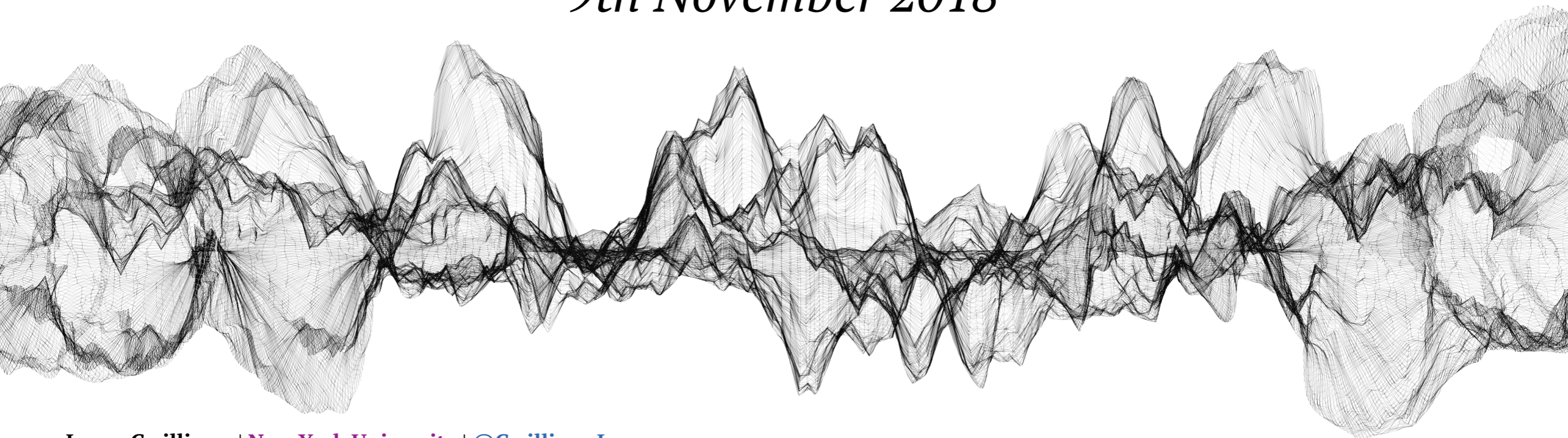




NEW YORK UNIVERSITY

Towards a mechanistic account of speech comprehension

Laura Gwilliams
9th November 2018



Levels of analysis

1. Phonemes within words

- Responses to phoneme ambiguity, phonetic features and acoustic properties (**bottom-up**)
- Neural signatures of ambiguity resolution, when provided with lexical information (**top-down**)

2. Words within sentences

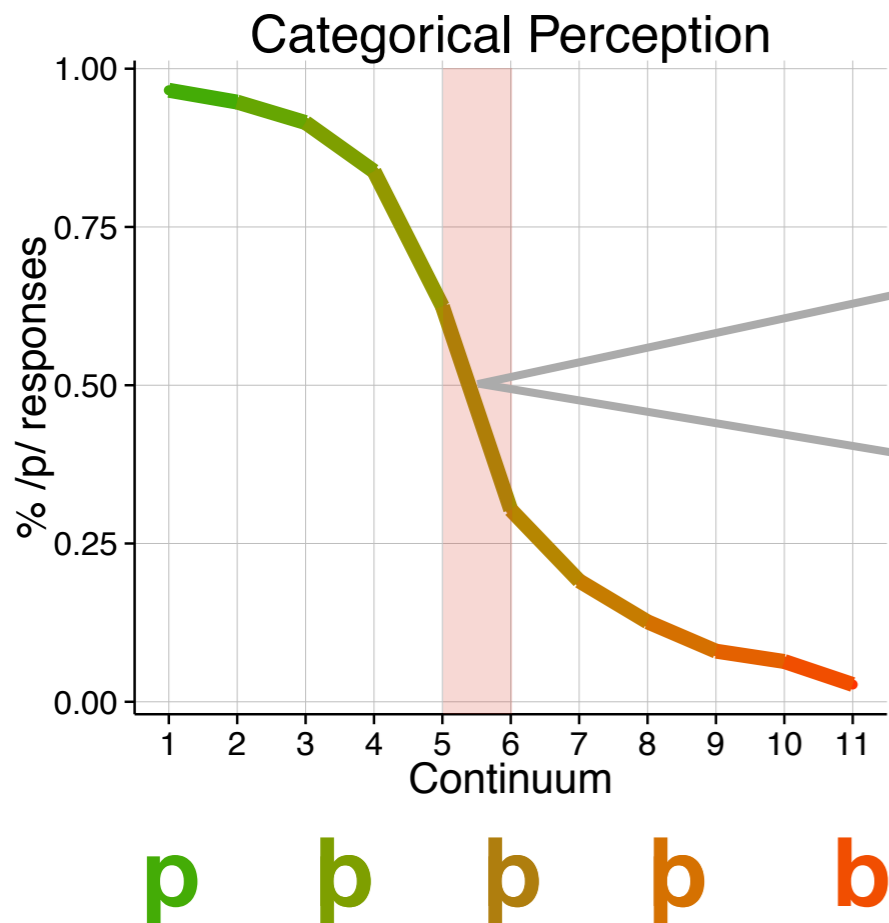
- Which linguistic properties encoded in brain activity?
- What are the relative time-courses of processing each property?
- What is the computational architecture?

Future Influences on Perception

- Speech is an inherently **noisy and ambiguous** signal
- To fluently derive meaning, listeners must **integrate top-down** contextual information to guide their interpretation
- Top-down input occurring *after* an acoustic signal can be integrated to **affect the perception of earlier sounds**
(Bicknell et al., submitted; Connine et al., 1991; Samuel, 1981; Szostak & Pitt, 2013; Warren & Sherman, 1974)

Future Influences on Perception

(this is a parakeet)



p

b

ee t

ai d

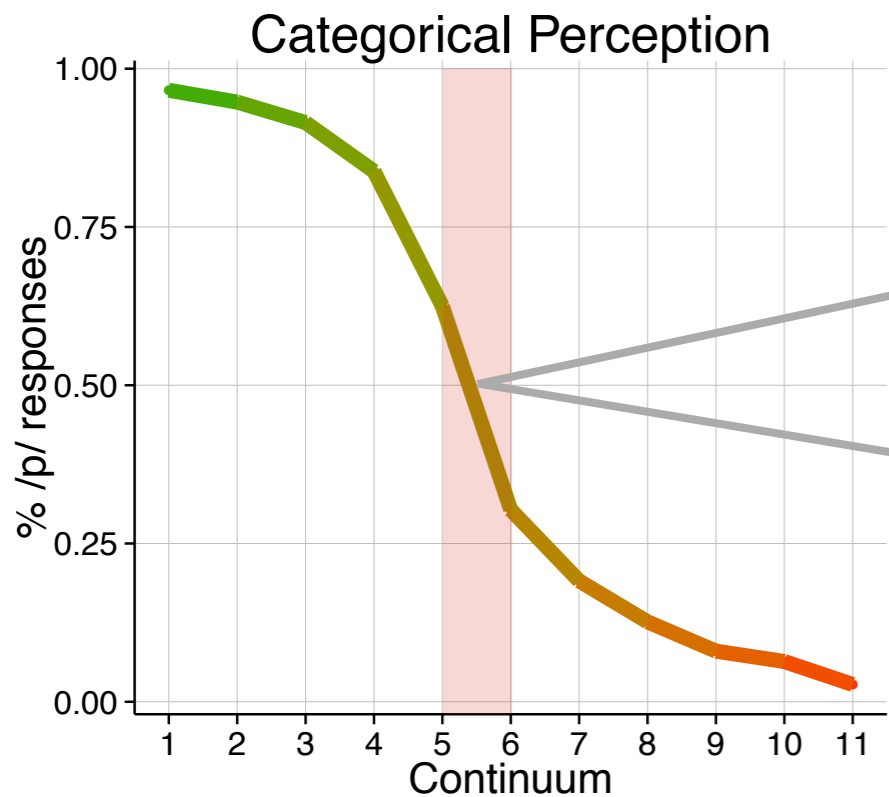


(this is a barricade)

“**P**oint of **D**isambiguation” (POD)

Future Influences on Perception

(this is a parakeet)



p a r a k e e t
b a r a k a i d



(this is a barricade)

“**P**oint of **D**isambiguation” (POD)

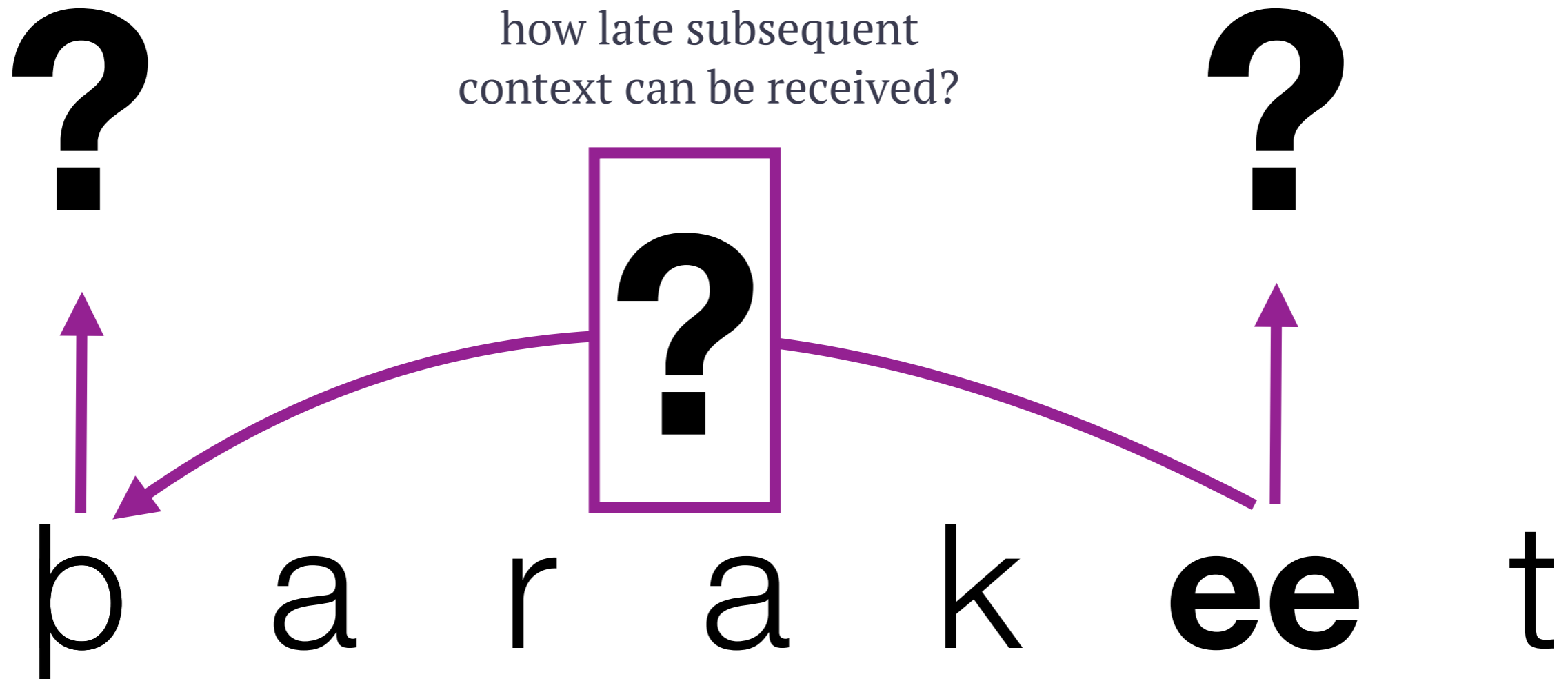
Today's Questions

i could focus on the ambiguity resolution part here, rather than the original response to ambiguity. then, tie in the ambiguity response part later, linking it with AI?

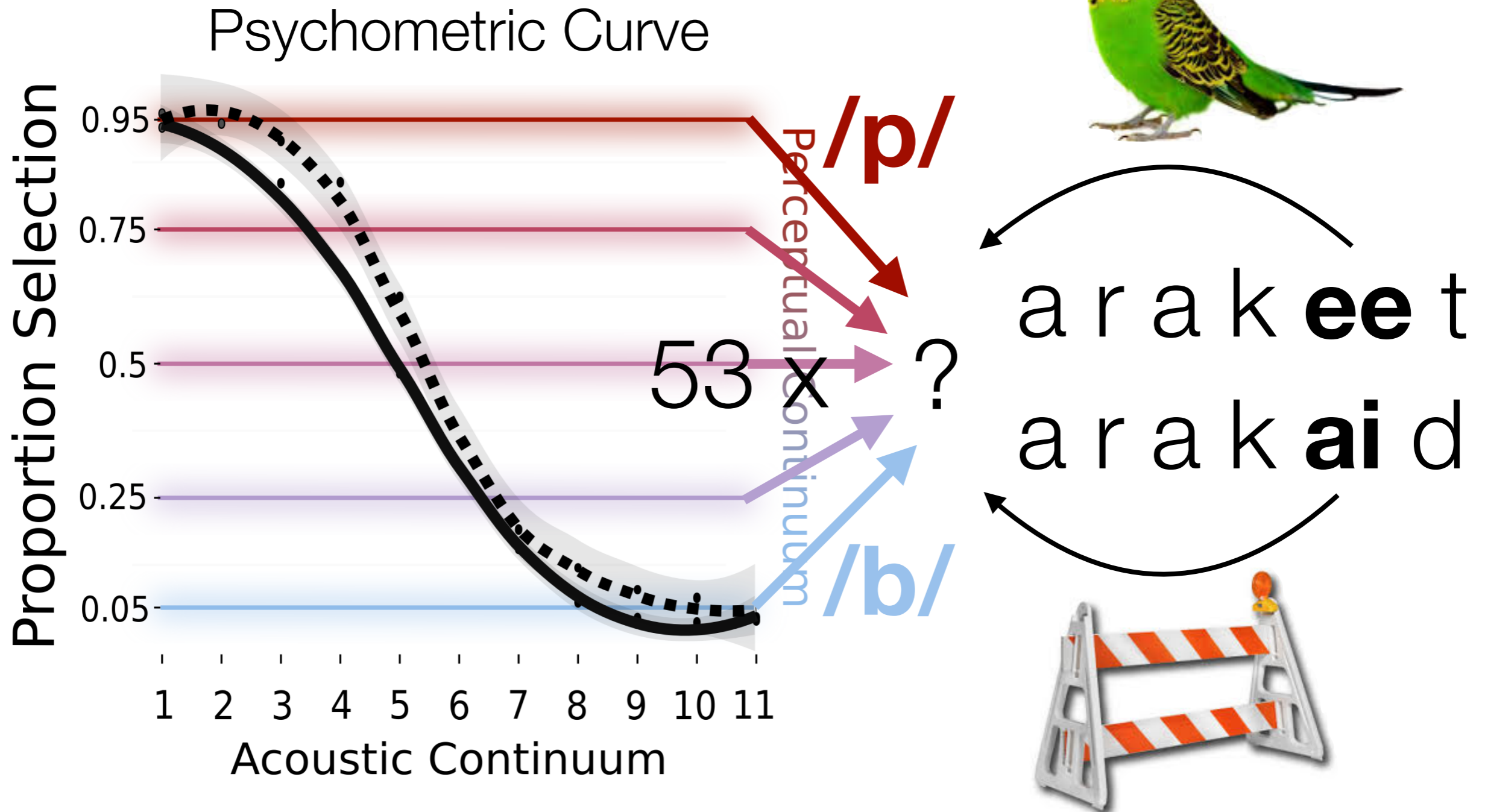
How does the auditory cortex **respond** to phonological ambiguity?

What are the neural signatures of ambiguity **resolution**?

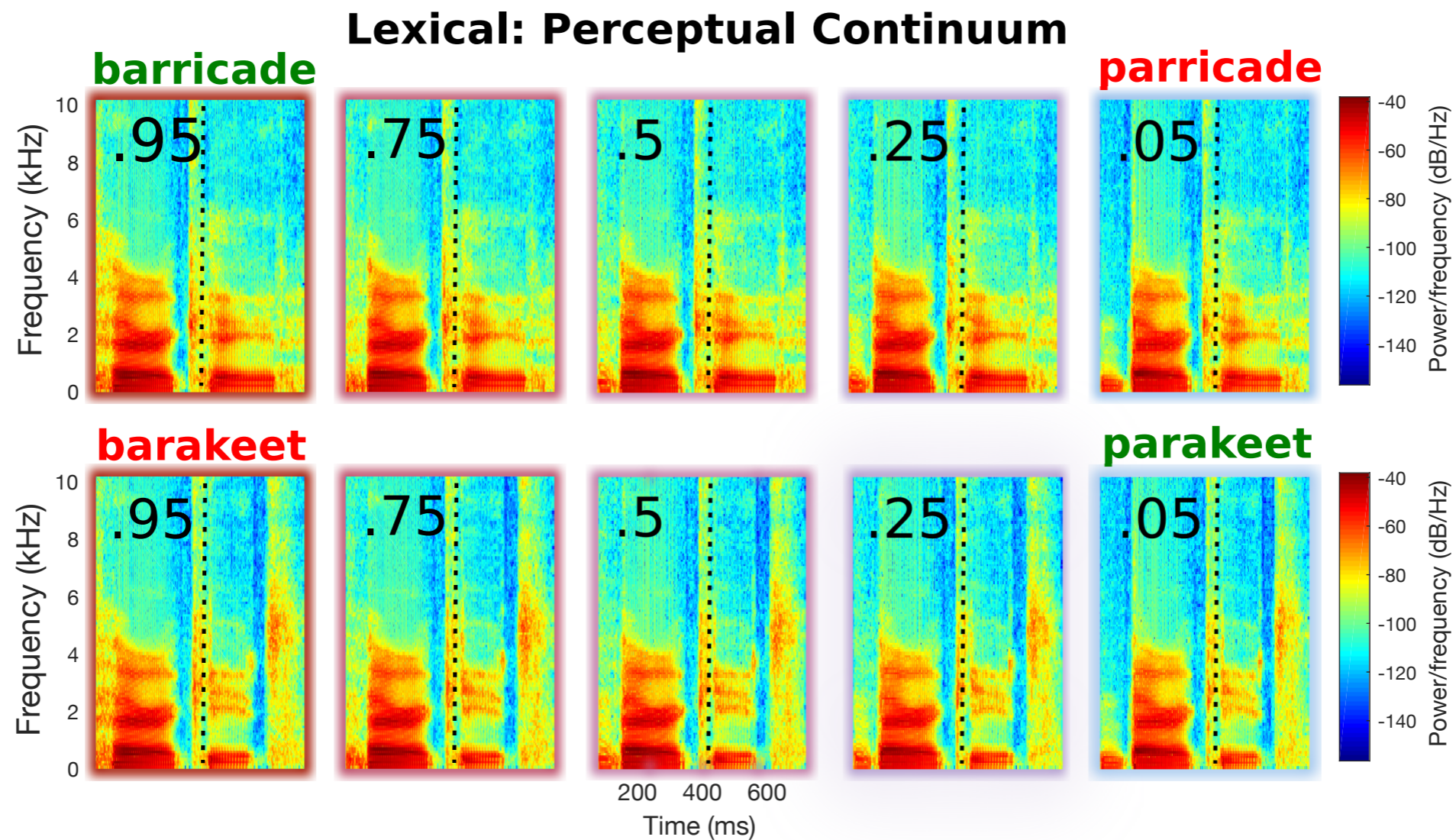
What is the **time-limit** on how late subsequent context can be received?



