# Towards Objective Measures of Speech Perception

Jonathan Z. Simon

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

University of Maryland

# Objective Measures of Speech Perception

- What do I mean by ***objective measure***?

  ▶ EEG/MEG measures of cortical activity

  ▶ Stimulus: naturalistic, long-duration speech

  ▶ Not addressed here:

    - subcortical activity
    - other non-invasive measures (fNIRS, fMRI)
    - other forms of speech

# Objective Measures of Speech Perception

- What do I mean by **speech perception**?

  ▶ Beyond intelligibility

  ▶ Allow for role of cognition

  ▶ Role of attention

  ▶ Importance of language in speech perception

  ▶ Importance of speech meaning (semantics)

  ▶ Processing effort? (not addressed here)

# Outline

- Background & motivation

  ▸ Neural responses in time

  ▸ Response prediction from a stimulus via Temporal Response Function (TRF)

  ▸ Stimulus reconstruction from responses

- Towards objective measures of

  ▸ Speech intelligibility

  ▸ Lexical processing of speech

  ▸ Semantic processing of speech

# Outline

- Background & motivation

  ▸ Neural responses in time

  ▸ Response prediction from a stimulus via Temporal Response Function (TRF)

  ▸ Stimulus reconstruction from responses

- Towards objective measures of

  ▸ Speech intelligibility

  ▸ Lexical processing of speech

  ▸ Semantic processing of speech

# EEG & MEG Responses in Time

AM at 3 Hz

3 Hz phase-locked response

Activity time-locked to temporal modulations of sounds

response spectrum

3 Hz

6 Hz

0          10

Frequency (Hz)

Ding & Simon, J Neurophysiol (2009)

# Predicting EEG/MEG <u>Responses</u>

# Predicting EEG/MEG Responses

Temporal Response Function (TRF)

Stimulus

Response

Resp. Stimulus

# Predicting EEG/MEG <u>Responses</u>

Temporal
Response
Function
(TRF)

Stimulus

Response

Resp. Stimulus

# Predicting EEG/MEG <u>Responses</u>

Stimulus

Temporal
Response
Function
(TRF)

Response

Resp.   Stimulus

# Predicting EEG/MEG Responses

Temporal
Response
Function
(TRF)

Stimulus

Response

Stimulus

Resp.

# Predicting EEG/MEG Responses

**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)

# Predicting EEG/MEG <u>Responses</u>

**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)

# Stimulus Reconstruction in Time



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

Speech envelope

Continuous MEG recording

Time [seconds]

# Stimulus Reconstruction in Time



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

Speech envelope

"Decoder"

Continuous MEG recording

1    2    3    4    5

Time [seconds]

# Stimulus Reconstruction in Time

# Stimulus Reconstruction in Time

# Stimulus Reconstruction in Time



**Stimulus Reconstruction:**

Stimulus and response are known;

Find the best matrix in time **_and space_** to produce the **stimulus** from the **response**

# Stimulus Reconstruction in Time



**Stimulus Reconstruction:**

Stimulus and response are known;

Find the best matrix in time **and space** to produce the **stimulus** from the **response**

# Stimulus Reconstruction in Time



**Stimulus Reconstruction:**

Stimulus and response are known;

Find the best matrix in time **and space** to produce the **stimulus** from the **response**

# Stimulus Reconstruction in Time



**Stimulus Reconstruction:**

Stimulus and response are known;

Find the best matrix in time **and space** to produce the **stimulus** from the **response**

# Stimulus Reconstruction in Time



**Stimulus Reconstruction:**

Stimulus and response are known;

Find the best matrix in time *and space* to produce the **stimulus** from the **response**

# Cortical Representations of Continuous Speech

- For long duration continuous speech

- Encoding & decoding (complementary)

- Linear model

- Acoustics: spectrotemporal **envelope**

- Envelope rates: ~ 1 - 10 Hz

# Outline

- Background & motivation
  - ▸ Neural responses in time
  - ▸ Response prediction from a stimulus via Temporal Response Function (TRF)
  - ▸ Stimulus reconstruction from responses
- Towards objective measures of
  - ▸ Speech intelligibility
  - ▸ Lexical processing of speech
  - ▸ Semantic processing of speech

