/speech

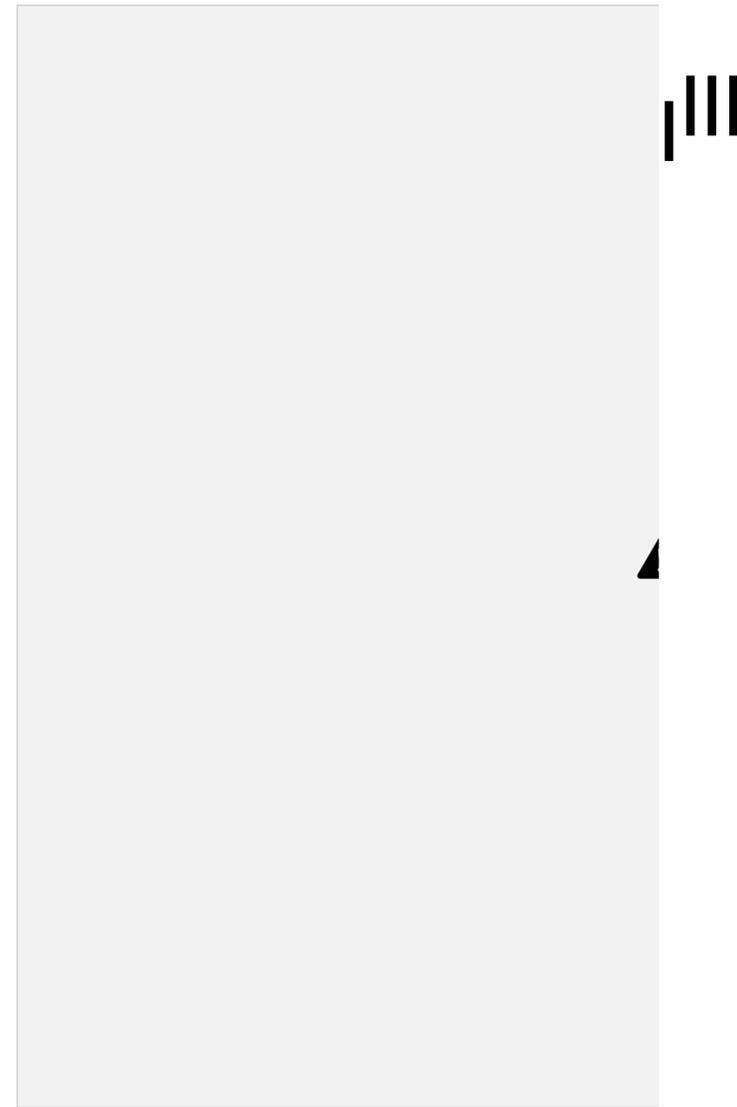
Mammalian Auditory Brainstem

bioRxiv preprint doi: <https://doi.org/10.1101/2022.10.14.512309>; this version posted Jan

available under a [CC-BY-NC-ND 4.0 International](#) (which was not certified by peer review) is the author/funder, who has granted bioRxiv a li

V

μV) Potential (



Butler & Lomber (2013) Shan et al. (2022) Time (ms)

Figure 2. The grand averaged broadband click-evoked ABR SEM (n=22). Waves I, III and V are annotated. All individual subject supplemental material **Figure S1.**

Brainstem Responses in Humans

(which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](#).

μV) Potential (
V

I III

Time (ms)

- typically a response to a punctate stimulus
- characterized by 3 robust peaks
- wave I: cochlear nerve
- wave III:

cochlear nucleus

- wave V: inferior colliculus (IC)

EEG

Figure 2. The grand averaged broadband click-evoked ABR waveforms. Shaded area shows ± 1

1

Auditory Brainstem Response (ABR)

SEM (n=22). Waves I, III and V are annotated. All individual subject responses are shown in supplemental material **Figure S1**.

bioRxiv preprint doi: <https://doi.org/10.1101/2022.10.14.512309>; this version posted January 4, 2023. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.
bioRxiv preprint doi: <https://doi.org/10.1101/2022.10.14.512309>; this version posted January 4, 2023. The copyright holder for this preprint bioRxiv preprint doi:

<https://doi.org/10.1101/2022.10.14.512309>; this version posted January 4, 2023. The copyright holder for this preprint

Auditory

Brainstem Responses in Humans

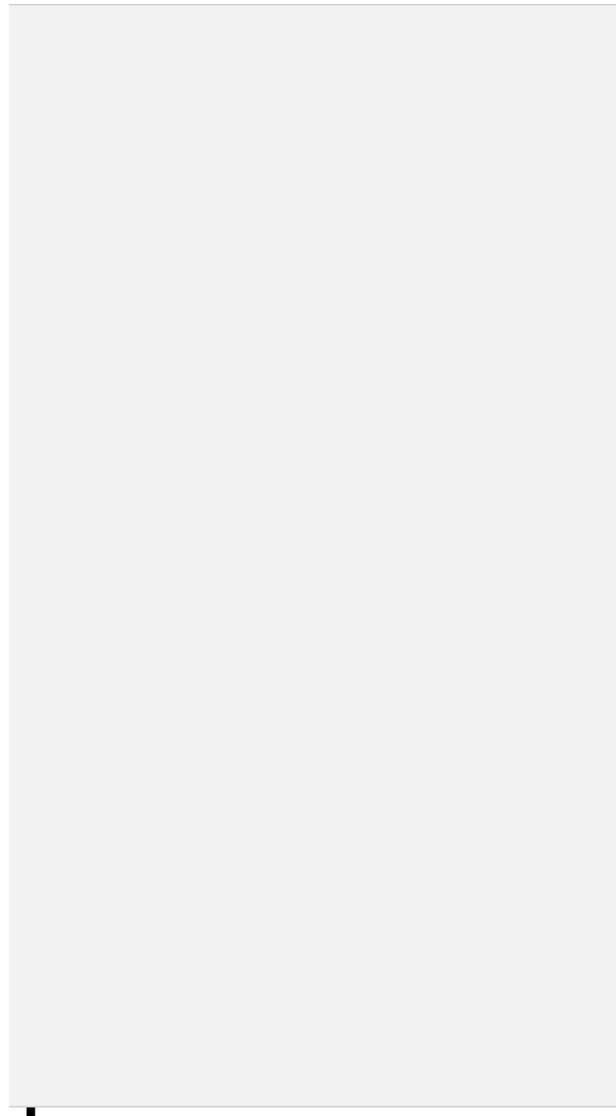
available under a [CC-BY-NC-ND 4.0 International license](#).

(which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made

(which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](#).

available under a [CC-BY-NC-ND 4.0 International license](#).

- **also** for continuous speech stimuli



μV) Potential (

function deconvolution (TRF) • of response with stimulus
 et al., 2014)
 • **still** characterized by 3 robust peaks • wave I: cochlear nerve
 • stimulus representation here: auditory nerve model (Zilany

• *temporal response* obtained by

Figure 5. General music- and speech-evoked ABR waveforms using the ANM as the regressor •

wave II: cochlear nucleus

in deconvolution. A. The grand averaged general music- and speech-evoked ABR waveforms.

Time (ms)

EEG

Time (ms)

Wave I, III and V are annotated. The waveforms were low passed with a cutoff at 1500 Hz. The

bioRxiv preprint doi: <https://doi.org/10.1101/2022.10.14.512309>; this version posted October 14, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

● wave V: inferior colliculus (IC)

shading areas show ± 1 SEM (n=22). B. Two examples of individual responses (subject 12 and

(which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

Figure 2. The grand averaged broadband click-evoked ABR waveforms. Shaded area shows ± 1

1

Figure 2. The grand averaged broadband click-evoked ABR waveforms. Shaded area shows ± 1

1

subject 18).

SEM (n=22). Waves I, III and V are annotated. All individual

subject responses are shown in

available under aCC-BY-NC-N

Figure 5. General music- and speech-evoked ABR waveforms using the ANM as the regressor

SEM (n=22). Waves I, III and V are annotated. All individual subject responses are shown in

supplemental material **Figure S1**.

supplemental material **Figure S1**.

in deconvolution. A. The grand averaged general music- and speech-evoked ABR waveforms.