Design & Materials

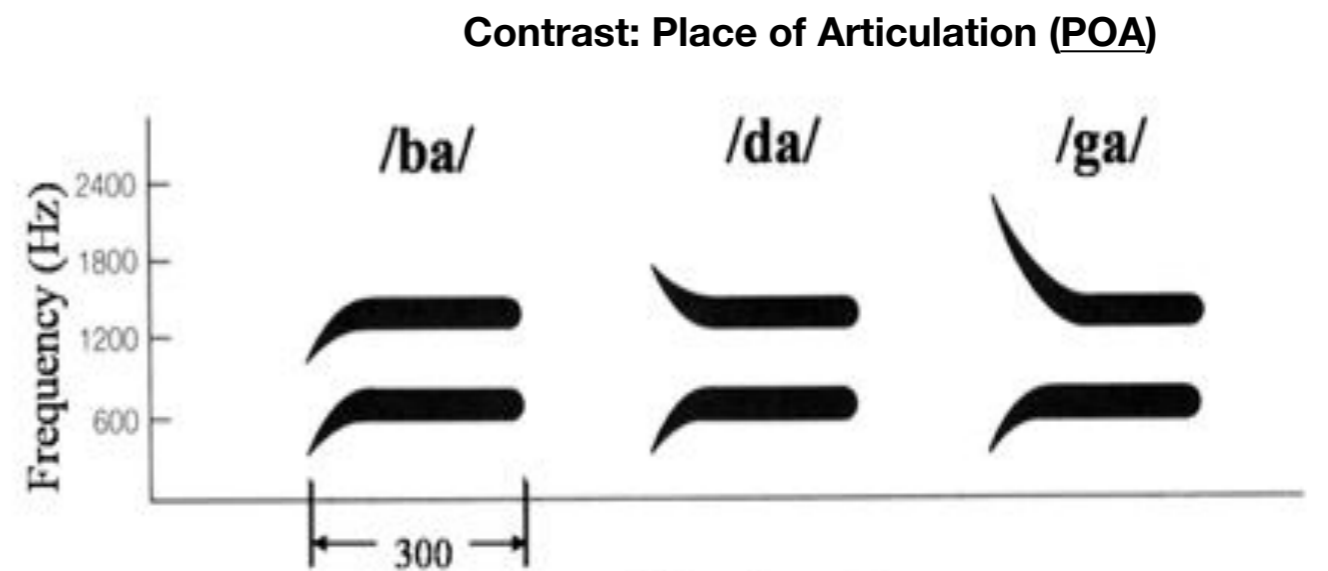
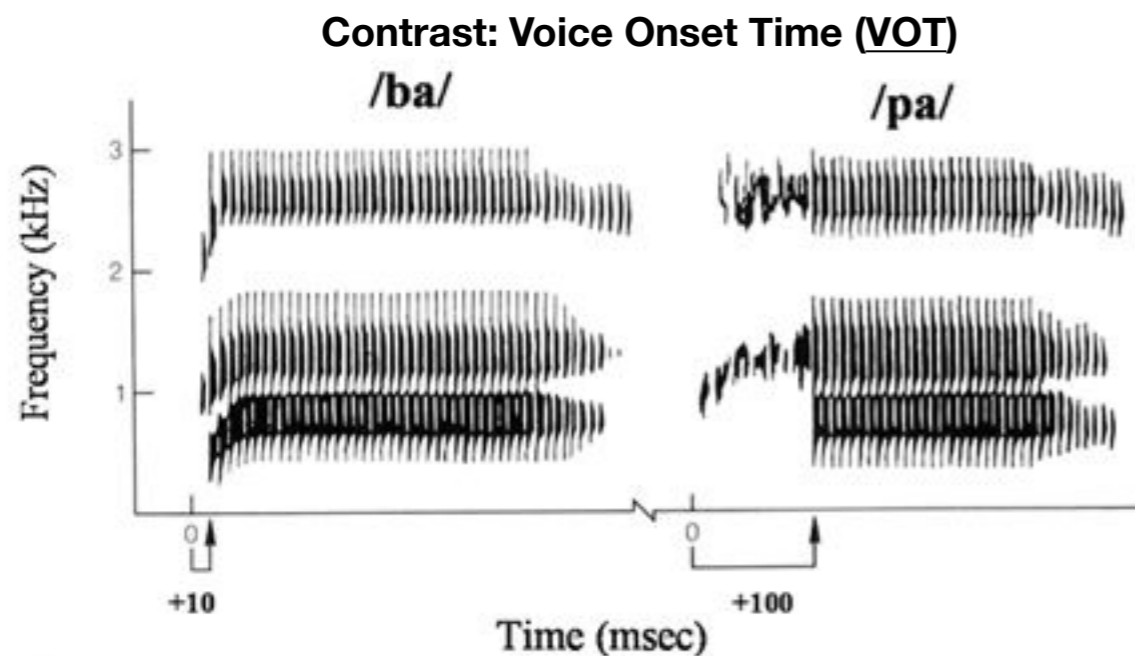


Design & Materials



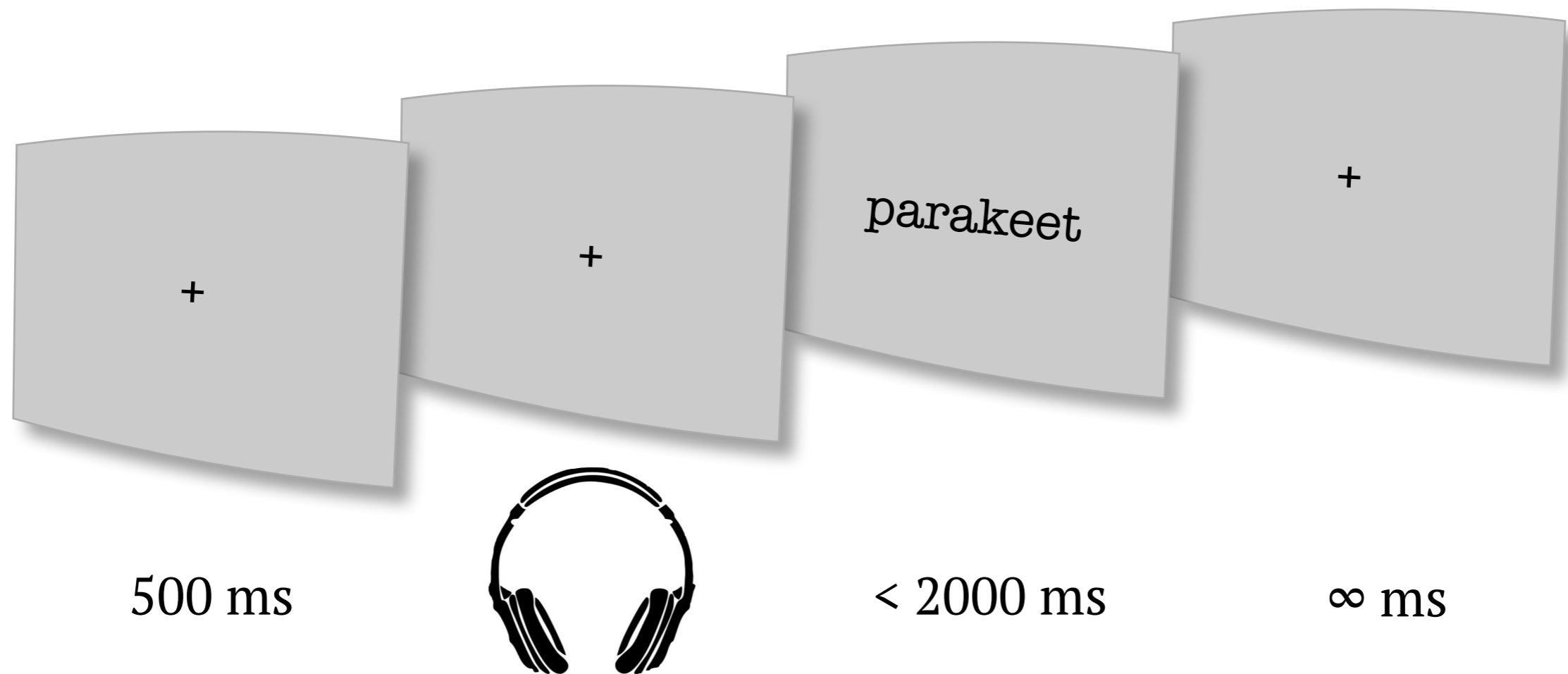
- Point of Disambiguation (POD) ranged 3-8 phonemes / 150-750 ms
- VOT (31 pairs) {p-b, t-d, k-g} and POA (22 pairs) {t-k, p-t}

Design & Materials

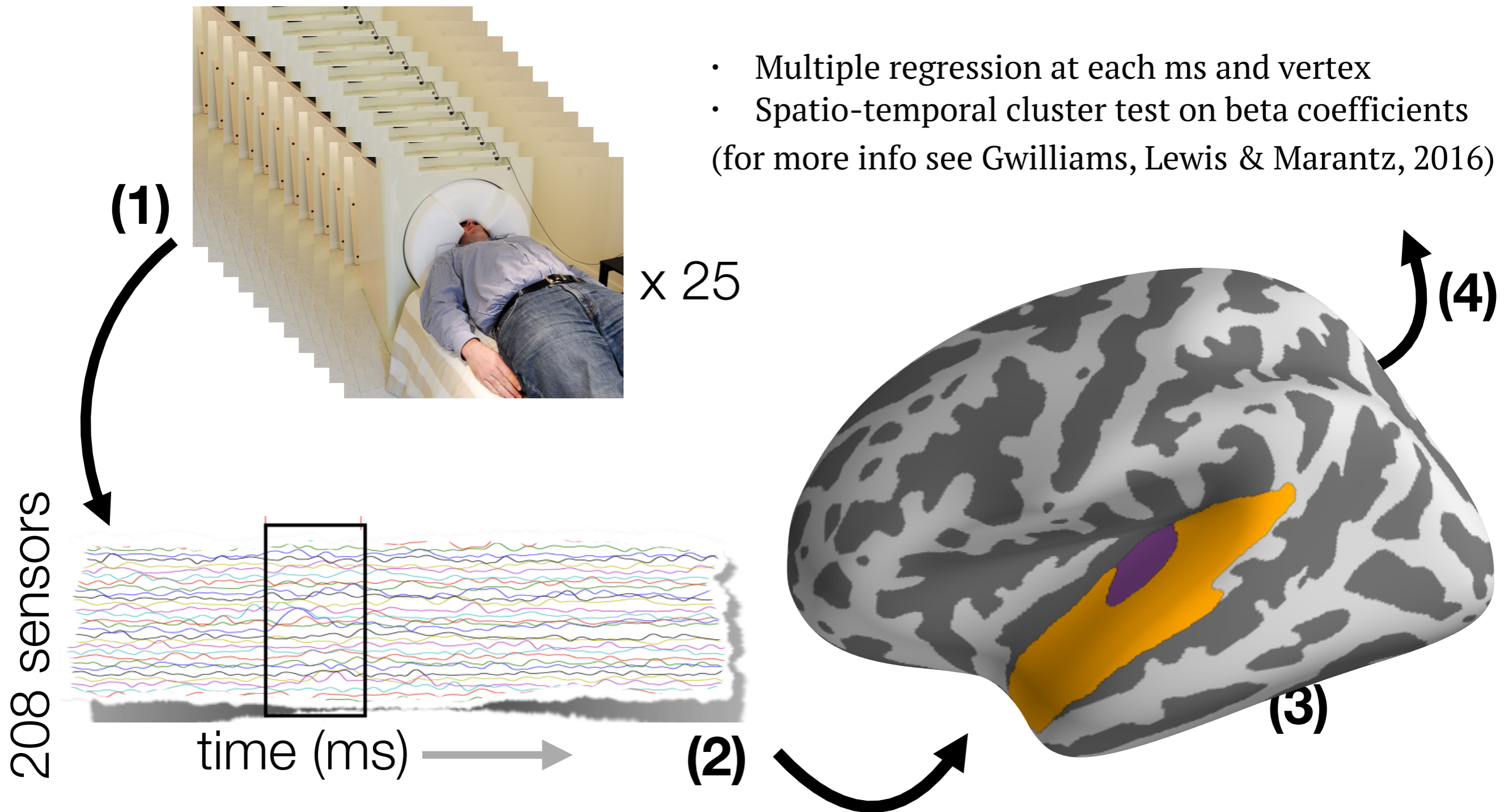


- Point of Disambiguation (POD) ranged 3-8 phonemes / 150-750 ms
- VOT (31 pairs) {p-b, t-d, k-g} and POA (22 pairs) {t-k, p-t}

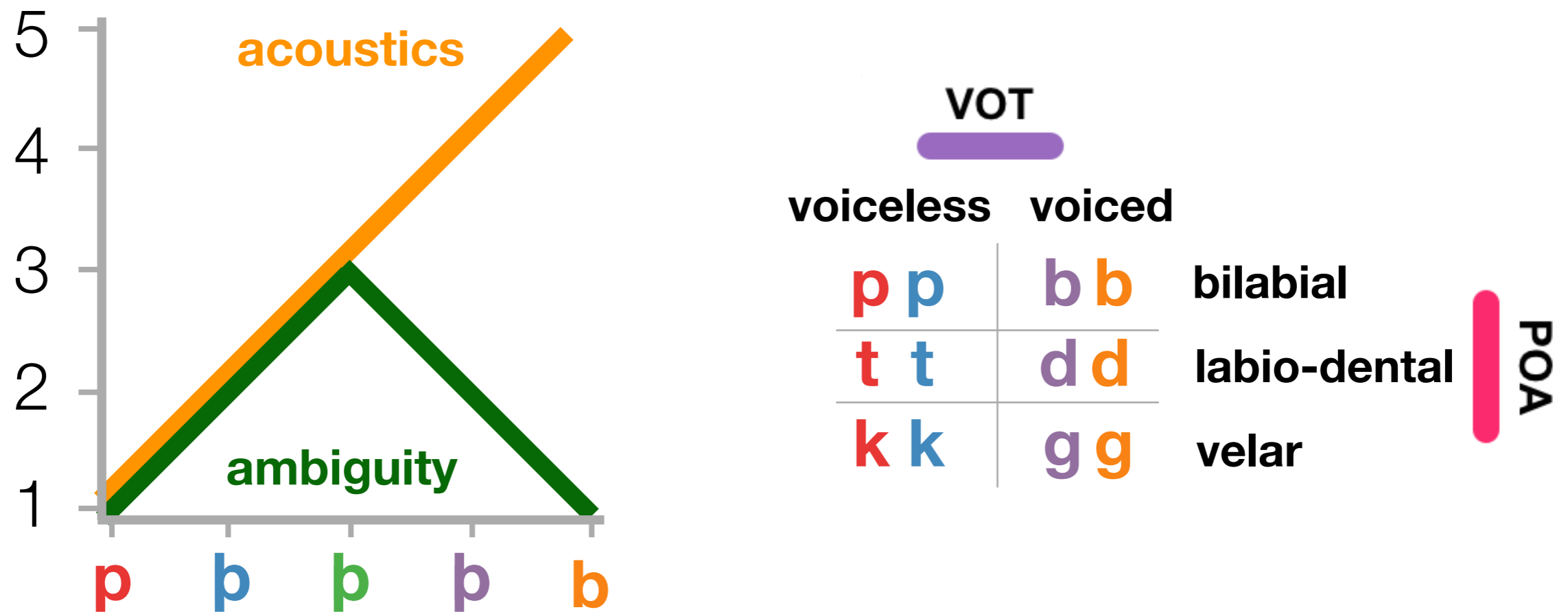
Design & Materials



Procedure & Analysis



Four Experimental Variables



Today's Questions

How does the auditory cortex **respond** to phonological ambiguity?