# Outline

- Background & motivation
  - ▸ Neural responses in time
  - ▸ Response prediction from a stimulus via Temporal Response Function (TRF)
  - ▸ Stimulus reconstruction from responses

- **Towards objective measures of**
  - ▸ **Speech intelligibility**
  - ▸ Lexical processing of speech
  - ▸ Semantic processing of speech

# Stimulus Reconstruction
# ➜ Intelligibility



Neural quality of time-locked response to speech acoustics

Behavioral ability to correctly report heard words

Vanthornhout et al., JARO (2018)

# Stimulus Reconstruction ➜ Intelligibility



Vanthornhout et al., JARO (2018)

# Stimulus Reconstruction
# ➡ Intelligibility



Objective CT vs. Behavioral SRT

r = 0.69, p = 0.001

Correlation Threshold (dB)

Speech Reception Threshold (dB)

Vanthornhout et al., JARO (2018)

# Stimulus Reconstruction ➡ Intelligibility



**But *what* is being measured neurally?**

Cortical responses, but where? (or when?)

Objective CT vs. Behavioral SRT

Correlation Threshold (dB)

+4

0

−4

−8

−12

−10    −8    −6

Speech Reception Threshold (dB)

r = 0.69, p = 0.001

Vanthornhout et al., JARO (2018)

# Stimulus Reconstruction in Time



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

Speech envelope

"Decoder"

Continuous MEG recording

1      2      3      4      5

Time [seconds]

**But *what* is being measured neurally?**

Cortical responses, but where? (or when?)

# Stimulus Reconstruction
# ➜ Intelligibility

**Integration window span indicates latencies of interest**

- choose window for reconstruction
- not based on highest correlation (of reconstructed stimulus)
- based on reconstruction **monotonicity** as a function of SNR.

Vanthornhout et al., JARO (2018)

# Stimulus Reconstruction
# ➜ Intelligibility

**Integration window span indicates latencies of interest**

- choose window for reconstruction
- not based on highest correlation (of reconstructed stimulus)
- based on reconstruction **monotonicity** as a function of SNR.

### Reconstruction Monotonicity by SNR



**Integration window choice: 0 ms to 75 ms**

**early auditory cortex**

**pre-attentive**

Vanthornhout et al., JARO (2018)

# Stimulus Reconstruction ➜ Intelligibility

- Continuous speech envelope reconstruction (neurometric) threshold predicts behavioral speech reception threshold (SRT).

- Uses long duration continuous speech

- Based on robust *acoustic* speech representation

- Early auditory cortex most critical (pre-attentive)

# Stimulus Reconstruction ➜ Intelligibility

- UPDATES from the Francart Lab

  ▸ Response prediction (~~stimulus reconstruction~~)

  ▸ Theta band

  ▸ Speech Envelope ➜ Spectrogram

  ▸ Added new representation: phonetic features*

*Role of phonetic features vs. spectrogram onsets?

Lesenfants et al., Hear Res (2019)

# Phonetic Features vs. Spectrogram Onsets

- + 'phonetic features' representation increases EEG response prediction: Di Liberto et al. (2015).

- Adding only *acoustic* spectrogram *onsets* gives same predictive benefits as phonetic features for MEG responses: Daube et al. (2019).

- Also seen in Simon lab: Brodbeck et al. (2018).

➡ Phonetic features too correlated with acoustic onsets, in natural speech, to isolate them

# Stimulus Reconstruction ➜ Intelligibility

- UPDATES from the Francart Lab

  ▶ Age really matters: Decruy et al. (2019)



Not just linear but quadratic uptick

Cognitive decline also matters

In agreement with Presacco et al. (2016a, 2016b).

# Outline

- Background & motivation

  ▸ Neural responses in time

  ▸ Response prediction from a stimulus via Temporal Response Function (TRF)

  ▸ Stimulus reconstruction from responses

- Towards objective measures of

  ▸ Speech intelligibility

  ▸ Lexical processing of speech

  ▸ Semantic processing of speech

# Outline

- Background & motivation
  - ▸ Neural responses in time
  - ▸ Response prediction from a stimulus via Temporal Response Function (TRF)
  - ▸ Stimulus reconstruction from responses
- **Towards objective measures of**
  - ▸ Speech intelligibility
  - ▸ **Lexical processing of speech**
  - ▸ Semantic processing of speech