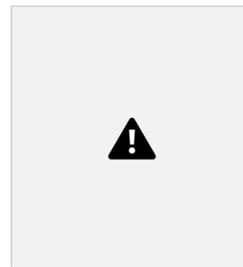
Shan et al. (2022) *Music and Speech Elicit Similar Subcortical Responses...* bioRxiv

Wave I III and V are annotated. The waveforms were low assed with a cutoff at 1500 Hz. The

Temporal Response Functions

Temporal
Response
Function (TRF)
Stimulus^s

Response^e
Stimulus^s ...
Response^e



Stimulus signal

Response signal:

TRF Model

Estimation & Fit

Estimated TRF

Estimated kernel

Linear kernel estimation:

Temporal Response Function (TRF) estimation:

Stimulus and response are known; find the best TRF

Stimulus and response are known; find the best linear kernel

to produce the response from the stimulus:

to produce the response from the stimulus:

Resp·

Stim·

Actual response

Resp·

Predicted response (Stimulus * kernel)

(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

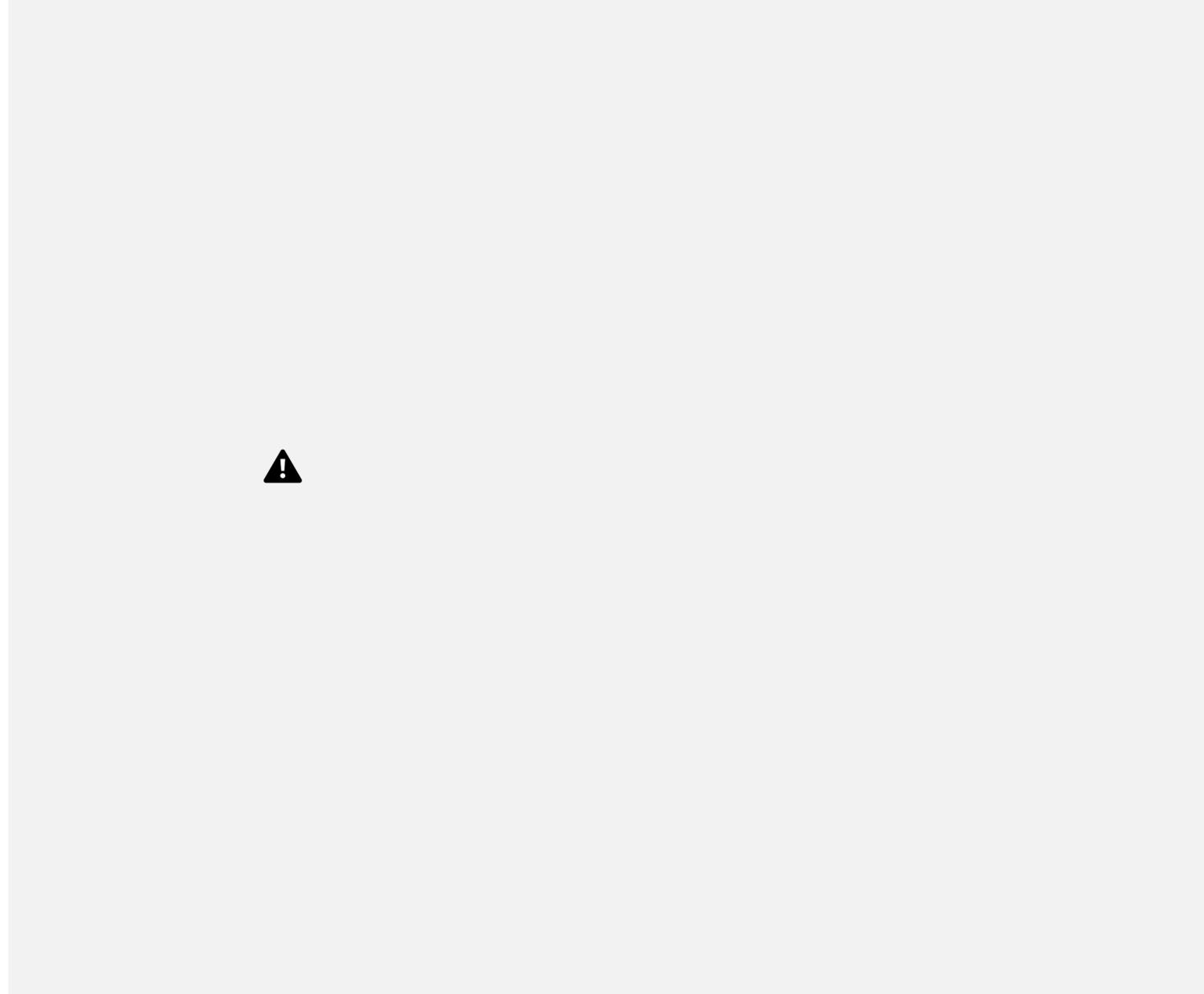
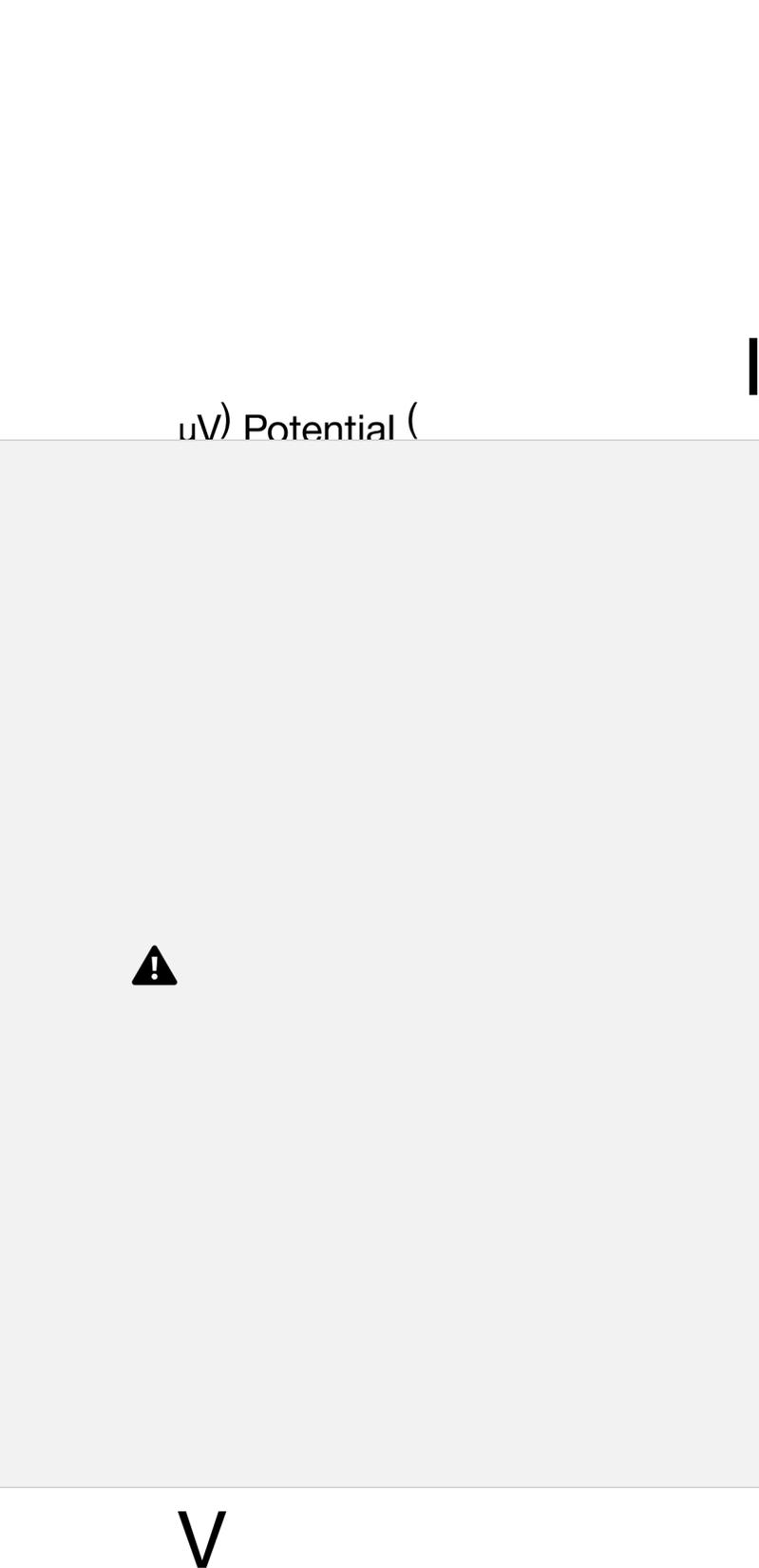
bioRxiv preprint doi: <https://doi.org/10.1101/2022.10.14.512309>; this version posted January 4, 2023. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.
bioRxiv preprint doi: <https://doi.org/10.1101/2022.10.14.512309>; this version posted January 4, 2023. The copyright holder for this preprint bioRxiv preprint doi:

<https://doi.org/10.1101/2022.10.14.512309>; this version posted January 4, 2023. The copyright holder for this preprint **Auditory**

Brainstem Responses in Humans

available under aCC-BY-NC-ND 4.0 International license.
(which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.
(which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.
available under aCC-BY-NC-ND 4.0 International license.