?



p

a

r

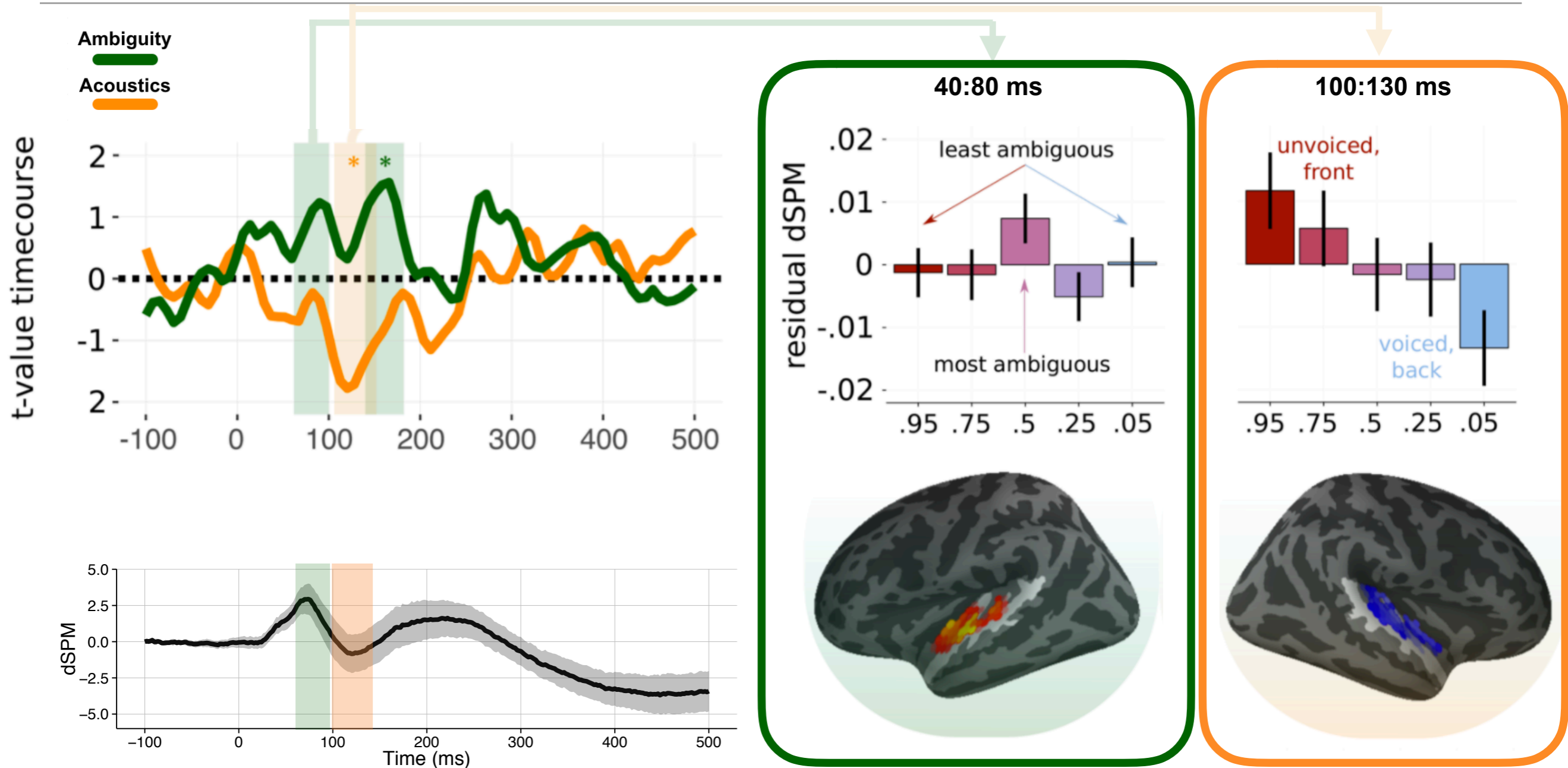
a

k

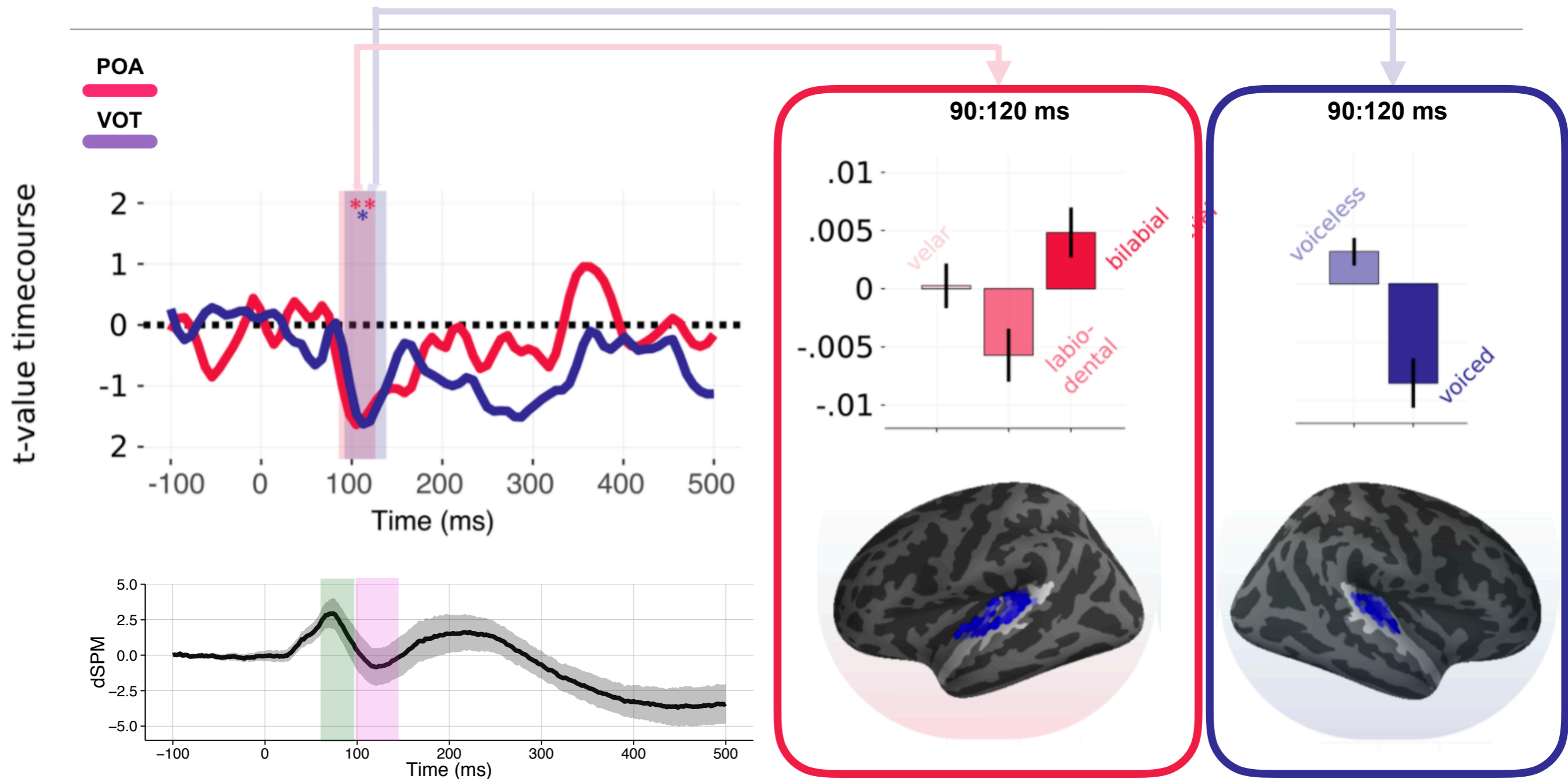
ee

t

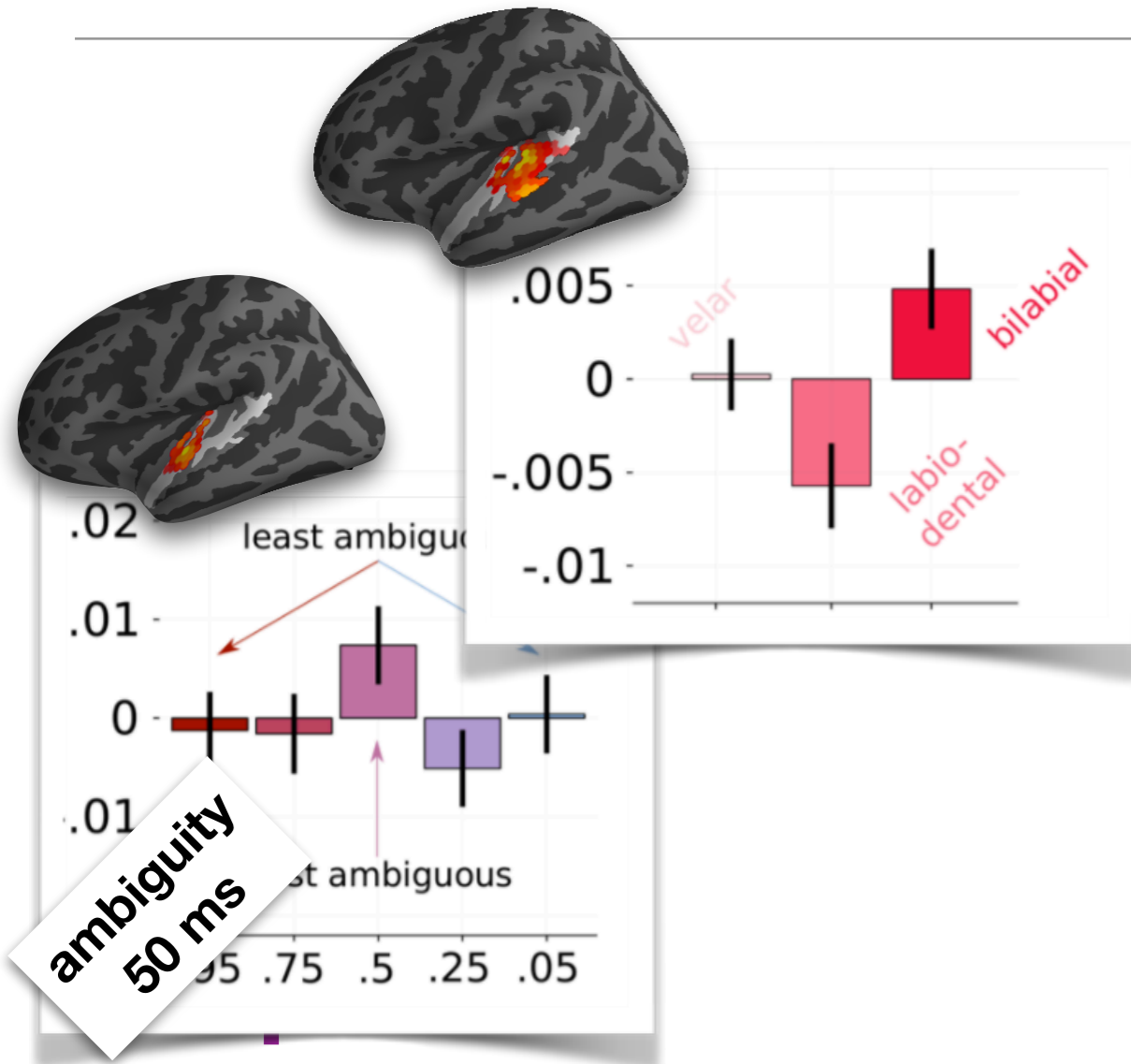
Subphonetics at Onset



Phonetic Features at Onset



Interim Conclusion

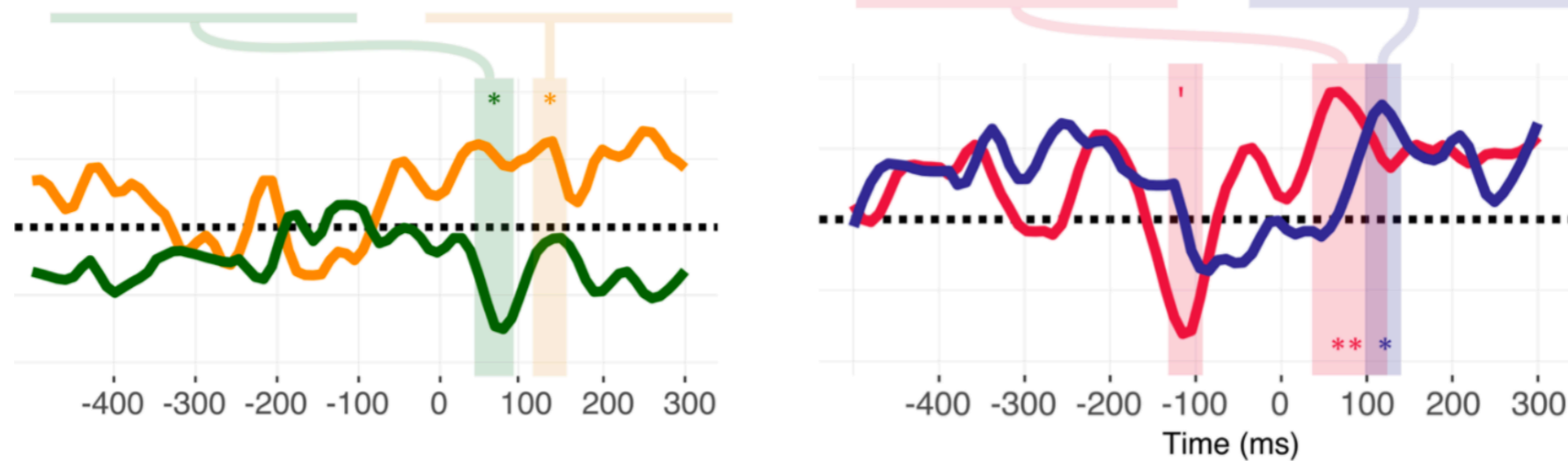


?