# Lexical Processing

- Processing by early auditory cortex critical

- Using more than global speech envelope helps

- Another level of speech perception:

  ▸ Transforming speech sounds into words

  ▸ "Lexical processing"

    - Language-based but not via word meaning

Brodbeck et al., Curr Biol (2018)

# Acoustic to Lexical Speech Processing



his          noble          mind          forgot          the   cakes

**Acoustic Envelope (8 bands)**

Acoustic envelope (8 bands)

**Acoustic Onset (8 bands)**

# Acoustic to Lexical Speech Processing

# Acoustic to Lexical Speech Processing



his     noble     mind     forgot     the  cakes

h ɪ z  n oʊ  b əl  m aɪ  n d f  ɝ g ɑ t  ð i  k eɪ  k s

**Phoneme Onset**

**Phoneme Surprisal**

$$surprisal_i = -\log_2\left( \sum_{word \in cohort_i} freq_{word}(i) \middle/ \sum_{word \in cohort_{i-1}} freq_{word}(i-1) \right)$$

**Cohort Entropy**

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

**SUBTLEX:**
**51 million words**
movie subtitle database

# Surprisal

Number of times a word that starts with this sequence occurs in SUBTLEX

K EY …
52908
(90 words)

Number of words that start with this sequence

Surprisal

K EY M …
23875 (45%)
(4 words)

K EY S …
16048 (30%)
(13 words)

K EY K …
2598 (5%)
(3 words)

K EY N …
1337 (3%)
(13 words)
⋮

"came", "Cambridge", …

"case", "cases", "caseworker",
"casein", …

"cake", "caked", "cakes"

"cane", "canine", "Canaan",
"Kane", "Keynesian", …

# Entropy

## Cohort entropy

‣ How unpredictable is the current word?



L EY K …

lake (95%)    lakes (5%)

K EY K …

cake (88%)    cakes (11%)    caked (1%)

B EY K …

baker (29%)    bacon (25%)    baked (14%)    bake (14%)

Entropy

# Word onsets

## Do we…

‣ Anticipate word boundaries based on context?

‣ Infer them later based on consistency?



**"The catalogue in a library"**

(Norris & McQueen, 2008)

# Acoustic Results



Acoustic
Envelope

# Acoustic Results

Acoustic
Envelope



Brodbeck et al., Curr Biol (2018)

# Acoustic Results



Brodbeck et al., Curr Biol (2018)

# Acoustic Results



Acoustic Envelope

Acoustic Onset

$5.8 \times 10^{-03}$

$\Delta z$

$0$

$1.1 \times 10^{-02}$

$0$

**

$1$

$0$

$-1$

30 ms

110 ms

70 ms

140 ms

cf. Daube et al., Curr Biol (2019)

- Onset explains more variance
- Latency(ies) as expected
- Strongly bilateral
- Onset stronger in right hemisphere

Brodbeck et al., Curr Biol (2018)

# Neural Lexical Processing

Phoneme
Surprisal

Cohort
Entropy

$1.1 \times 10^{-03}$

$0$

$1.3 \times 10^{-03}$

$0$

**

Brodbeck et al., Curr Biol (2018)

# Neural Lexical Processing

Phoneme Surprisal

Cohort Entropy



Brodbeck et al., Curr Biol (2018)

# Neural Lexical Processing



Brodbeck et al., Curr Biol (2018)

# Neural Lexical Processing

- Rapid transformation to lexical
- Word boundaries identified
- Surprisal = local measure of phoneme prediction error (predictive coding?)
- Cohort entropy = global measure of lexical competition across cohort
- Strongly left hemisphere dominant



Brodbeck et al., Curr Biol (2018)

# Listening at the Cocktail Party

# Acoustic Attention

2 competing speakers, equal loudness, attend to one

# Acoustic Attention

2 competing speakers, equal loudness, attend to one



- Onset Representation Dominates
- Attended Dominates Later

Brodbeck et al., Curr Biol (2018)

# Lexical Attention



Brodbeck et al., Curr Biol (2018)

# Lexical Attention



Attended lexical model

Unattended lexical model

Phoneme Onset

Word Onset

110 ms

Phoneme Surprisal

Cohort Entropy

130 ms

0    100    200    300    400    500

0    100    200    300    400    500

Time [ms]

- Only attended speech processed lexically
- Lexical processing slowed by ~15 ms

Brodbeck et al., Curr Biol (2018)

# Lexical Processing

- Speech perception at level of transforming speech sounds into words

- "Post-acoustic" phoneme processing

- Word-based

- Attention required (?)