- **also** for continuous speech stimuli



- *temporal response function (TRF)*
- obtained by

deconvolution of response with stimulus

stimulus representation here: auditory nerve model (Zilany et al.,

2014)

● still

characterized by robust

peaks ● wave I: cochlear nerve

Figure 5. General music- and speech-evoked ABR waveforms using the ANM as the regressor ●

wave III: cochlear nucleus

in deconvolution. A. The grand averaged general music- and speech-evoked ABR waveforms.

Time (ms)

EEG

Time (ms)

Wave I, III and V are annotated. The waveforms were low passed with a cutoff at 1500 Hz. The

bioRxiv preprint doi: <https://doi.org/10.1101/2022.10.14.512309>; this ve

● wave V: inferior colliculus (IC)

shading areas show ± 1 SEM (n=22). B. Two examples of individual responses (subject 12 and

(which was not certified by peer review) is the author/funder, who has gra

Figure 2. The grand averaged broadband click-evoked ABR waveforms. Shaded area shows ± 1

1

Figure 2. The grand averaged broadband click-evoked ABR waveforms. Shaded area shows ± 1

1

subject 18).

SEM (n=22). Waves I, III and V are annotated. All individual

subject responses are shown in

available under a [CC-BY-NC-N](#)

Figure 5. General music- and speech-evoked ABR waveforms using the ANM as the regressor

SEM (n=22). Waves I, III and V are annotated. All individual subject responses are shown in supplemental material **Figure S1**.

supplemental material **Figure S1**.

in deconvolution. A. The grand averaged general music- and speech-evoked ABR waveforms.

Shan et al. (2022) *Music and Speech Elicit Similar Subcortical Responses...* bioRxiv

Wave I III and V are annotated. The waveforms were low assed with a cutoff at 1500 Hz. The

print doi: <https://doi.org/10.1101/2022.10.14.512309>; this version posted January 4, 2023. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made

bioRxiv preprint doi: <https://doi.org/10.1101/2022.10.14.512309>; this version posted January 4, 2023. The copyright holder for this preprint **Auditory**

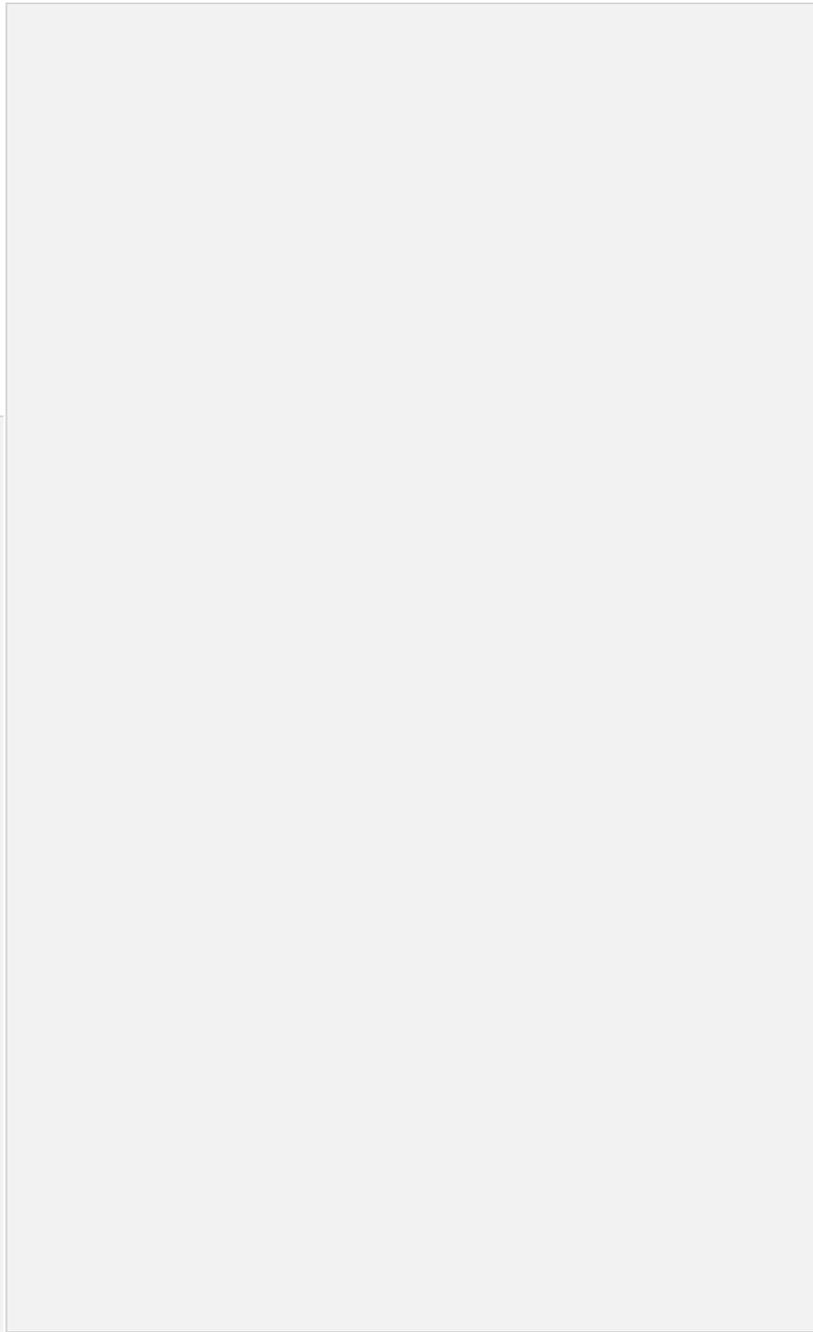
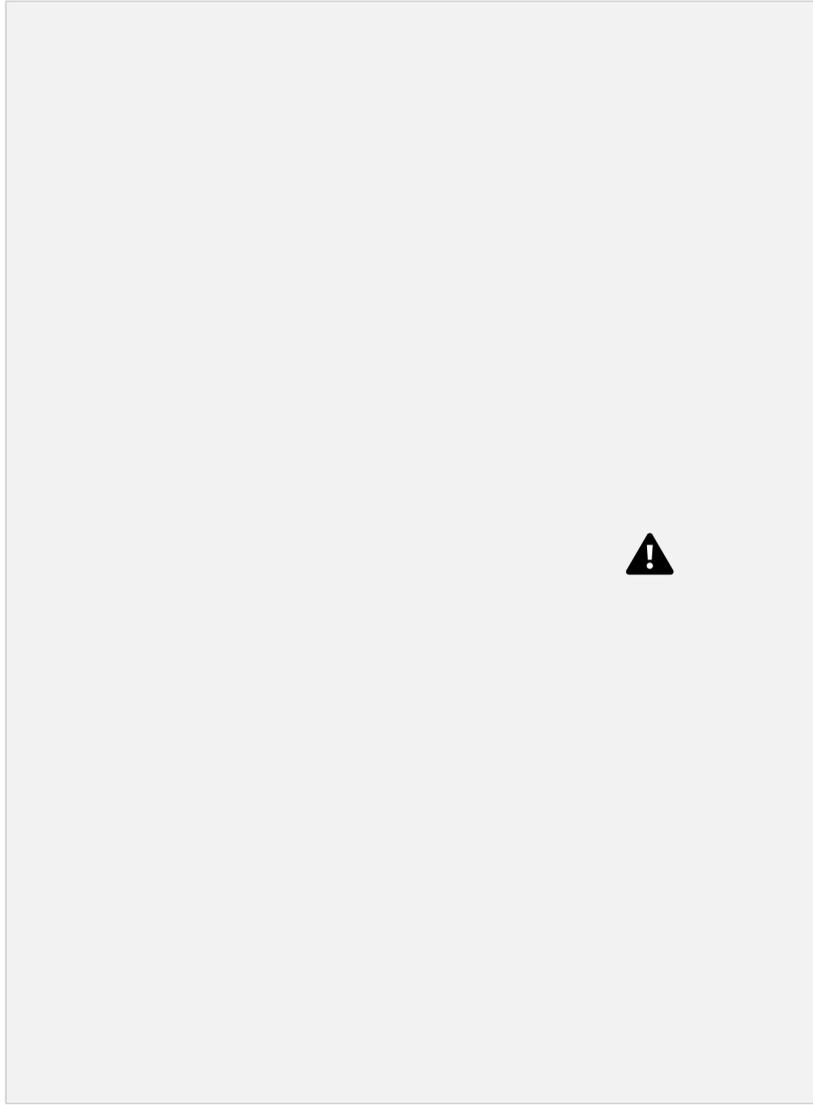
Brainstem Responses in Humans

available under a [CC-BY-NC-ND 4.0 International license](#).