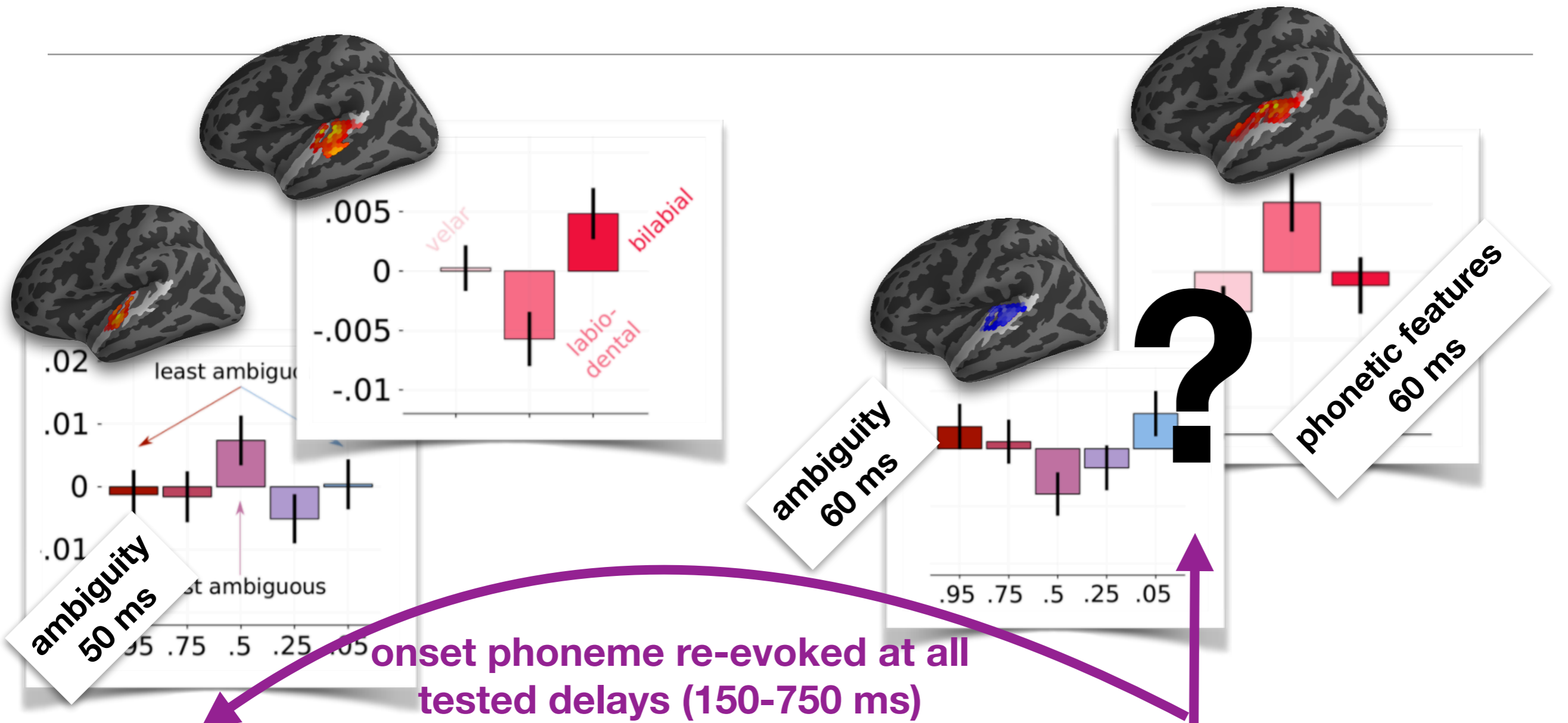


p a r a k e e t

Ambiguity at POD

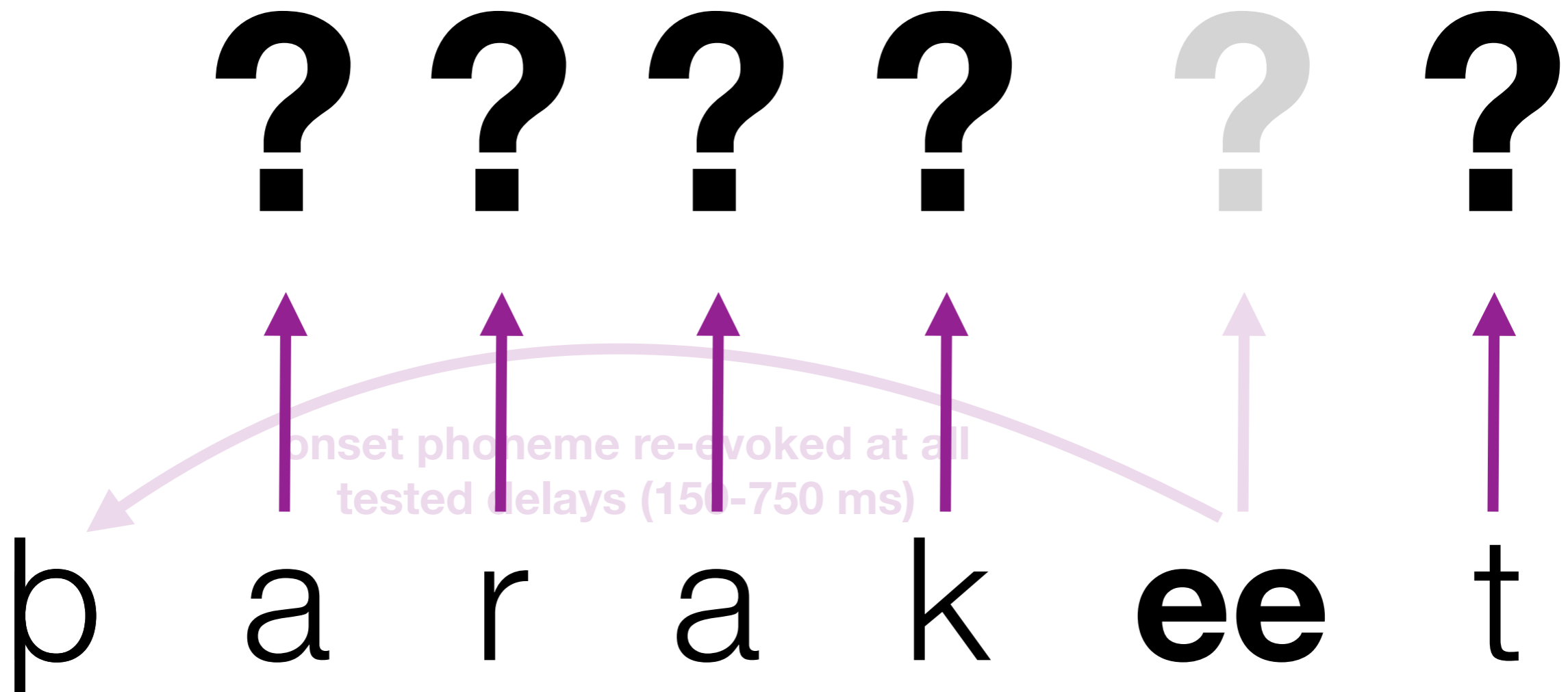


Interim Conclusion

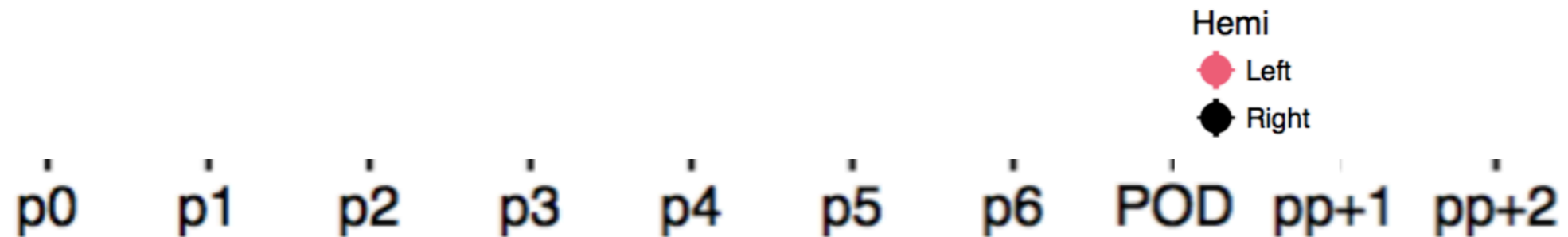


b a r a k e e t

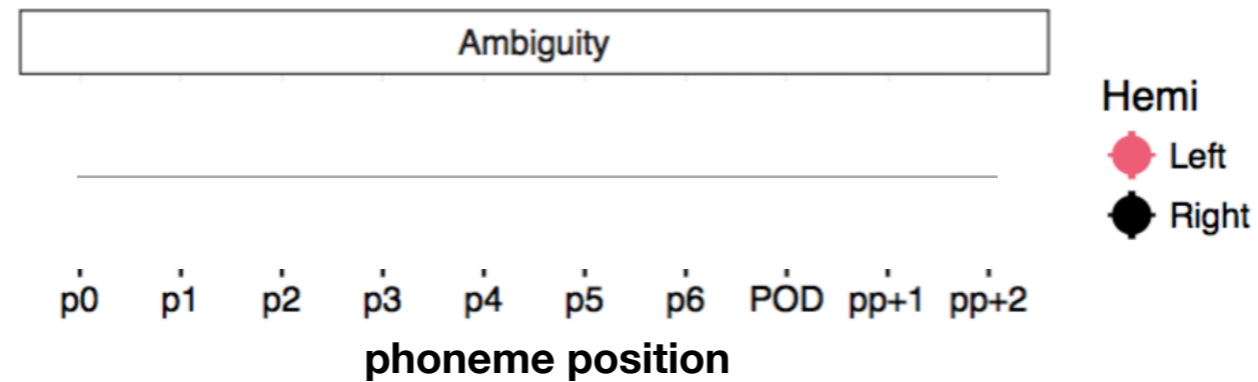
Interim Conclusion



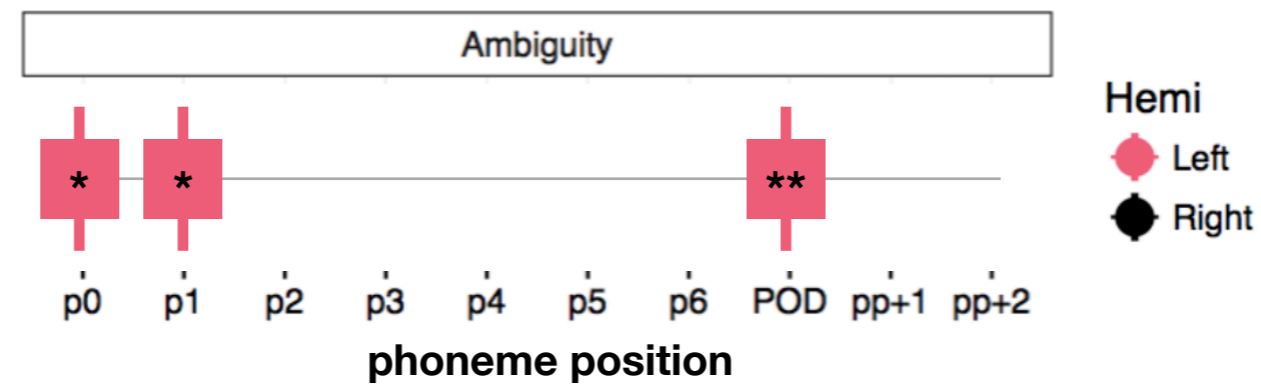
Reactivation in Intermediate Positions



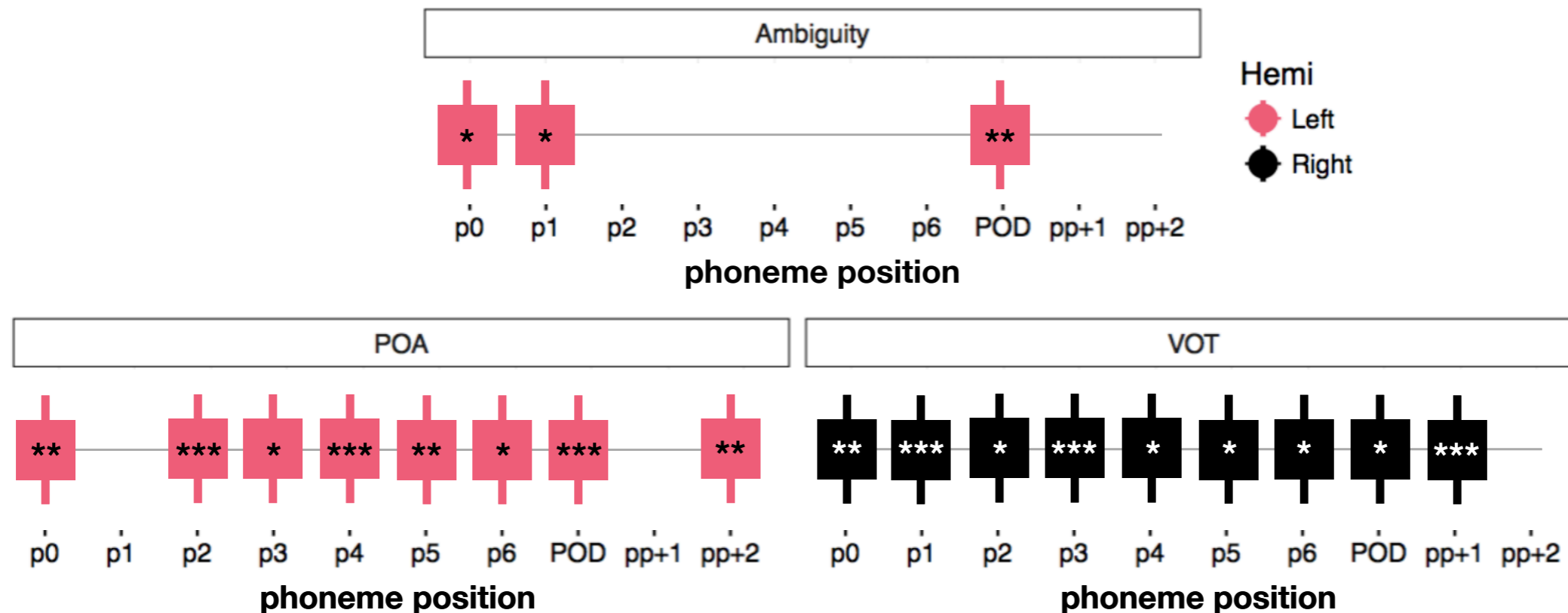
Reactivation in Intermediate Positions



Reactivation in Intermediate Positions



Reactivation in Intermediate Positions



ballet

prove

pin

bath

pacify

bond

palate

book

beef

p

pants

balance

bind

paddle

boast

poke

panda

ballet

prove

pin

bath

pacify

bond

palate

beef

book

b

b

p



pants

balance

bind

paddle

boast

poke

panda

ballet

prove

pin

bath

pacify

bond

palate

beef

book

b

b

p



pants

balance

bind

paddle

boast

poke

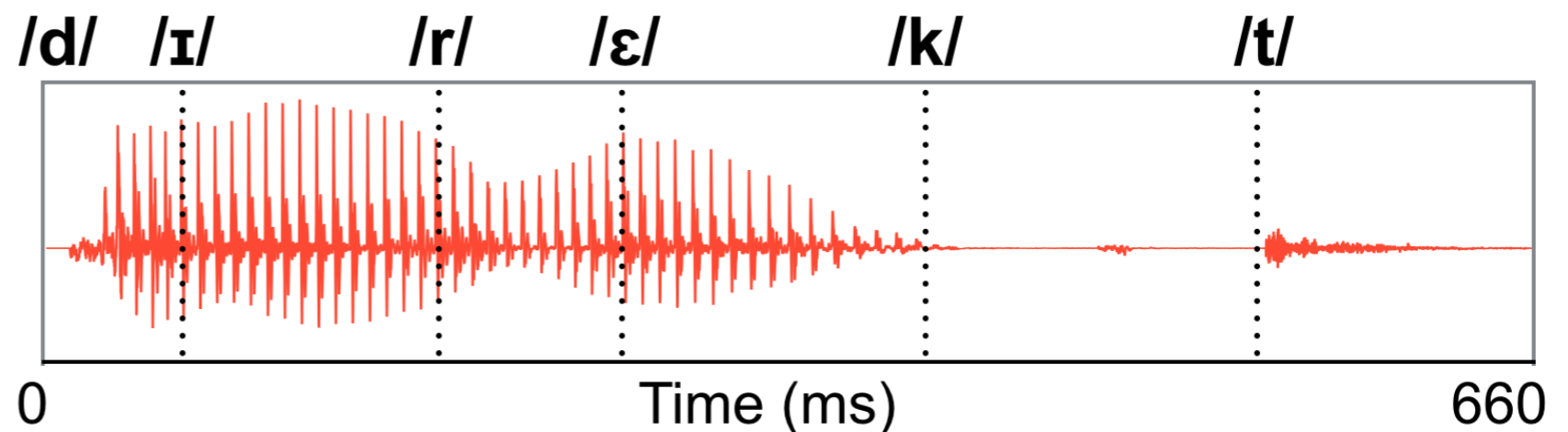
panda

Testing for phonological commitment

- **Surprisal:**

No commitment
Commitment

$$-\log_2 \left(P(\varphi_a | A) \frac{f(\varphi_a, \varphi_2, \dots, \varphi_t)}{f(\varphi_a, \varphi_2, \dots, \varphi_{t-1})} Q_a^t + \dots \right)$$



- **Entropy:**

No commitment
Commitment

$$P(w|C, A) = P(w|C_a) P(\varphi_a | A) + P(w|C_b) P(\varphi_b | A)$$

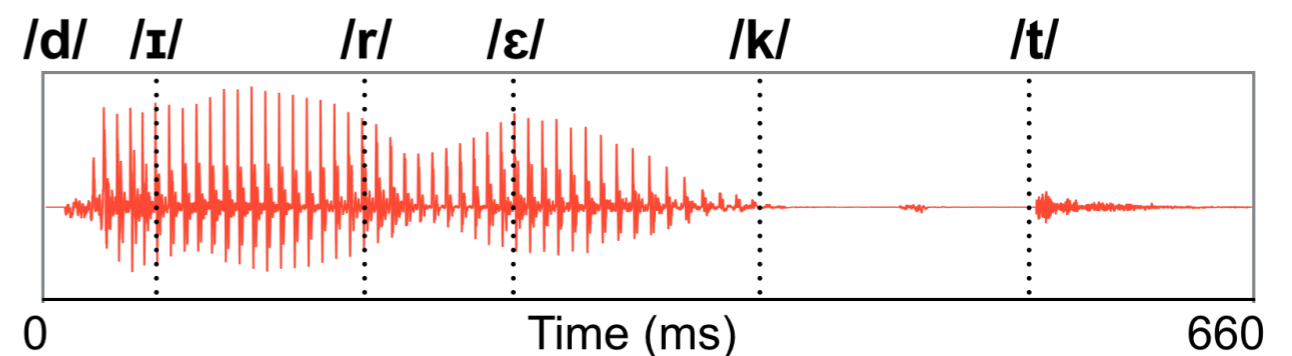
Model Setup

- **Critical variables:**

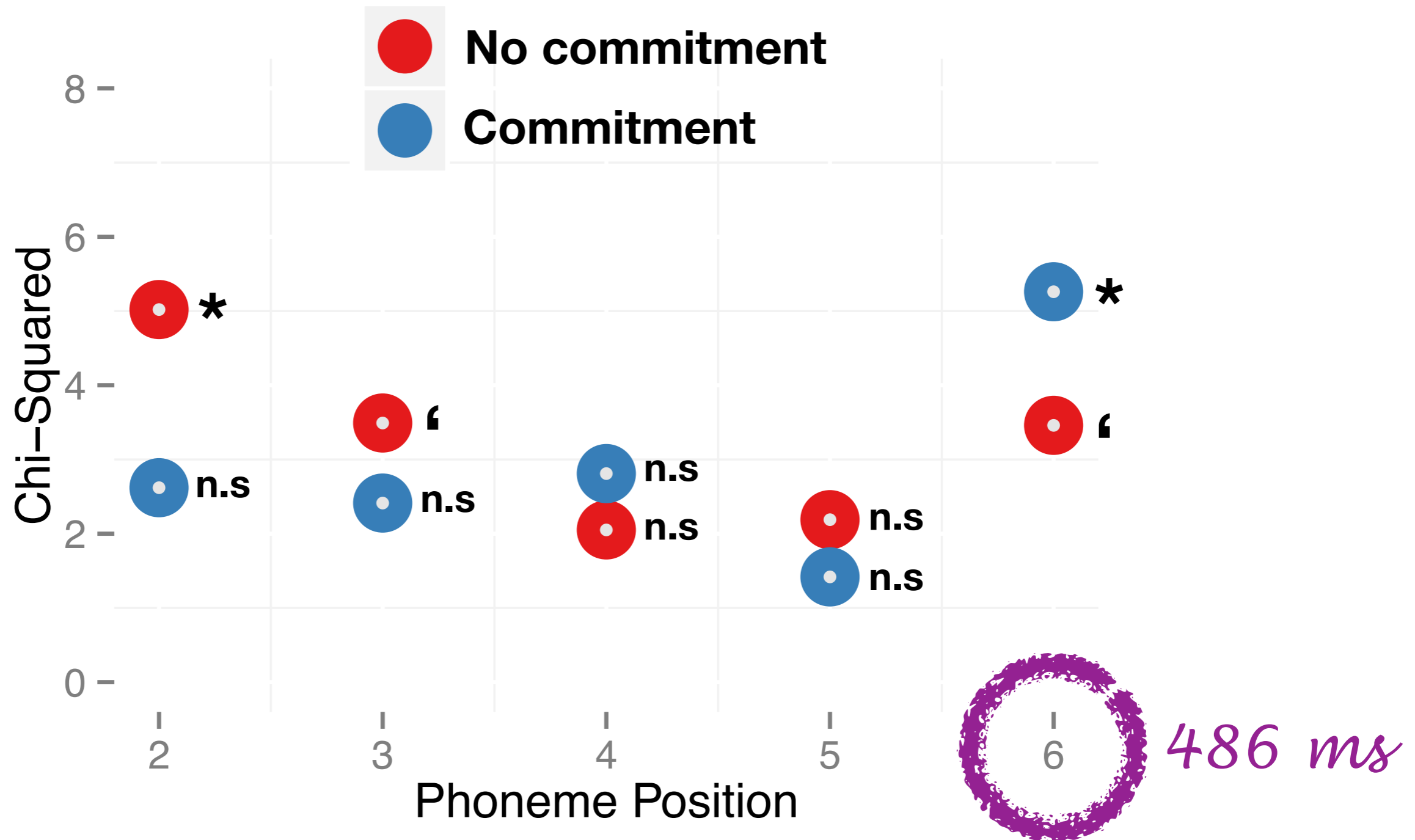
no commitment entropy
no commitment surprisal
commitment entropy
commitment surprisal