- Surprisingly early

Brodbeck et al., Curr Biol (2018)

# Outline

- Background & motivation

  ‣ Neural responses in time

  ‣ Response prediction from a stimulus via Temporal Response Function (TRF)

  ‣ Stimulus reconstruction from responses

- Towards objective measures of

  ‣ Speech intelligibility

  ‣ Lexical processing of speech

  ‣ Semantic processing of speech

# Outline

- Background & motivation
  - ▸ Neural responses in time
  - ▸ Response prediction from a stimulus via Temporal Response Function (TRF)
  - ▸ Stimulus reconstruction from responses

- **Towards objective measures of**
  - ▸ Speech intelligibility
  - ▸ Lexical processing of speech

  - ▸ **Semantic processing of speech**

# Semantic Processing

- Speech perception includes perceiving the meaning of the speech

- Computational language models give several semantic measures: *semantic dissimilarity*

- Analysis of *Semantic-dissimilarity*-based TRF

  ▶ potential basis of objective measure of perception of speech meaning

Broderick et al., Curr Biol (2018)

# Semantic Processing

- Speech perception includes perceiving the meaning of the speech

- Computational language models give several semantic measures: *semantic dissimilarity*

- Analysis of *Semantic-dissimilarity*-based TRF

  ▸ potential basis of objective measure of perception of speech meaning

Broderick et al., Curr Biol (2018)

# Semantic Processing

## Semantic-dissimilarity TRF



Broderick et al., Curr Biol (2018)

# Semantic Processing

## Semantic-dissimilarity TRF



Time (ms)

Broderick et al., Curr Biol (2018)

# Semantic Processing

## Semantic-dissimilarity TRF



Attended Speech

Time (ms)

Broderick et al., Curr Biol (2018)

# Semantic Processing



Semantic-dissimilarity TRF

Unattended Speech

Attended Speech

Time (ms)

# Semantic Processing

Semantic-dissimilarity TRF



Unattended Speech

Attended Speech

Time (ms)

Broderick et al., Curr Biol (2018)

- This TRF reflects processing of semantics
- This semantic processing depends on attention

# Summary

- Speech perception takes many forms
- Cortical processing of speech takes many forms
- Many potential ways to link the two
  - Faithful representation of speech acoustics
  - Processing speech sounds into words (lexical)
  - Semantic level processing
  - Cognitive aspects of perception allowed
- Cortical (temporal) processing of continuous speech processing: both encoding & decoding

# Thank You

# Acknowledgements

**Current Lab Members & Affiliates**

*Christian Brodbeck*
*Alex Presacco*
Proloy Das
Jason Dunlap
Theo Dutcher
*Alex Jiao*
Dushyanthi Karunathilake
Joshua Kulasingham
Natalia Lapinskaya
Sina Miran
David Nahmias
Peng Zan

**Past Lab Members & Affiliates**

Nayef Ahmar
Sahar Akram
Murat Aytekin
Francisco Cervantes Constantino
Maria Chait
Marisel Villafane Delgado
Kim Drnec
*Nai Ding*
Victor Grau-Serrat

Julian Jenkins
Pirazh Khorramshahi
Huan Luo
Mahshid Najafi
Krishna Puvvada
*Jonas Vanthornhout*
Ben Walsh
Yadong Wang
Juanjuan Xiang
Jiachen Zhuo

**Collaborators**

Pamela Abshire
Samira Anderson
Behtash Babadi
Catherine Carr
Monita Chatterjee
Alain de Cheveigné
Stephen David
Didier Depireux

Mounya Elhilali
*Tom Francart*
Jonathan Fritz
Michael Fu
Stefanie Kuchinsky
Steven Marcus
Cindy Moss
David Poeppel
Shihab Shamma

**Past Undergraduate Students**

Nicholas Asendorf
Ross Baehr
Anurupa Bhonsale
Sonja Bohr
Elizabeth Camenga
Katya Dombrowski
Kevin Hogan
Andrea Shome
James Williams