(which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](#).

V

Magnitude (AU)



EEG

Figure 5. General music- and speech-evoked ABR waveforms using the ANM as the regressor

deconvolution. A. The grand averaged general music- and speech-evoked ABR waveforms.
Time (ms)

I, III and V are annotated. The waveforms were low passed with a cutoff at 1500 Hz. The shaded areas show ± 1 SEM (n=22). B. Two examples of individual responses (subject 12 and

subject 18).

1

Figure 5. General music- and speech-evoked ABR waveforms using the ANM as the regressor

SEM (n=22). Waves I, III and V are annotated. All individual subject responses are shown in supplemental material **Figure S1.**

in deconvolution. A. The grand averaged general music- and speech-evoked ABR waveforms. Wave I III and V are annotated. The waveforms were low assed with a cutoff at 1500 Hz. The

Thalamic Response in Humans Butler &

Lomber (2013)

Middle Latency Response (MLR)



Middle Latency Response
(MLR)

Cortical Responses

≈ Auditory Thalamus [Medial
Geniculate Body]
EEG

Polonenko & Maddox (2021) *Exposing distinct subcortical components...* eLife

-.

-1.0

-.
-1.0

Thalamus & Brainstem in Humans

N1

-100 0 100 200 300 400 500 600 700 -100 0 100 200 300 400 500 600 Time (ms)

Time (ms)

s and

tem^{ABR} s and

s and

ABR

tem^{ABR} tem^{ABR}

ABR

V

8

V

V

V

6

MLR

Target: Quiet

Target: Easy

Target: Hard

8
6
4
2
0

-2
|
|||
8

8
4
6
6
2
4
4

0
2
2
-2
0
-4
-2
-2

Pa

Na

-10 -4

-4

0 10 20 30 40 50 60

70 Na

Na
Na

Pa

■

Pa

Pa

Pb

~~MLR~~

■

■

Pb Pb

Pb

Distractor:
Easy

Distractor:
Hard

Time (ms)

EEG

-4

Lerud et al. (2025) *Continuous and Concurrent Auditory TRFs...* ARO Poster

J.P. Kulasingham, C. Brodbeck and A. Presacco et al.

Thalamo-cortical Response in Humans

TRF of 70-100 Hz speech envelope

MEG

. Kulasingham, C. Brodbeck and A. Presacco et al. NeuroImage 222 (2020) 117291

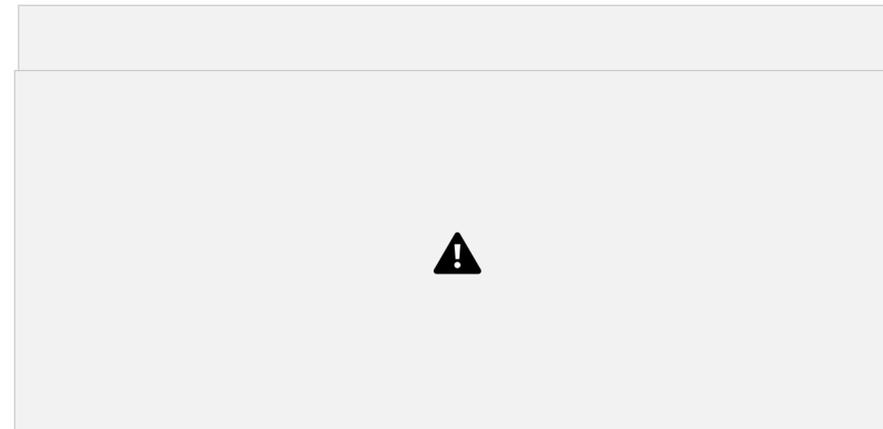
TRF of 70-100 Hz speech carrier

40 ms peak latency

⇒ primary 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



Thalamo-cortical Response in Humans

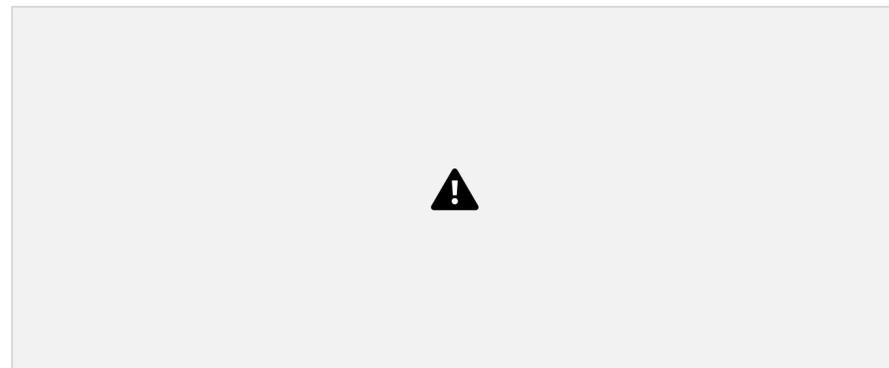
Attend Male

2×10^{-4}

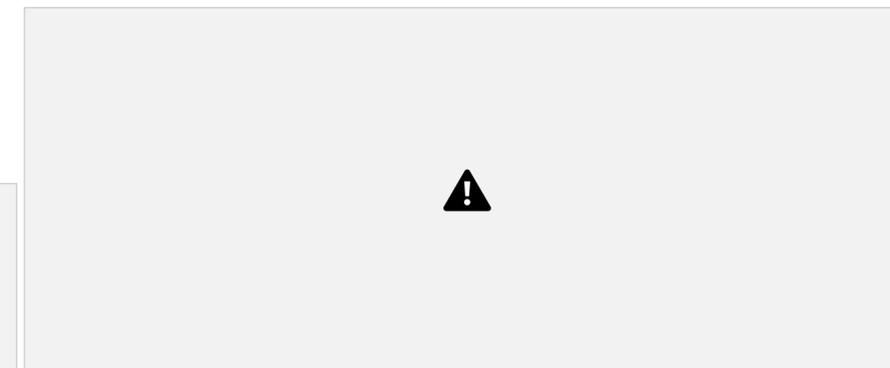
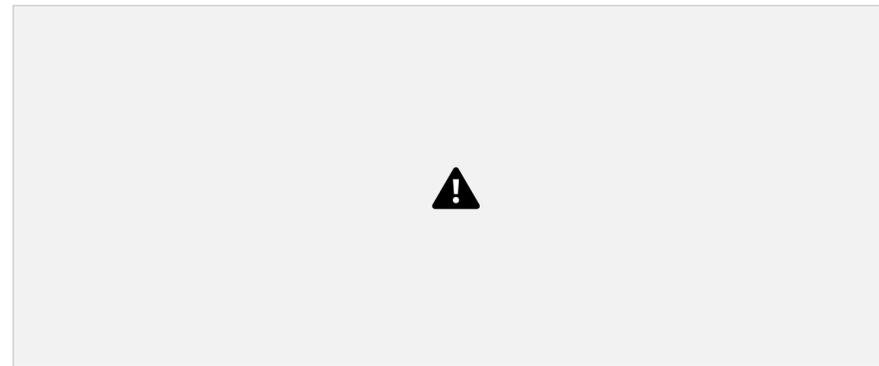
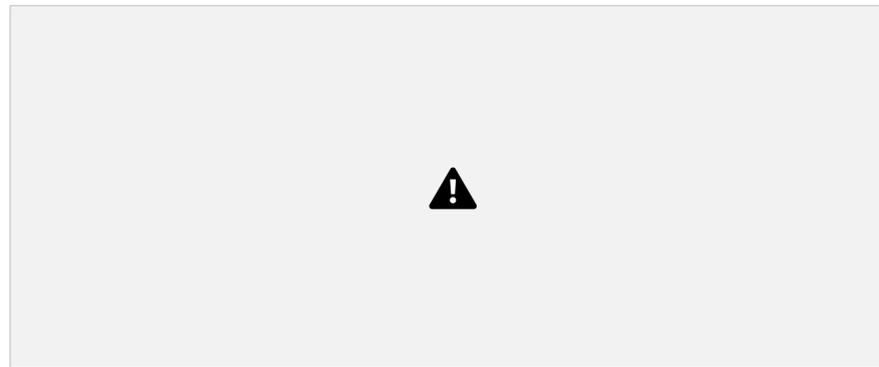
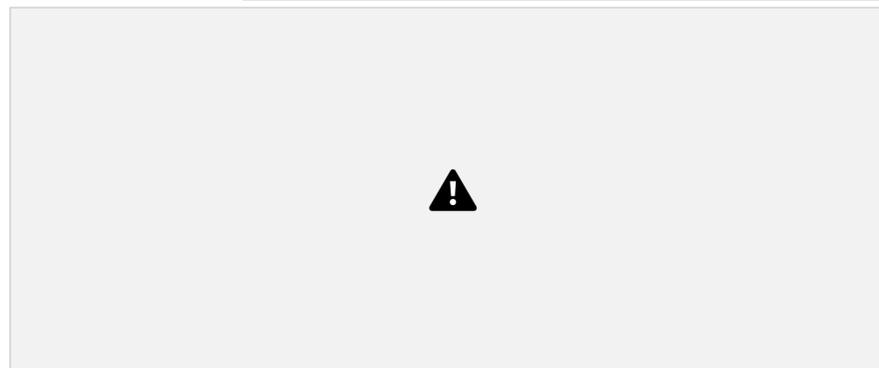
TRF Amplitude (std units)