- **Control variables:**

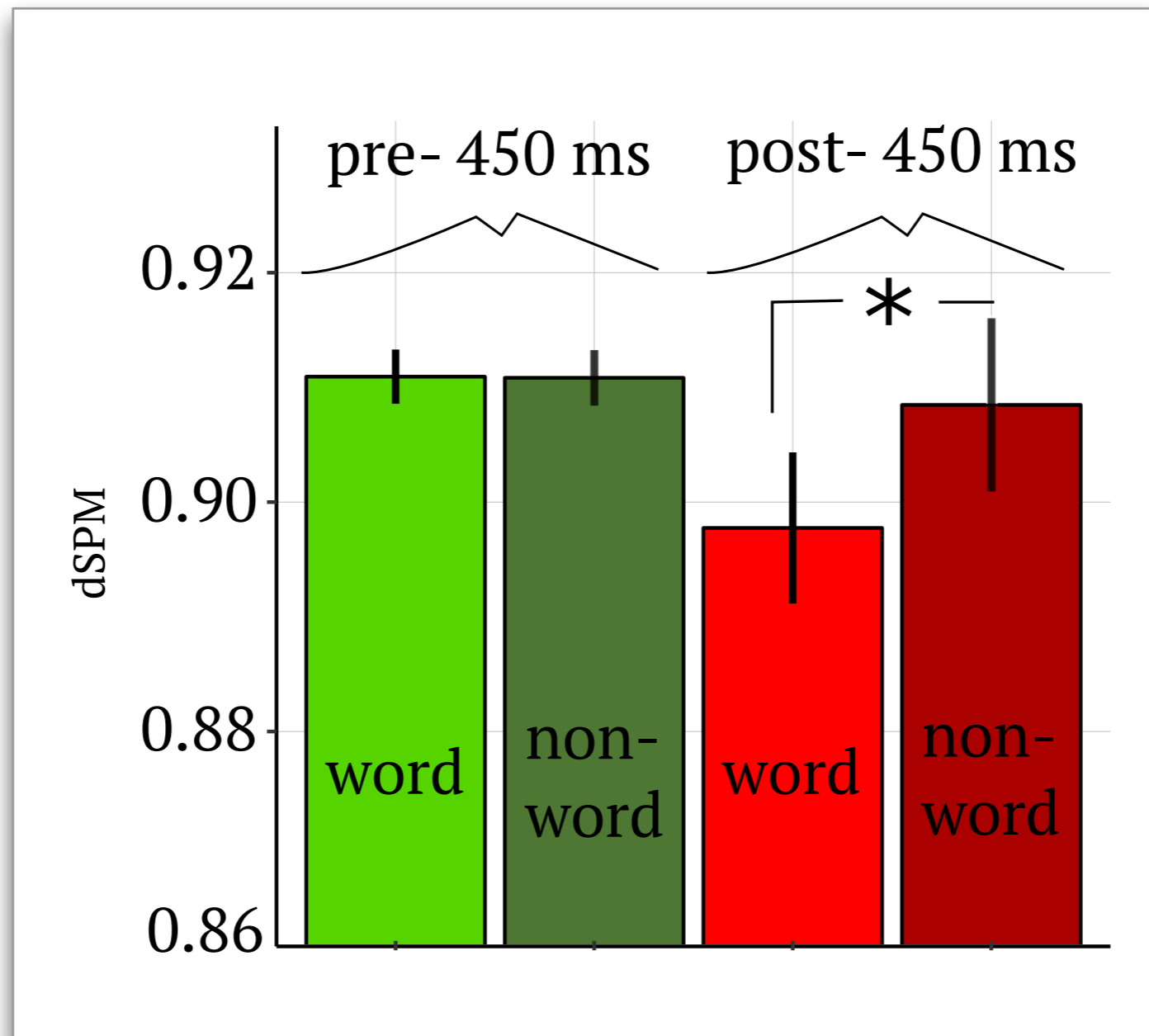
phoneme latency (ms)
phoneme latency (number of phonemes)
trial number
block number
stimulus amplitude
phoneme pair
ambiguity



Results



Further test of commitment



Interpretation

Processing hierarchy: Scott and Johnsrude, 2003; Hickock and Poeppel, 2004; Liebenthal et al., 2005; Rauschecker and Scott, 2009

lexical access



phonological commitment

/p/

b a r a k e e t

acoustic-phonetic maintenance

Interpretation

Processing is not purely feedforward, or feed “up”: TRACE model: McClelland and Elman, 1986; McMurray et al. 2009. cf. MERGE: Norris et al. 2000

lexical access



phonological commitment

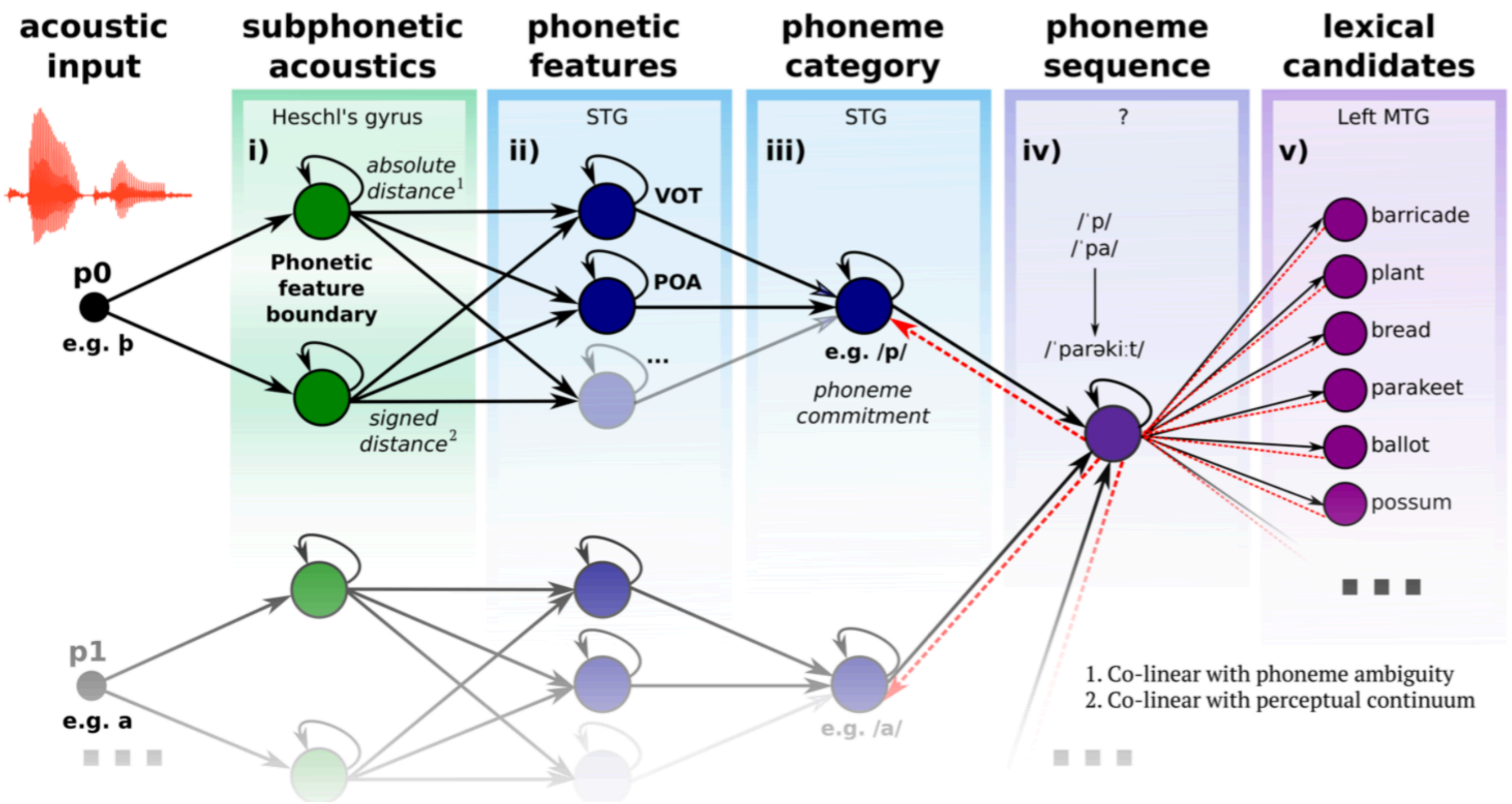
/p/

b a r a k e e

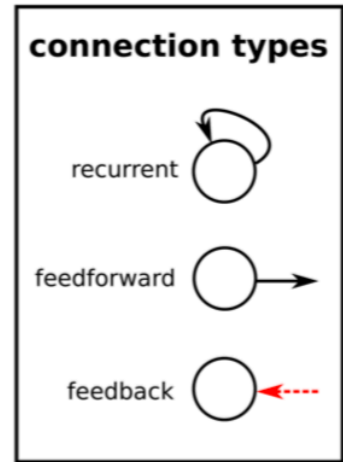
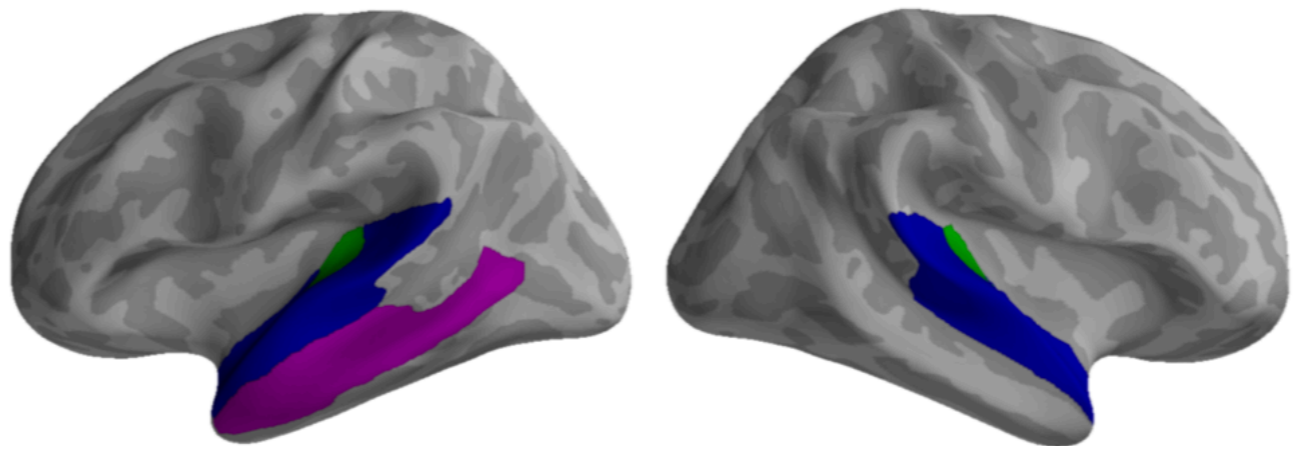
acoustic-phonetic maintenance

making commitment to a category does not cost — the system can flexibly avoid committing to a category; avoiding the function of exposure

but what happens if you continuously jump on a trampoline because you made a wrong choice? well, you fall when talking to someone with a different accent.



↓
phonemes
↓



Levels of analysis

1. Phonemes within words

- Responses to phoneme ambiguity (**bottom-up**)
- Neural signatures of ambiguity resolution, when provided with lexical information (**top-down**)

2. Words within sentences

- Which linguistic properties encoded in brain activity?
- What are the relative time-courses of processing each property?
- What is the computational architecture?

Fimas & Corbit (1972)

Cutler et al. (1986), Barry (1980-1984)

Taft & Forster (1975), Taft (1979)

Pinker & Prince (1988)

Marslen-Wilson & Welsh (1978)

sentence structure



phrasal structure

the fat cat dis | appear | ed

lemmas

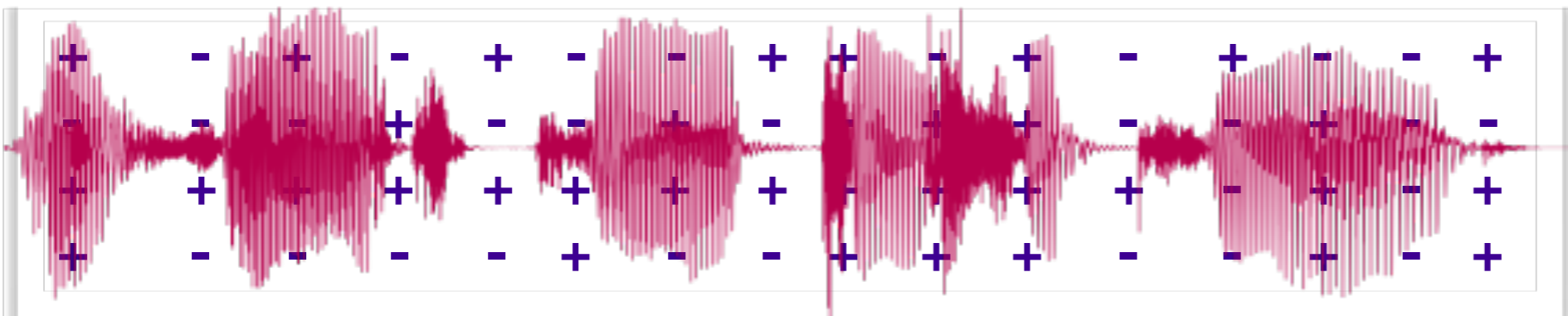
morphemes

dah fat kat dis ah pee ud

syllables

DH AH F AE T K AE T D IH S AH P IH R D

phonemes



phonetic features

acoustics

S

1) which linguistic units are encoded in brain activity?

2) what is the relative time-course?



VP

appear | ed

phrasal structure

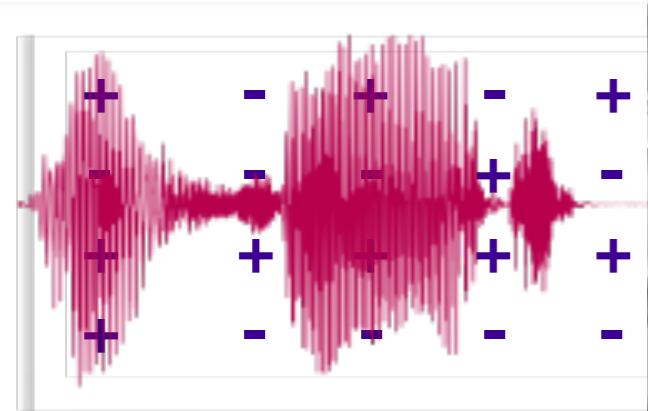
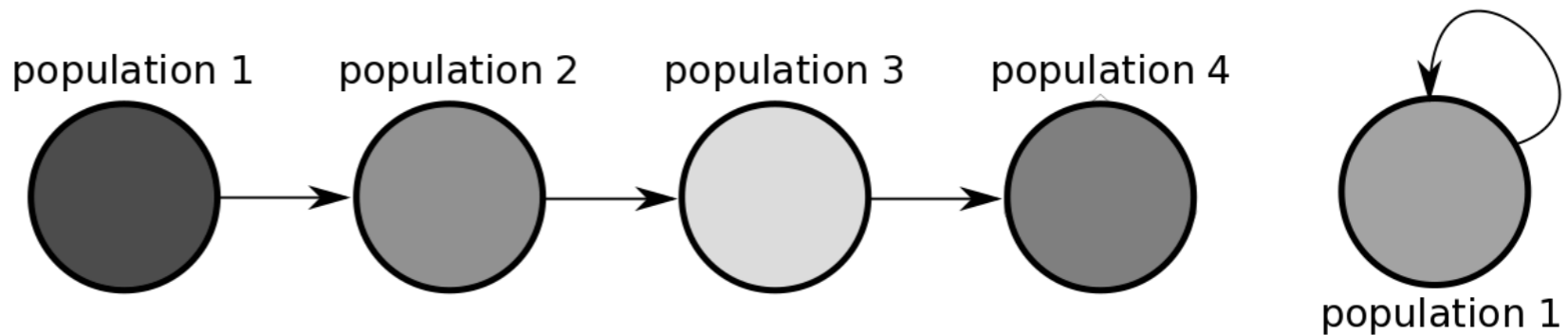
lemmas

morphemes

dah fa

DH AH F AE T

3) what is the computational architecture?



Setup

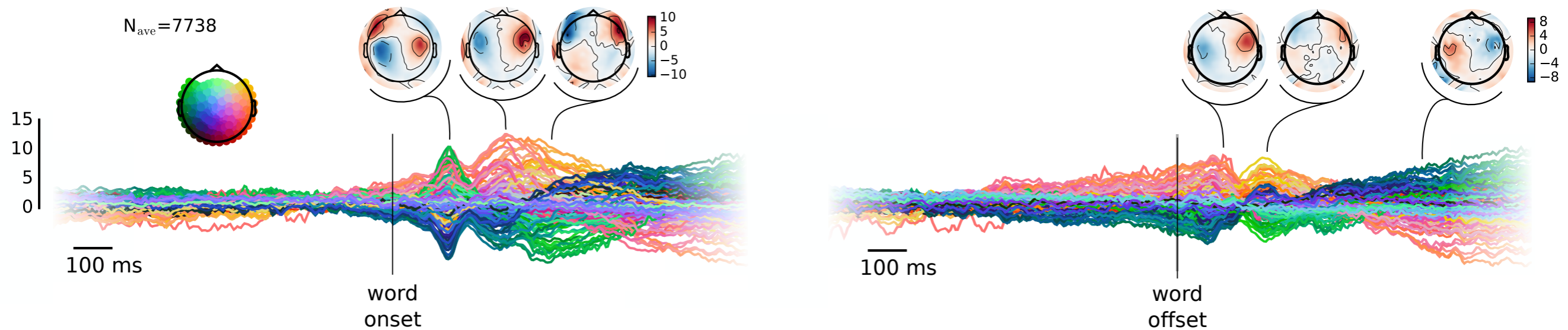
- 18 participants
 - Listening to four narrative stories (twice)
 - 2 x one hour recordings
 - KIT 208 channel MEG system
 - Engagement task
-
- ~40,000 phonemes per participant
 - ~8,000 words per participant



جامعة نيويورك أبوظبي

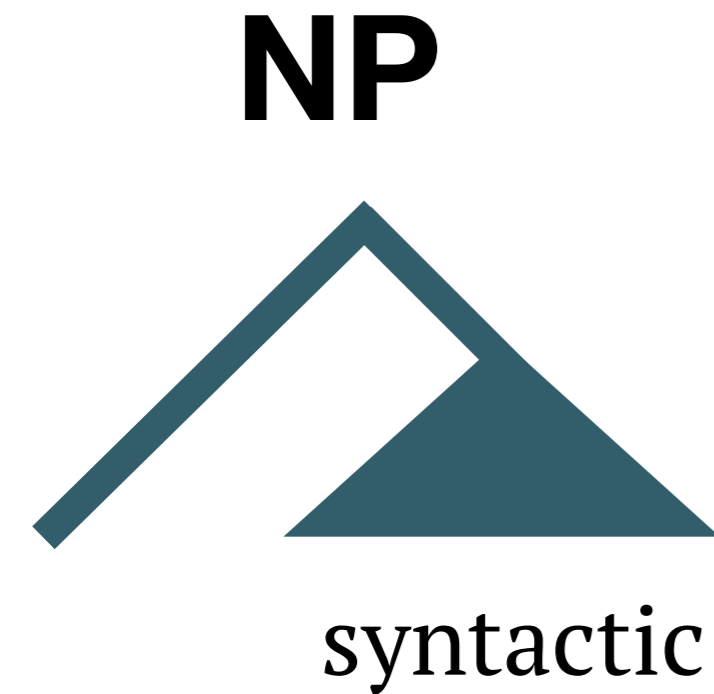
 NYU | ABU DHABI

Event-locked average response



explain this plot/ colour of channels in the topography

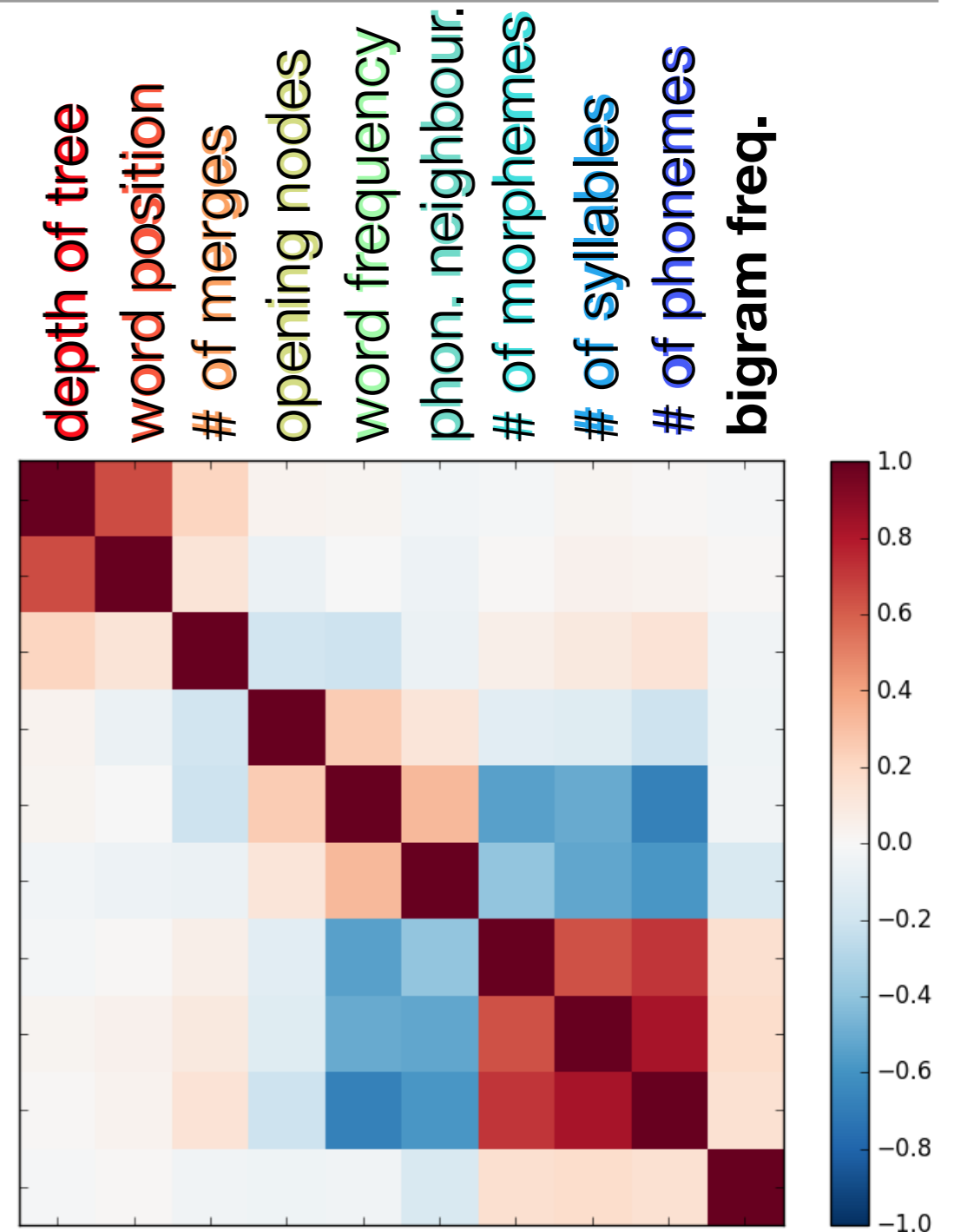
Stimulus features



word use

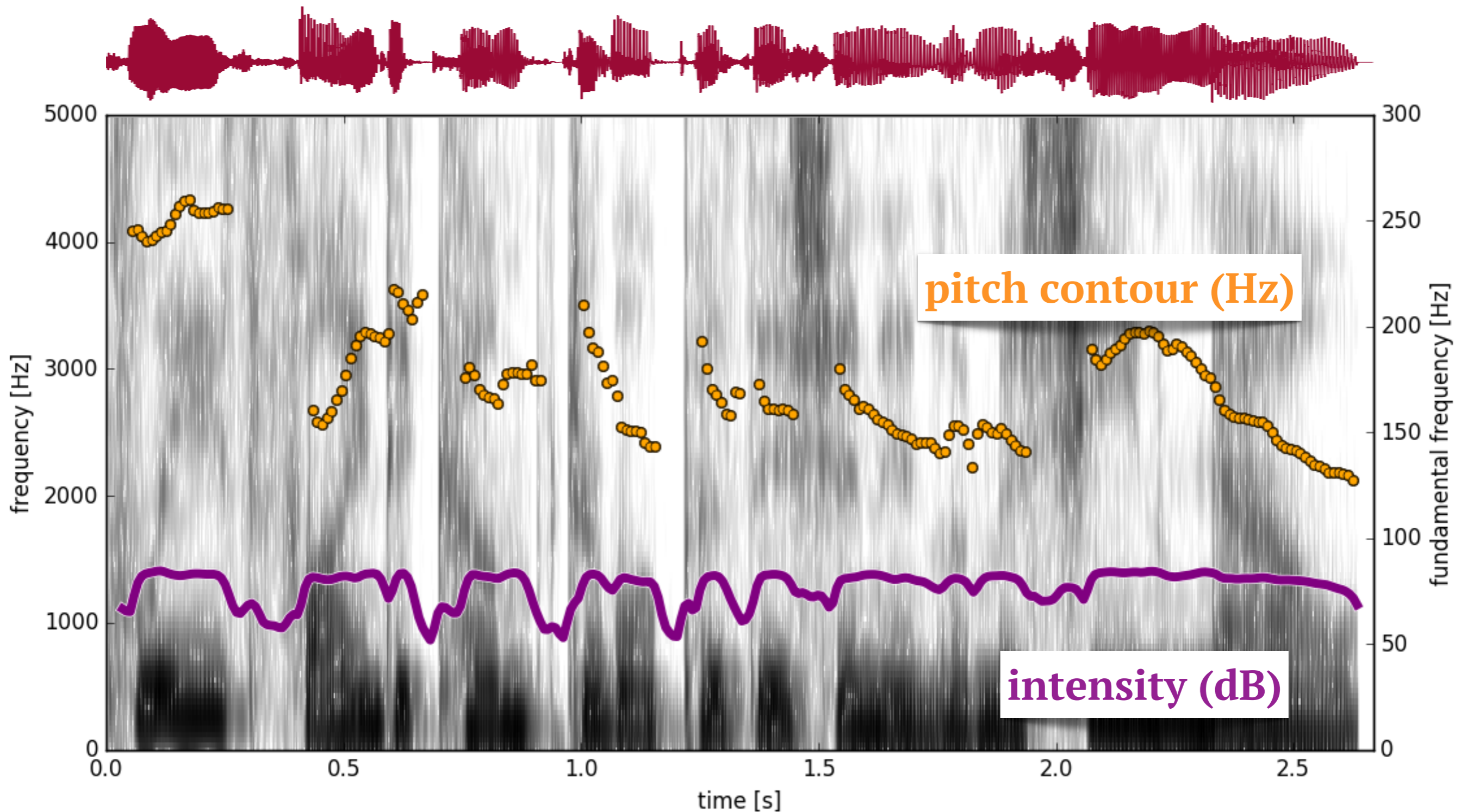
word structure

- depth of tree
- word position
- # of merges
- opening nodes
- word frequency
- phon. neighbourhood
- # of morphemes
- # of syllables
- # of phonemes
- bigram frequency

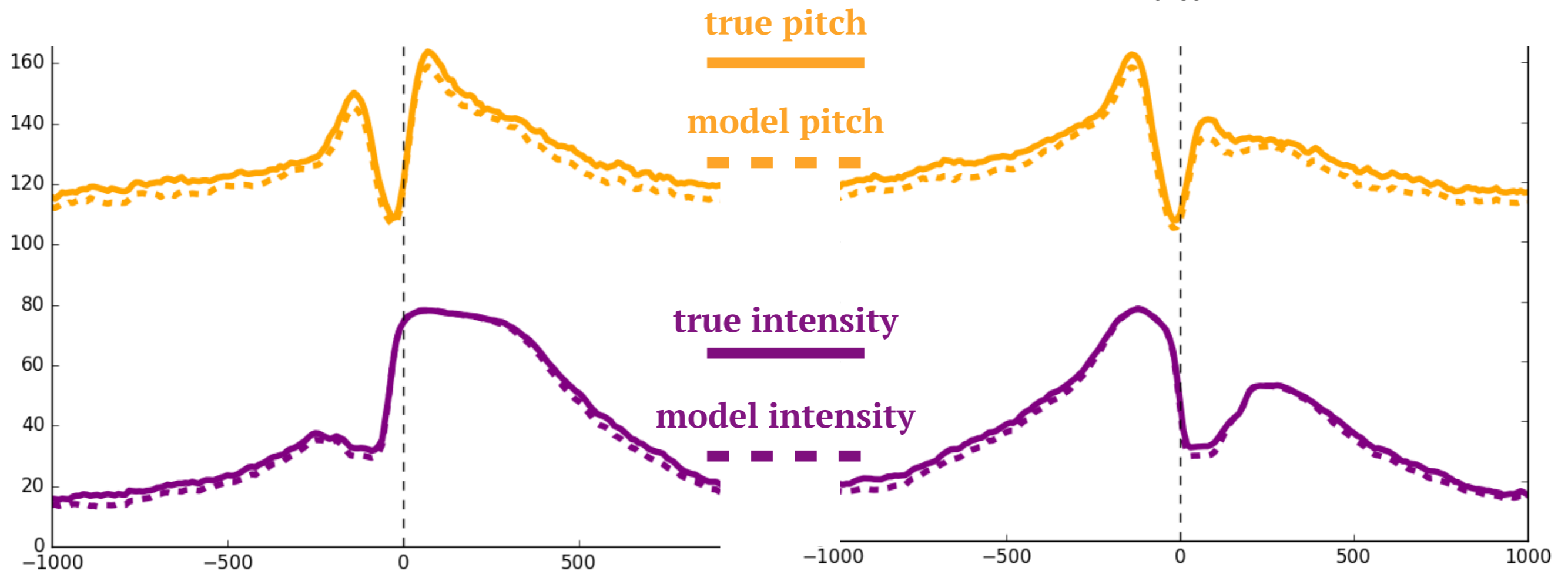
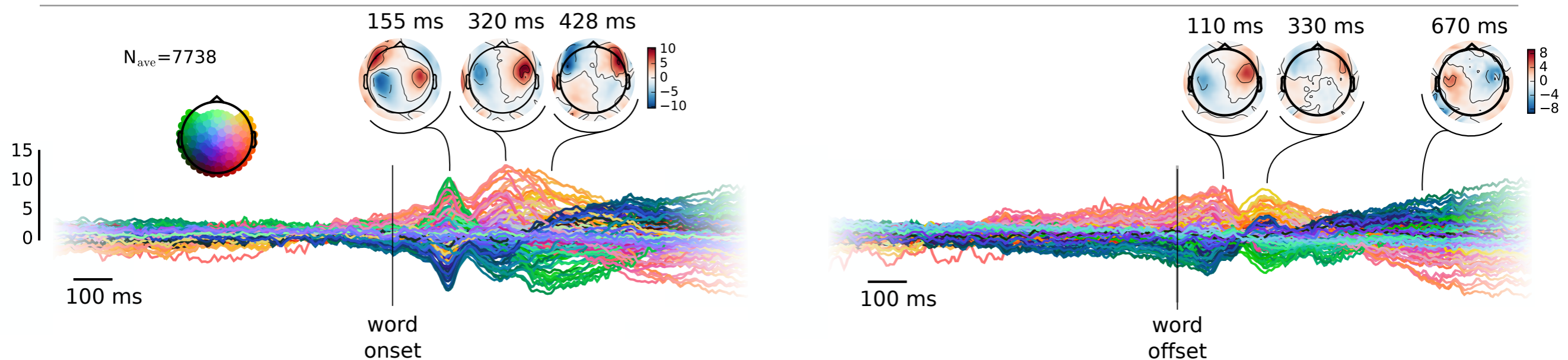