0

2×10^{-4} 0



Ignore Male



Primary cortex modulated by selective attention

Attend > Ignore

Commuri et al. (2023) ... *High-Gamma Band Depend on Selective Attention*, Front Neurosci

Outline

- Auditory neurophysiology in animals vs. non-invasive neural recordings in humans — where is there common ground? ➔ here, human recordings = electroencephalography (EEG) & magnetoencephalography (MEG) ●
Neural processing of same-species-vocalizations *and* neural

processing of speech

➔ speech as vocalization that is also a carrier for language

- Categorical perception & neural processing of elements of vocalization/speech

Vocalizations & Categorical Perception

- Vocalizations, including speech, are often perceived categorically

b

Proportion of trials reported as the same $e^{100806040}_{20}$ 'dad' ● Even in rhesus

monkeys

'bad'

Tsunada et al. (2011)

0 20 60 80 100

0 40

Test-stimulus morph (%)

Bizley & Cohen (2013)

Vocalizations & Categorical Perception

- Categorical perception adds robustness to communications

- Consequently, categorical perception is also a robust percept

h e l p h e l p h e l p h e l p Dilley & Pitt (2010)



Vocalizations & Categorical Perception

- Categorical perception adds robustness to communications
- Consequently, categorical perception is also a robust percept

h e l p h e l p h e l p h e l p Dilley & Pitt (2010)

Cortical Responses to Phonemes in Monkey

- Cortical neurons in anterolateral belt (ALB) respond *categorically* to phonemes. Tsunada et al. (2011)

b Proportion of trials reported as the same $e^{100806040}_{20}$

*'dad'*²⁰⁰

0% morph (bad) 60% morph 20%
morph 80% morph 40% morph
100% morph (dad) 50% morph

Firing rate (Hz)

'bad'

0 20 60 80 100

0 40

Test-stimulus morph (%)

150

100

50

0

Time from stimulus onset (ms)

0 500

Bizley & Cohen (2013)

Cortical Responses to Phonemes in Humans •

How does one separate human cortical responses to phonemes

from cortical responses to the sounds of phonemes? ●
Multivariable regression in the time-domain 

- multivariable-Temporal Response Functions

(mTRFs) Gammatone



**acoustic
features**

- Gammatone
- Gammatone
- Envelope
- Envelope
- Phoneme
- Surprisal
- Gammatone
- Gammatone
- Envelope Onset
- Envelope Onset

Phoneme
Onset



Envelope

Gammatone Envelope Onset

Phoneme Onset

*

*

*

.

signals

Measured

Neural signals

Measured Neural

*

*

Entropy

Phoneme Surprisal

*

Cohort

Cohort

*

Measured Neural signals

Measured Neural signals

phonemic features

- Onset
- Onset
- Phoneme
- Phoneme
- Surprisal

- Surprisal
- Unigram
- Surprisal
- Cohort
- Cohort
- Entropy
- Entropy
- GPT2
- Surprisal
- Entropy

- Phoneme
- Phoneme
- Onset
- Word

Word Onset

Unigram Surprisal

GPT2 Surprisal

*

*

.

*

.

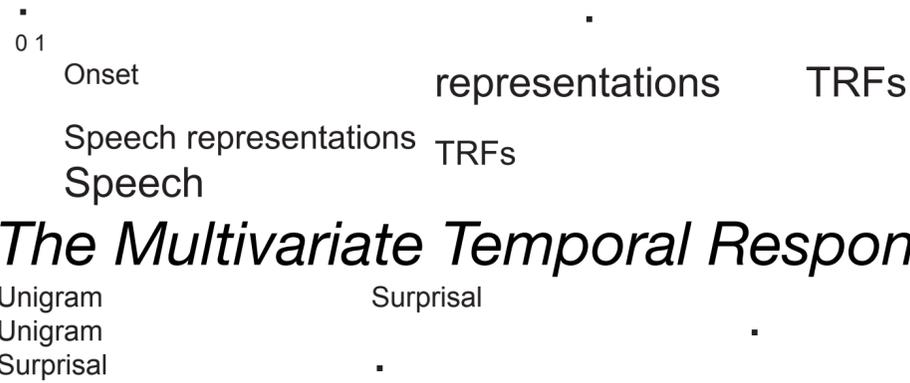
.

Predicted-Neural signals Predicted Neural signals

Word
Speech Representations

TRFs

Word Onset 0 1



Crosse et al. (2016) *The Multivariate Temporal Response Function (mTRF) Toolbox ...* ; Front Hum Neurosci

Brodbeck et al. 2023 *Eelbrain: A Python Toolkit for Time-Continuous Analysis ...* , eLife

Predicted Neural signals

Predicted Neural signals

Further Disentangling Phonemes

- Phonemes, while not identical to their underlying acoustics, are still strongly correlated with their underlying acoustics
- even mTRFs have trouble when predictors are too correlated
- Are there phoneme measures could we use that are less

correlated with the acoustics?

- Yes! based on linguistic statistical distributions:
 - phoneme surprisal
 - phoneme cohort entropy
- Also, might learn about neural processing of these measures

Surprisal Surprisal

Number of times a word that starts with
this

K EY M ...

Phoneme

sequence
occurs in
SUBTLEX

K EY ...
52908
(90 words)

Number of
words that
start with
this sequence

SUBTLEX:
23875 (45%) (4 words)

K EY **S** ...
16048 (30%) (13 words)

K EY **K** ...
2598 (5%) (3 words)

KEY N ...

1337 (3%) (13 words)

...

“came”, “Cambridge”, ...

“case”, “cases”, “caseworker”,
“casein”, ...

“cake”, “caked”, “cakes”

“cane”, “canine”, “Canaan”,
“Kane”, “Kynesian”



$$\sum_{word \in cohort_{i-1}} \frac{1}{|cohort_{i-1}|} \sum_{word \in cohort_{i-1}} freq_{word}(i-1)$$

51 million words
movie subtitle database

$$surprisal_i = -\log_2 \frac{freq_{word}(i)}{\sum_{word \in cohort_{i-1}} freq_{word}(i-1)}$$

Cohort Entropy

Cohort entropy

- How unpredictable is the current word?

L EY K ... K EY K ... B EY K ...

lake
(95%) Entropy

lakes (5%)
cake

(88%) cakes (11%)
caked (1%)
baker (29%)
bacon (25%)
baked (14%)
bake
(14%)



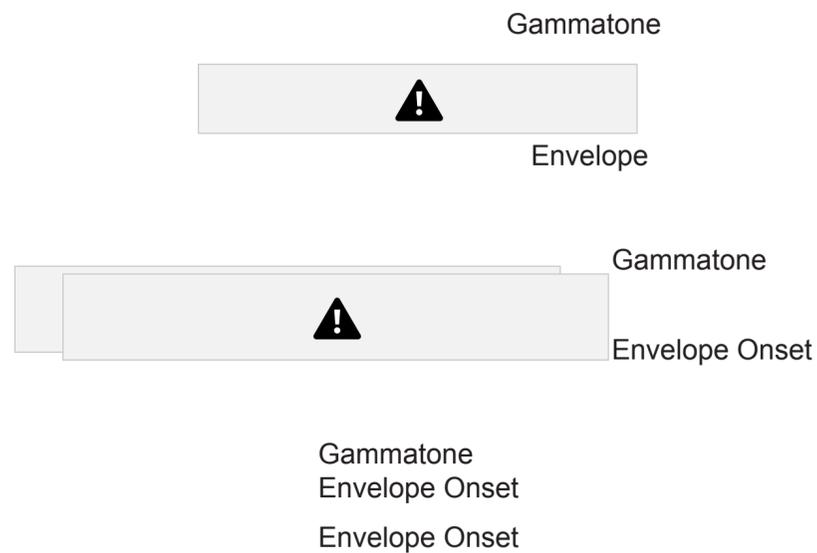
$$H_i^{cohort} = - \sum_{word \in cohort_i} p_{word} \log_2 p_{word}$$