stimulus similarity matrix

Pitch and intensity in the acoustic signal

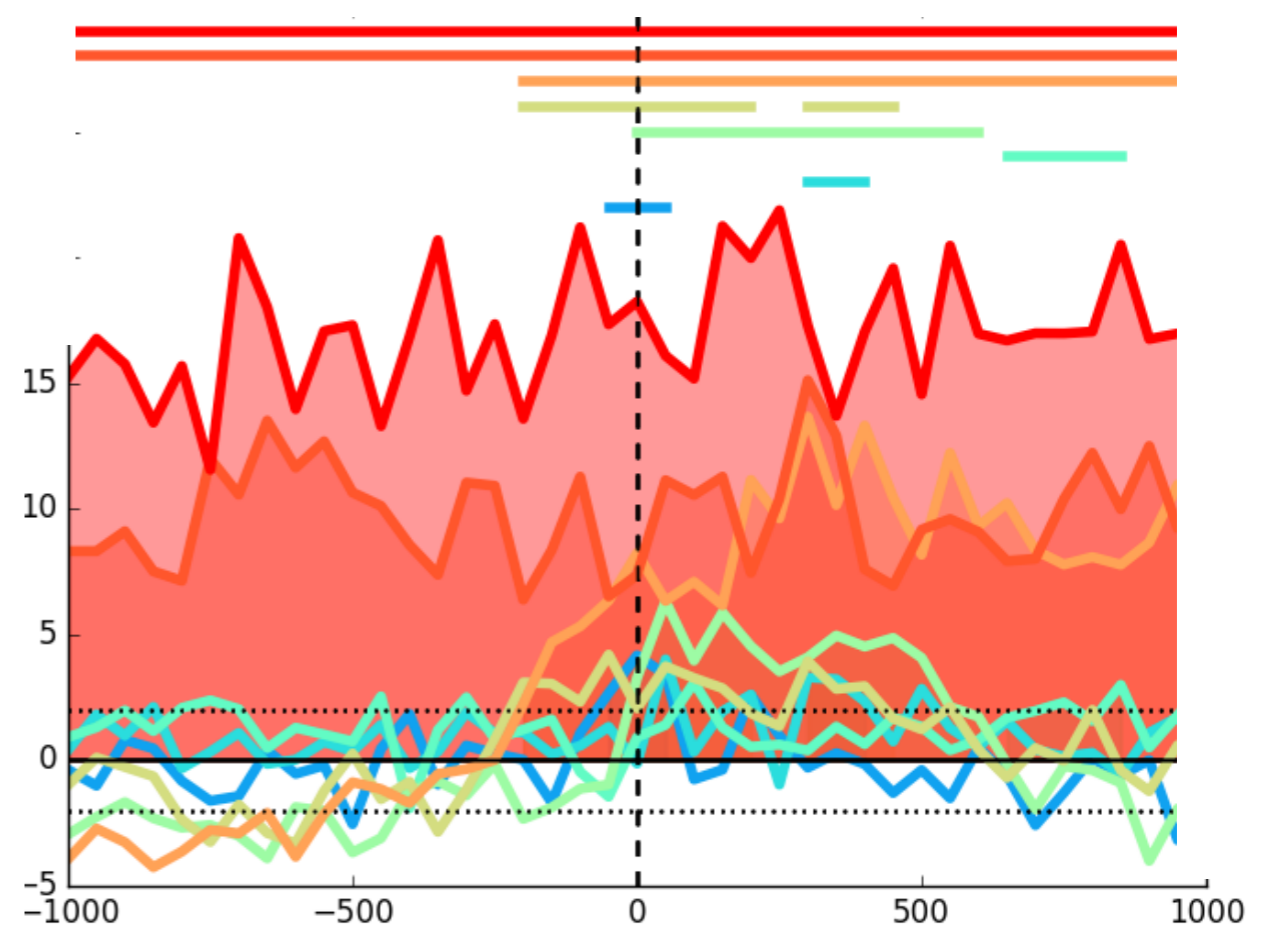
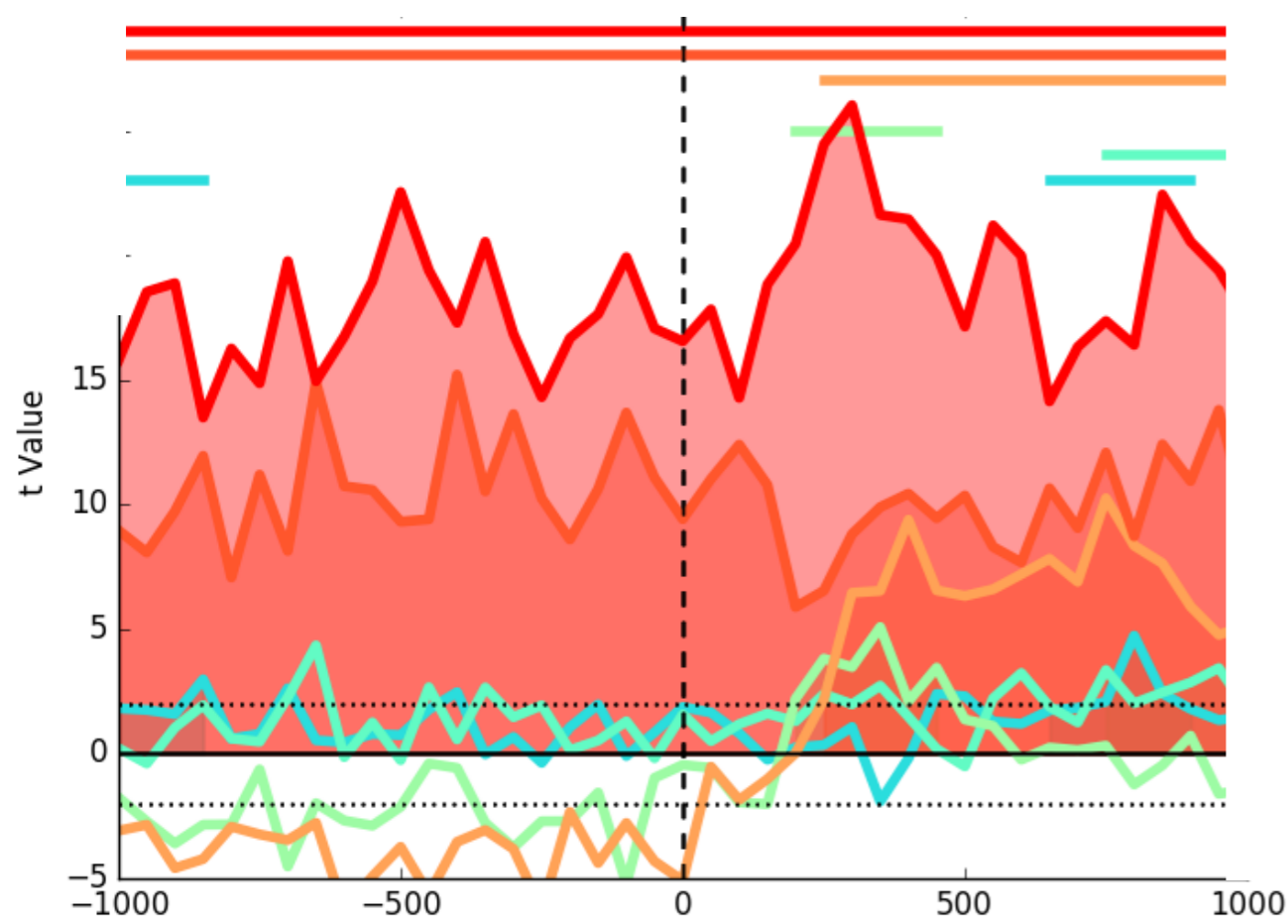
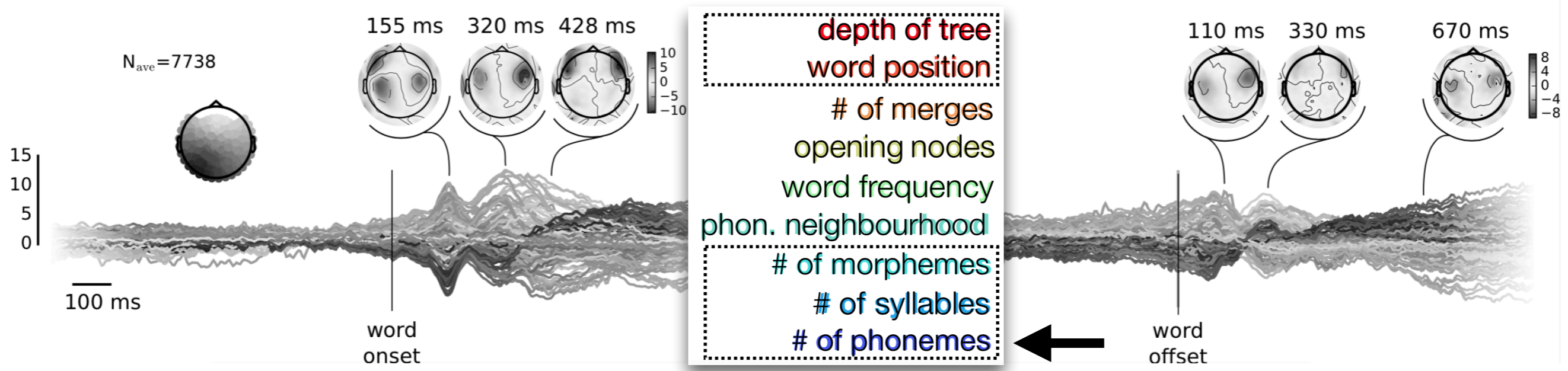


Pitch and intensity



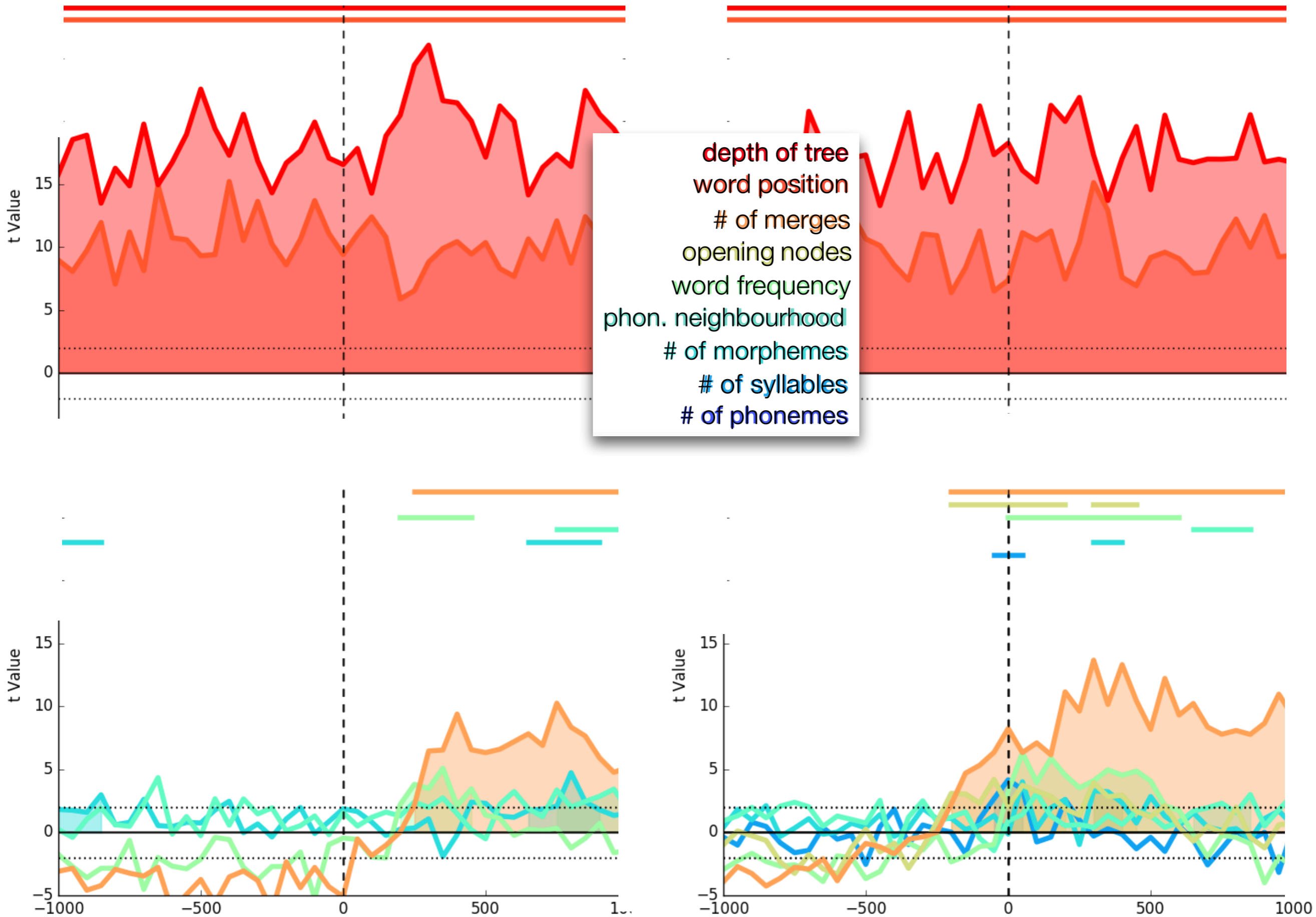
**which linguistic units are
encoded in brain activity?**

Deconfounded decoding accuracy



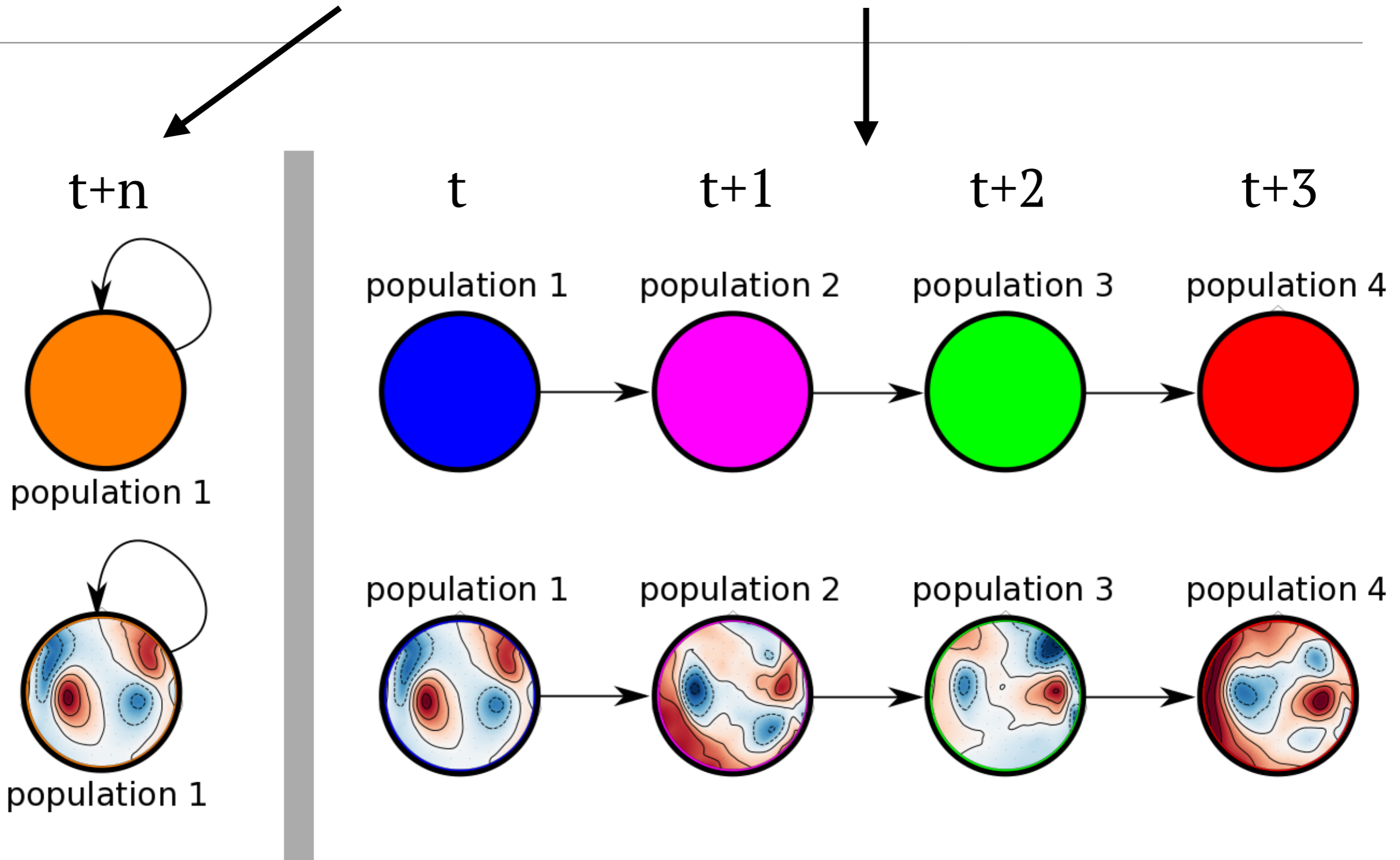
**what are the relative
time-courses?**

Deconfounded decoding accuracy

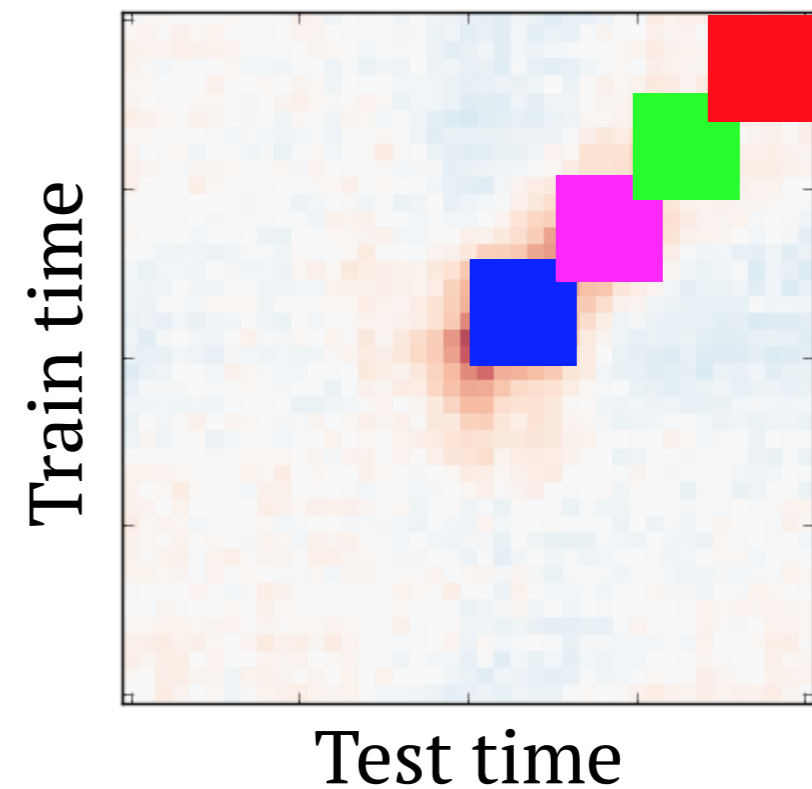
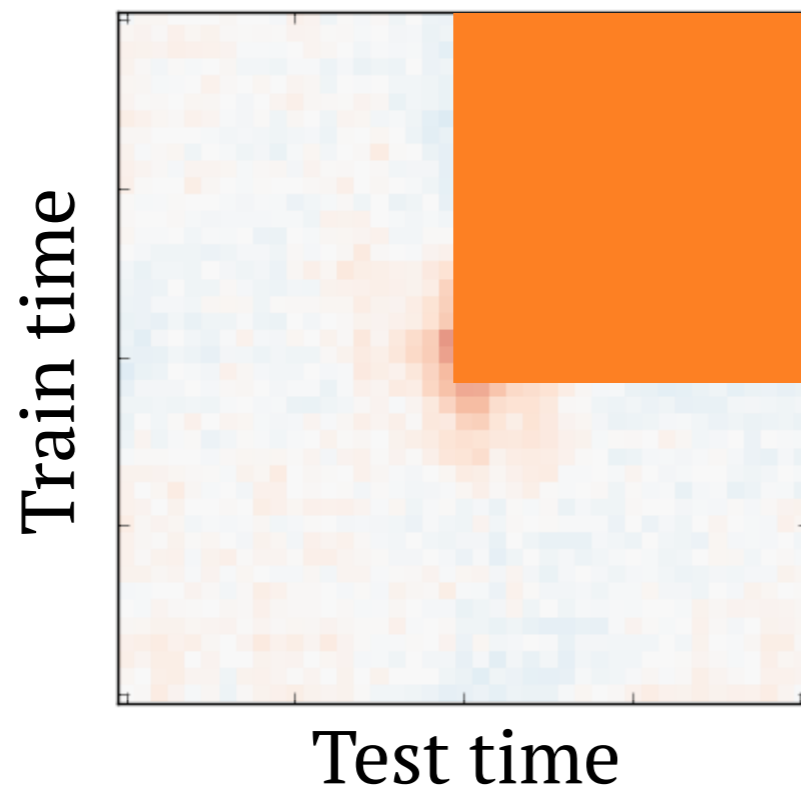
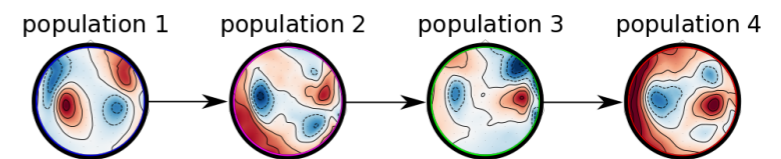
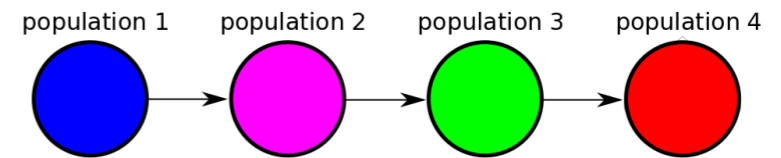
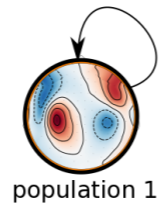
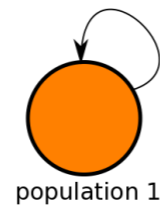