Cortical Responses to Phonemes in Humans •

How does one separate cortical responses to phonemes from

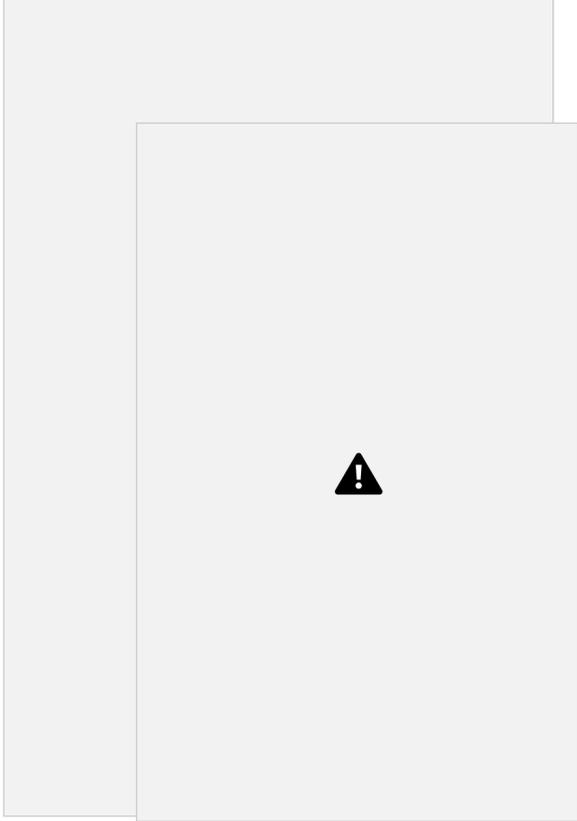
cortical responses to the sounds of phonemes? • Multivariable time-domain regression: 

- multi-Temporal Response Functions (mTRFs)



acoustic features

- Phoneme
- Onset
- Gammatone
- Gammatone
- Envelope
- Envelope
- Phoneme
- Surprisal
- Gammatone



Gammatone

Envelope

Gammatone Envelope Onset

Phoneme Onset

*



*

*

.

signals

.

.



.

Measured

Neural signals

Measured Neural

.



Entropy

Phoneme

Surprisal

*

*

*

Phoneme
Phoneme

Cohort

Cohort

*

Measured Neural signals

Measured Neural signals

Onset
Word

Surprisal

Unigram

Surprisal

Cohort

Cohort

Entropy

Entropy

GPT2

Surprisal

Entropy

Word

Onset

Unigram Surprisal

GPT2

Surprisal

*

*

.

*

.

.

phonemic features

Onset

Onset

Phoneme

Phoneme

Surprisal

Word

Speech

Representations

TRFs

Word

Onset

0 1

Predicted-Neural
signals Predicted
Neural signals

Onset
0 1

Unigram
Unigram
Surprisal
Surprisal
Speech representations
Speech representations
TRFs
TRFs

•
•
Predicted Neural signals Predicted Neural signals

Study Experimental Design

Speech-envelope Modulated
Noise

Scrambled words Narrative

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



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



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

Cortical Responses to Speech Acoustics in Humans

acoustic envelope onsets acoustic envelope + -

+

Scrambled
Narrative

0.1

Noise

Non-word

0.06

MEG

0 200 400 600 a.k.a. “speech tracking”

0 200 400 600

~60 ms: acoustic bottom-up processing

~120 ms: acoustic but attention-dependent

based STRFs are used to model the intertrial variability of the LFP. *B*: correlation between the

shape of the LFP. *B*: correlation between the

shape

shape

Are Human Cortical Latencies “Long”?

ance of the LFP. *B*: correlation between the shape

measures the similarity of tuning across

of STRFs measures the similarity of tuning across

of STRFs measures the similarity of tuning across

nals. Delta-, theta-, and alpha-variance neural signals. Delta-, theta-, and alpha-variance
neural signals. Delta-, theta-, and alpha-variance
e highly correlated, and the higher fre STRFs are highly correlated, and the higher fre
STRFs are highly correlated, and the higher fre
nds (gamma, high gamma, MUA) also quency bands (gamma, high gamma, MUA) also
quency bands (gamma, high gamma, MUA) also
luster of similarity to each other. *C*: **LFP-based STRFs (ferret A1)**
show a cluster of similarity to each other. *C*:
show a cluster of similarity to each other. *C*:
TRFs in each row are measured for the example STRFs in each row are measured for the
example STRFs in each row are measured for the

- **A note for auditory neurophysiologists**

rding site but using different LFP bands. same recording site but using different LFP bands.
same recording site but using different LFP bands.
indicate an increase in the neurophysi Red areas indicate an increase in the neurophysi
Red areas indicate an increase in the neurophysi
gnal following an increase in power of ological signal following an increase in power of

- **120 ms latency is not as “crazy late” as it**

ological signal following an increase in power of

ponding spectro-temporal stimulus fea the corresponding spectro-temporal stimulus fea

the corresponding spectro-temporal stimulus feature, and blue areas indicate a decrease. STRFs are normalized to have the same maximum absolute value.

might seem

normalized to have the same maximum absolute value. LFP variance STRFs are generally inhibitory in the alpha and beta bands, variable (inhibitory or excitatory) in the gamma band, and excitatory in the high gamma band. MUA STRFs are generally excitatory. The peak latency of the LFP

- **Even in primary auditory cortex (A1) of**

ferret, spectro-temporal receptive fields

ferret, spectro-temporal receptive fields

generally excitatory. The peak latency of the LFP

TRF is later in the alpha and beta bands

variance STRF is later in the alpha and beta bands

variance STRF is later in the alpha and beta bands

high gamma band. The frequency tun

than in the high gamma band. The frequency tun

(STRFs) made with speech stimuli

than in the high gamma band. The frequency tun

variance STRFs is generally similar to

ing of LFP variance STRFs is generally similar to

ing of LFP variance STRFs is generally similar to

MUA STRF but usually has additional

that of the MUA STRF but usually has additional

that of the MUA STRF but usually has additional

like the other signals, STRFs for mean

peaks. Unlike the other signals, STRFs for mean

have peaks with latency >100 ms

peaks. Unlike the other signals, STRFs for mean

peaks at multiple latencies. These

LFP show peaks at multiple latencies. These

LFP show peaks at multiple latencies. These

pically show an early peak (25 ms

STRFs typically show an early peak (25 ms

when made with Local Field Potential

STRFs typically show an early peak (25 ms

with negative polarity indicating depolar

latency) with negative polarity indicating depolar

latency) with negative polarity indicating depolar

l

t

l

i

o

v

w

e

e

p

d

e

b

a

y

k

a

(

p

\$

o

6

s

0

i

m

1

ization, followed by a positive peak (60 ms

ization, followed by a positive peak (60 ms

(LFP), not spikes

indicating hyperpolarization, and some (latency) indicating hyperpolarization, and 150
some

or longer-latency peaks.

(latency) indicating hyperpolarization, and some times other longer-latency peaks.

times other longer-latency peaks.

IANC^E

Ding et al. (2016) *Encoding of Natural Sounds by Variance of the Cortical Local Field Potential* J Neurophysiol-

Phonemic Responses in Humans

phoneme onset
phoneme

surprisal cohort



0 200 400 600

- Clear evidence of phoneme-driven responses, uncorrelated with acoustics
- Evidence of categorical neural processing of vocalization (speech)
- Low-level phoneme processing at **~80 ms** (not much later than 60 ms)
- Additional later processing at **~350 ms** with negative polarity

N400-like, associated with predictive coding (Eddine et al., 2024)

Karunathilake et al. (2025) *Neural Dynamics of the Processing of Speech Features ... J Neurosci*

Beyond Phonemes

- In human speech, phonemes building blocks of words
- Words and groups of words are used to convey meaning
- Animal vocalizations are often used to convey meaning

Vocalizations Convey Meaning

- In rhesus monkeys, some vocalizations transmit information regarding food quality

low-quality: “grunt”

high-quality: “harmonic arch” or “warble”

Grunt Harmonic arch Warble Baseline

Bizley & Cohen (2013)

Responses to Meaningful Vocalizations

~~Neurons in monkey ventral prefrontal cortex (VPFC) respond~~
categorically based on
meaning, not acoustics

- VPFC neurons encode transitions

between calls of different
abstract categories

~~• VPFC neurons do not encode~~

Bizley & Cohen (2013)

transitions between acoustically
distinct stimuli transmitting the

most Noise Exemplars reliably differ
exemplars, the mean z-score value was not

possibility that vPFC neurons transitions between stimuli that are in different classes (Ulanovsky, Las, &
different than zero ($p > .05$). DISCUSSION

same information

Gifford et al. (2005) *The neurophysiology of functionally meaningful categories ...* J Cog Neurosci

Humans Cortical Responses to Words

- Words often convey meaning in human speech



word-level features

envelope

Gammatone

Envelope

spectrogram onset

Gammatone

Envelope Onset

spectrogram

phoneme

Phoneme

onset

Onset

phoneme

Phoneme

surprisal

Surprisal

cohort

Cohort

entropy

Entropy

word

Word

Onset

onset

surprisal

Unigram

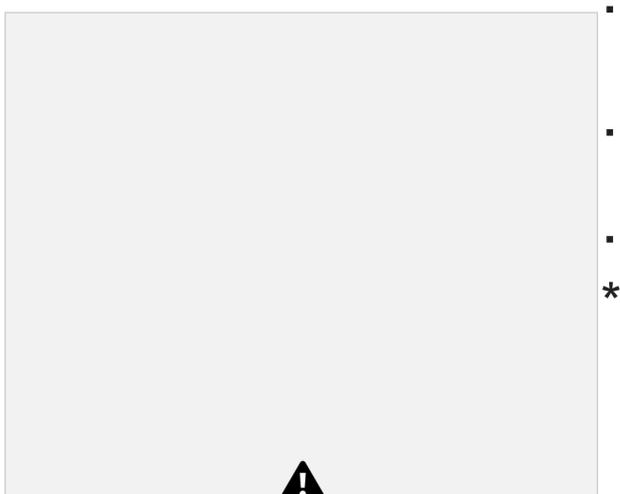
(no context) Surprisal

surprisal

GPT2

Surprisal

(GPT-2 model)



*

*

▪

▪

▪

0.1

Measured Neural signals Predicted Neural signals

Speech Representations TRFs^{TRFs}

Speech representations

Word Onsets

catalogue inner

eye

The cat a log in a lie **Do we...**

cattle

login

library

- ▶ Anticipate word boundaries based on context?
- ▶ Infer them later based on consistency?

eye

catalogue inner

R T The cat a log in a lie

login

library

cattle

Figure 1. Recognition of the phrase “The catalogue in a library,” as spoken by speaker of British English:

“The catalogue in a library”

[/ðəkætəlɒɡɪnələɪbrɪ]. The upper panel shows the competitive inhibition process that occurs among activated candidate words in an interactive-activation model, such as Shortlist A. Words that compete for the same stretch of input inhibit each other via direct, bidirectional inhibitory connections. Only a subset of the best-matching candidates is shown. The lower panel illustrates the path-based search through a word lattice used in automatic

Norris & McQueen, 2008

Word Surprisal (without context)

Frequency of words based on SUBTLEX

the

to

and

of

in

a

.

.

.

.

Word Surprisal (contextual)

yourself

a

out

it

that

one

if you happened to find
the

your

.
. .
. .

(via GPT-2)

Word Responses in Humans

onset

+

-

Non-word
Scrambled
Narrative

surprisal
(without
context)

+

-

MEG

+

0.2

word

0.1

Noise

- Clear evidence of word-driven responses, uncorrelated with acoustics
- Evidence of categorical neural processing of vocalization (speech) •
- Low-level phoneme processing at **~100 ms** (not much later than 80 ms) •
- Additional N400-like processing at **~450 ms**, c.f. predictive coding
(Eddine et al.,

2024) Karunathilake et al. (2025) *Neural Dynamics of the Processing of Speech Features ...* J Neurosci

Contextual Word Surprisal Results

(without context) word surprisal

0.15 + +
word surprisal

(contextua
l)

Noise

+
0.15 + +

Narrative

—
Non-word
Scrambled

—
MEG

600 0 200 400 600 0 200 400 600

0 200 400 600 0

- Context-based surprisal is more robust than naive surprisal
- N400 like response in both predictors, c.f. predictive coding

(Eddine et al., 2024)

Karunathilake et al. (2025) *Neural Dynamics of the Processing of Speech Features ...* J Neurosci

Neural Speech Processing Progression

Top-down **Bottom-up** Structured meaning

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

- Phonemic and word-based cortical processing are categorical

450

Word-based

- Contextual word surprisal not unrelated to

semantics

- Long latency stages (consistent with top-down processing) in line with predictive

processing models

350 120

Lexical

Sub-Lexical **Phonemic**

Acoustic

Speech

Stimuli
100 80

60

0 0

time (ms) time (ms)

Karunathilake et al. (2025) *Neural Dynamics of the Processing of Speech Features ...* J Neurosci

Application: Is Distorted Speech Intelligible?

- Even very clear speech may be unintelligible
- More common: very distorted speech may still be intelligible
- Can neural *categorical encoding of speech features* be used to determine when the brain processes speech sounds as intelligible?

Intelligibility Experimental Design

- Manipulate intelligibility but
 - keep acoustics unchanged - Speech acoustics:
 - Vocoded speech
 - Clear speech Vocoded speech

(a)

speech clarity rating speech clarity rating Trial 1

three-band noise
vocoded speech

- Intelligibility manipulated via priming

⚠️ PRE ⚠️ CLEAN ⚠️ POST

~20 s ~20 s ~20 s Intelligibility rating (0-5)? Intelligibility rating (0-5)?

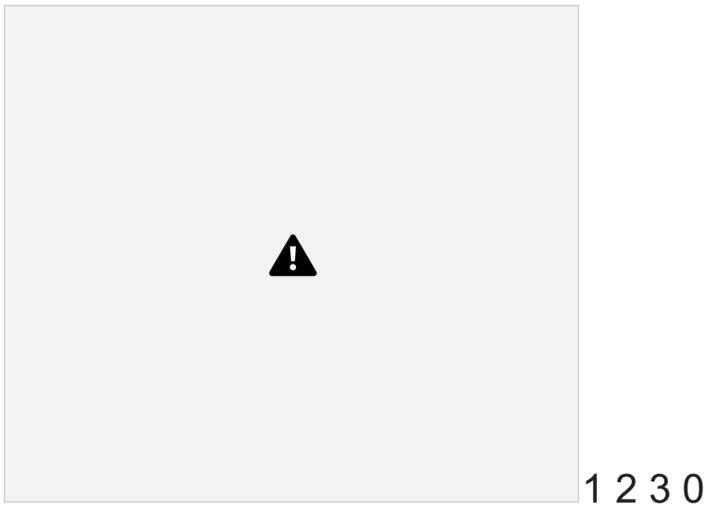
·
·
·
Trial 36

● Hypothesized intelligibility

4
measure(s)

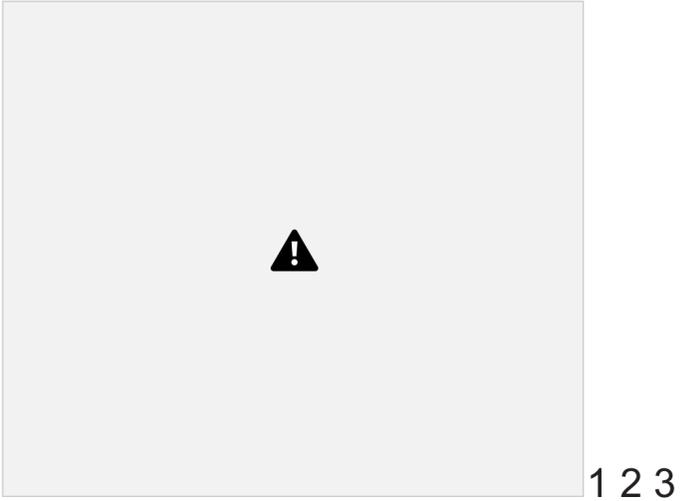
2
“Slice an apple through at its equator, and you will find five small chambers arrayed in a perfectly symmetrical

0
Vocoded speech Clear speech 0

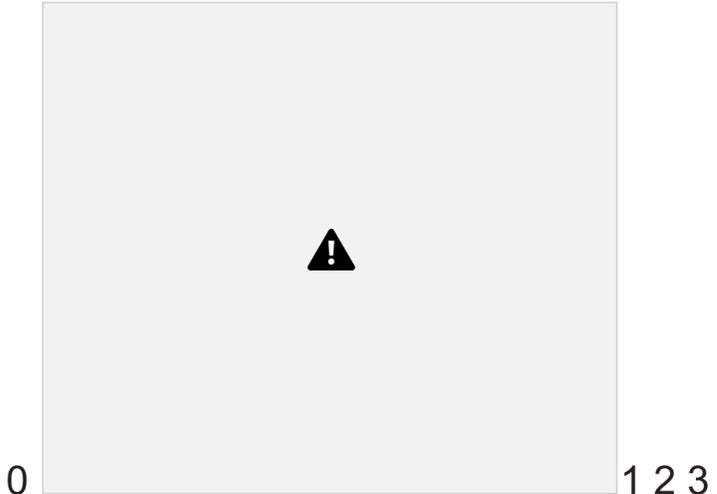


Frequency (kHz)

- word boundaries



Vocoded speech



0

-40

1 2 3

-80

-120

starburst—a pentagram.”

Time (s)

Time (s)

Time (s)

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

Intelligibility Behavioral

Results (a) ***

Speech clarity
increases from Pre
condition to Post
condition

Speech Clarity Rating

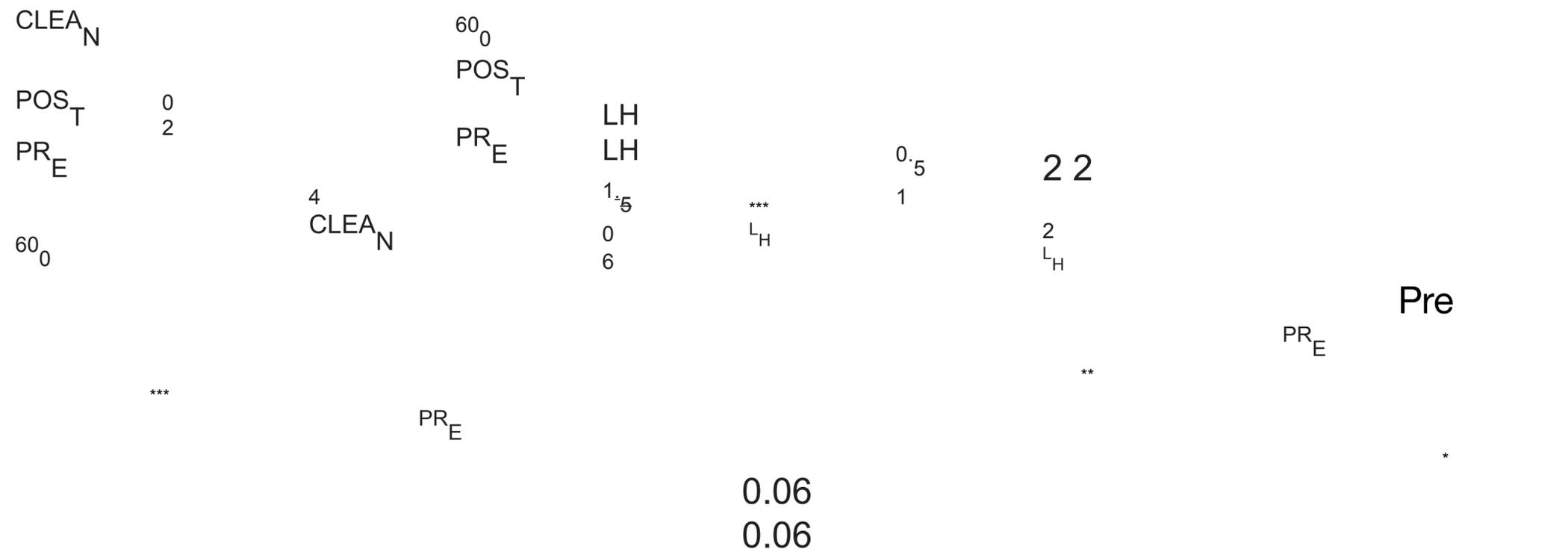
Intelligibility Rating
4

PRE POST
Pre Post Condition

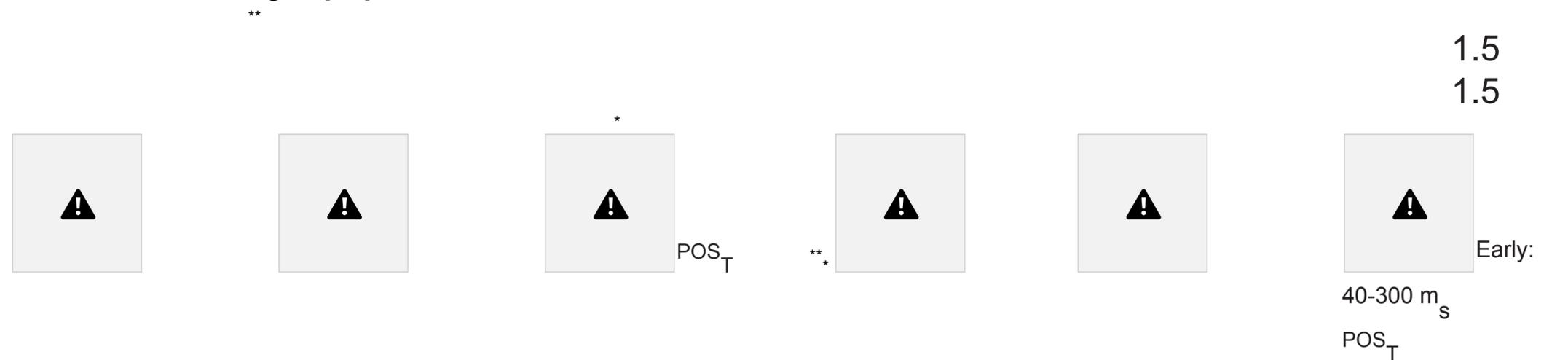
Karunathilake et al. (2023) *Neural Tracking Measures of Speech Intelligibility...*, PNAS

0
40

Intelligibility Neutral Results



- Word onset TRF shows both early (+)



Early: 40-300 m_s

and late (-) processing stages

- Only in

left hemisphere

- Late

processing shows a stage

Word onset TRF late

CLEA_N

Early: 40-300 m_s

CLEA_N

Post

Clean 1

R_H



response

e

0.0.5

*** PR_E

increases Pre → Post

R_H



POS_T CLEA_N

early 0

~~single subject~~

PR_E

0

0

0.0.5

POS_T CLEA_N

*

early

L_H

~~Early: 40-300~~

ms

0 200 400 600

0 200 400 600

PR_E

L_H

**

POS_T

**

left right

**



Late: 330-650 m_s

POS_T

**



1

**

PR_E

MEG0.5

0.5

0 0

PR_E PR_E

Late: 330-650 m_s

larger change than early

CLEA_N
2
R_H

● Acoustic responses: no



Late: 330-650 m_s

LH RH

late PR_E
R_H

~~Late: 330-650 ms~~
left right POS_T

PRE
PRE
CLEA_N
0.0₆

● Word
Surprisal
response also
increases

Pre → Post

~~Word-based
Responses:~~

● Neural

1
0

0.5

0 200 400 600

**



RH

POST
POST

0
*
LH
*

**
**

*

ive meas ure of intelligibility PRE CLEAN POST PRE CLEAN PRE CLEAN
object POST PRE CLEAN POST CLEAN POST

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

Early: 40-300 ms Late: 330-650 ms

Summary

Phil. Trans. R. Soc. B | Volume 375 | Issue 1789 | 6 January 2020

- Investigating *neural speech & language*

The Royal Society is a self-governing Fellowship of many of the world's most distinguished scientists drawn from all areas of science, engineering, and medicine. The Society's fundamental purpose, as it

processing in humans has broader impacts:

has been since its foundation in 1660, is to recognise, promote, and support excellence in science and to encourage the development and use of science for the benefit of humanity.

- clinical applications

The Society's strategic priorities emphasise its commitment to the highest quality science, to curiosity-driven research, and to the development and use of science for the benefit of society.

● animal communications

These priorities are:

- Promoting science and its benefits
- Recognising excellence in science
- Supporting outstanding science

What can animal communication teach us about human language?

● evolution of language?

- Providing scientific advice for policy
- Fostering international and global cooperation
- Education and public engagement

For further information on the Royal Society

The Royal Society
6 – 9 Carlton House Terrace
London SW1Y 5AG
T +44 20 7451 2500

● Categorical perception & **categorical neural**

W royalsociety.org

For further information on Philosophical Transactions
of the Royal Society B
T +44 20 7451 2602

processing in vocalization/speech

E philtransb@royalsociety.org

W royalsocietypublishing.org/journal/rstb

- seen for speech: phonemes, words, ... ● dissociable from acoustics
- provides new insight re: linguistics
- not available unless speech intelligible

0962-8436(20200106)375:1789

ISBN: 978-1-78252-429-8

ISSN 0962-8436

The Royal Society Registered Charity No 207043

ISSN 0962-8436 | Volume 375 | Issue 1789 | 6 January 2020

What can animal communication teach us about human language? Theme issue compiled and edited by Jonathan B. Fritz, William J. Idsardi and Gerald S. Wilkinson



These slides
available at:
ter.ps/simonpubs

thank you

<https://canl.umd.edu> Mastodon: @jzsimon@fediscience.org