**what is the computational
architecture?**

Recurrent vs. Feedforward

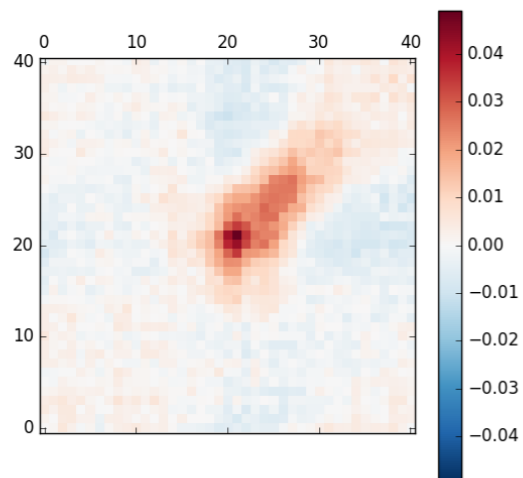


Recurrent vs. Feedforward

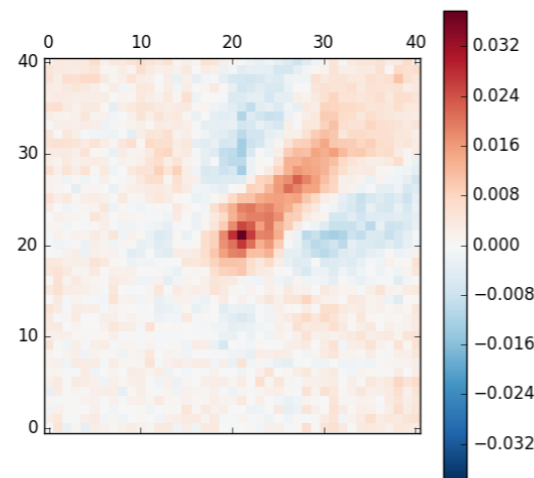


Recurrent vs. Feedforward

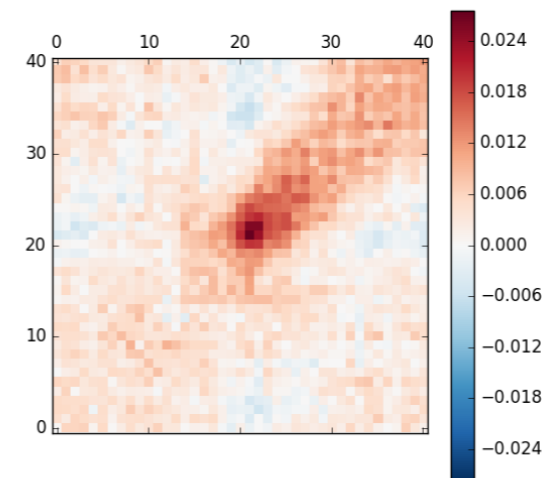
of syllables



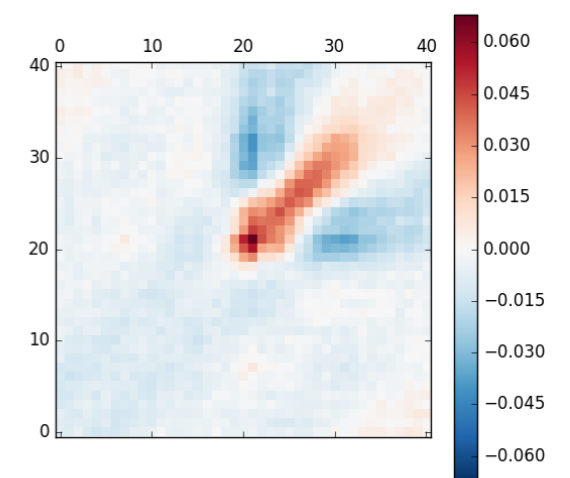
of morphemes



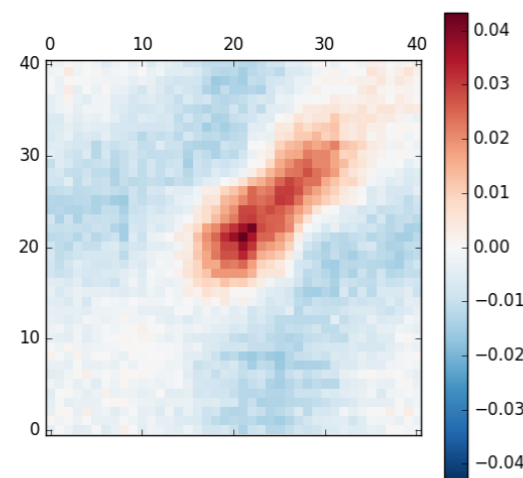
phon. neighbourhood



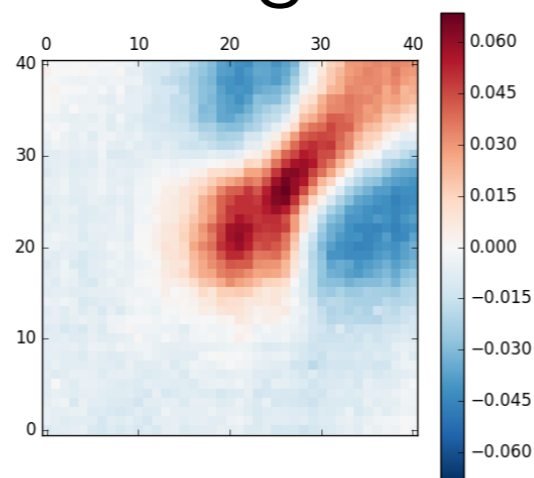
word frequency



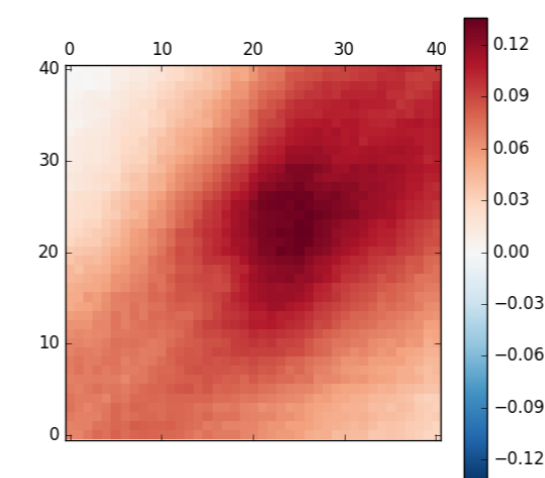
opening nodes



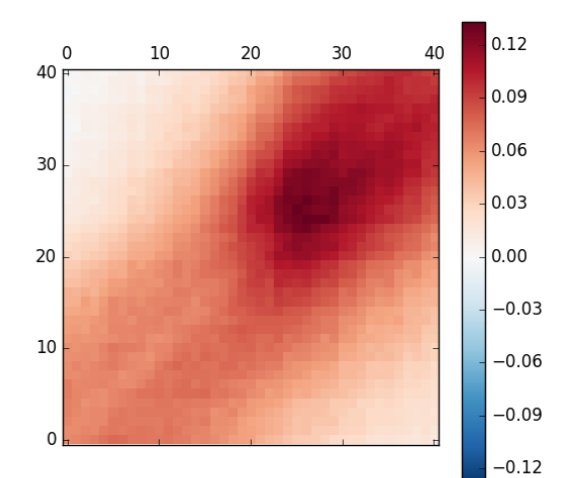
of merges



word number



depth



Discussion

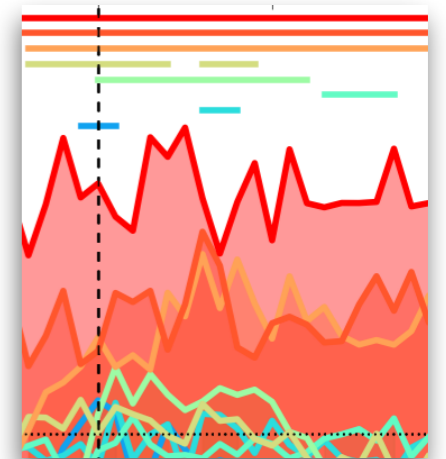
(1) Which linguistic units are encoded?

- Multiple features, **spanning the hierarchy**
- Including # of **syllables**; # of **morphemes**

depth of tree
word position
of merges
opening nodes
word frequency
phon. neighbourhood
of morphemes
of syllables

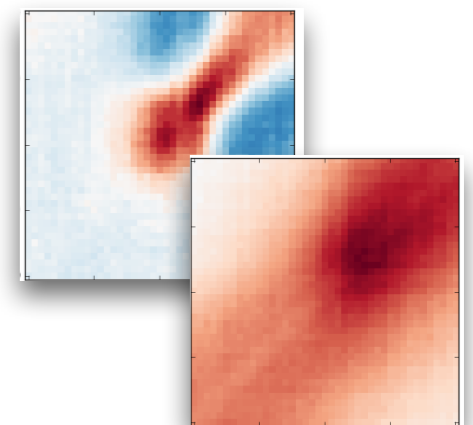
(2) What is the relative time-course?

- Overall a highly **parallel** architecture



(3) What is the computational architecture?

- Both **feedforward** and **recurrent** computations, depending on the linguistic property



Levels of analysis

1. Phonemes within words

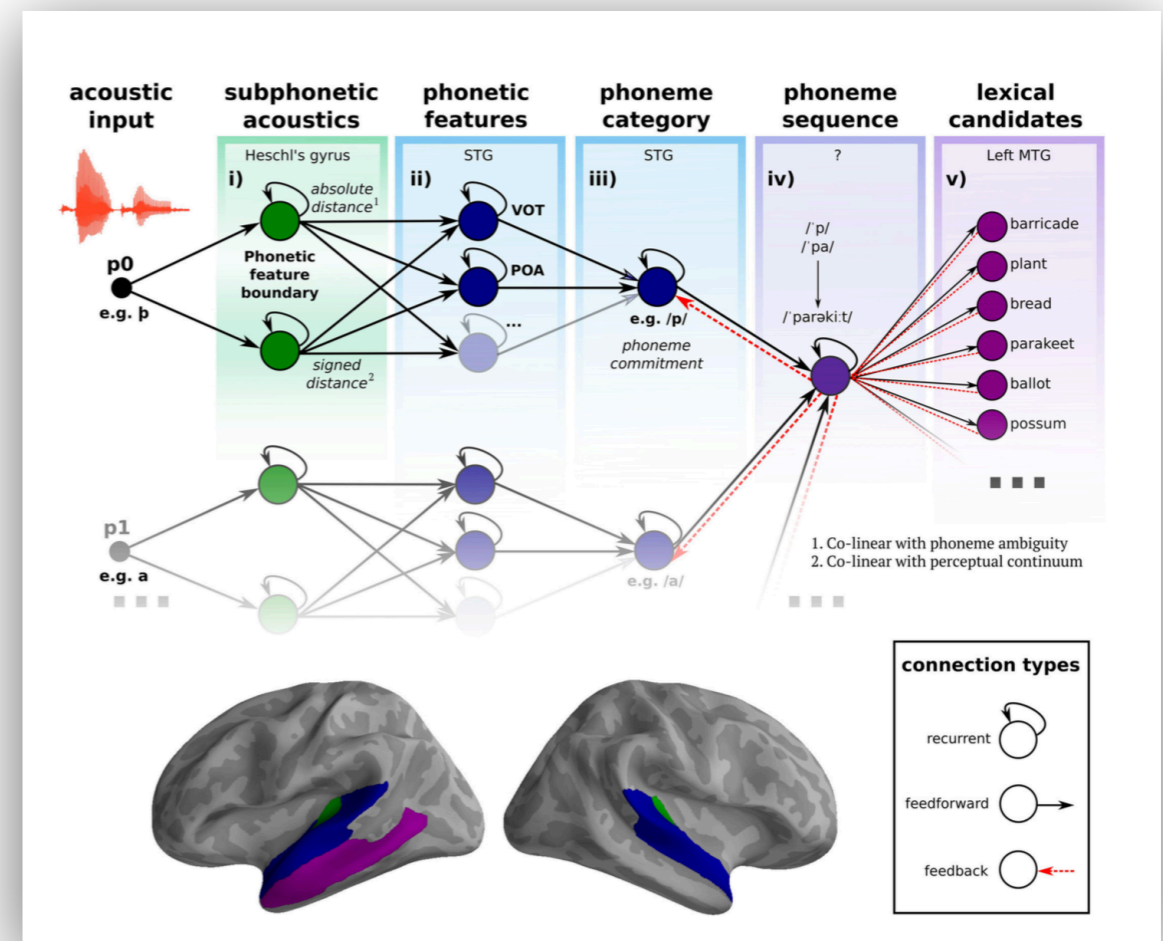
- Responses to phoneme ambiguity, phonetic features and acoustic properties (**bottom-up**)
- Neural signatures of ambiguity resolution, when provided with lexical information (**top-down**)

2. Words within sentences

- Which linguistic properties encoded in brain activity?
- What are the relative time-courses of processing each property?
- What is the computational architecture?

Take home message

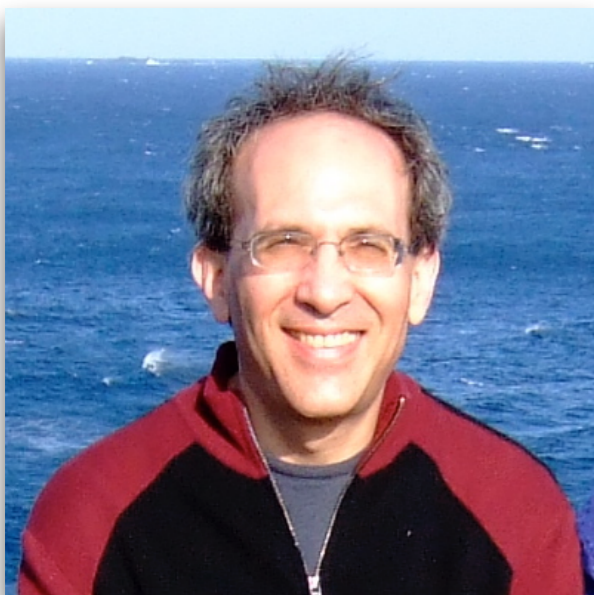
- **Feedforward, recurrent, and highly parallel processes** may be a general architecture supporting speech comprehension
 - words within sentences
 - phonemes within words



✉ laura.gwilliams@nyu.edu
🐦 [@GwilliamsL](https://twitter.com/GwilliamsL)

With big thanks to:

- My supervisors, **Alec Marantz** and **David Poeppel**, as well as everyone in the **Neuroscience of Language Lab** and **Poeppel Lab**!



Funding: G1001 Abu Dhabi Institute

Laura Gwilliams | [New York University](https://www.nyu.edu) | [@GwilliamsL](https://twitter.com/GwilliamsL)



NEW YORK UNIVERSITY

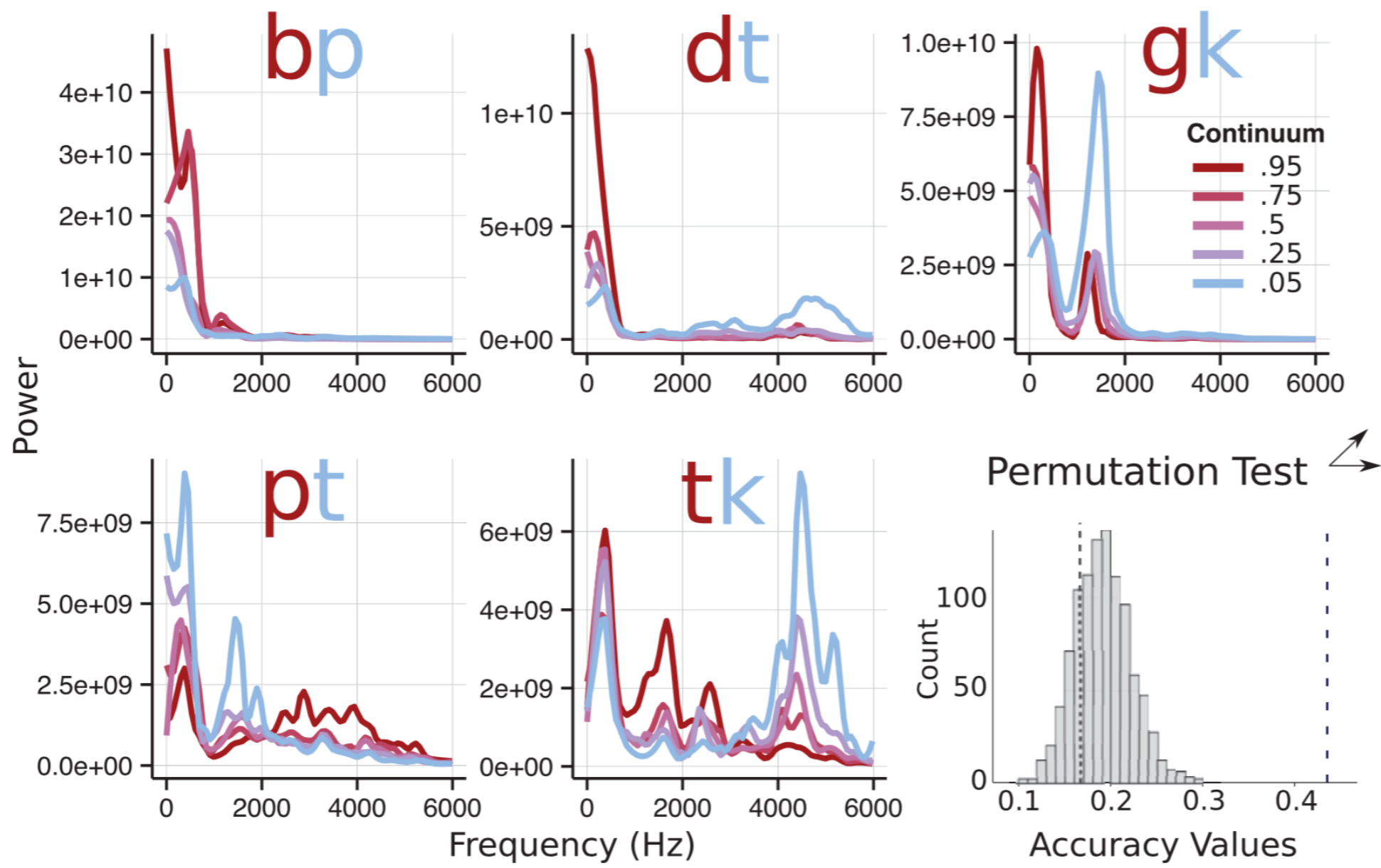
✉ laura.gwilliams@nyu.edu

🐦 [@GwilliamsL](https://twitter.com/GwilliamsL)

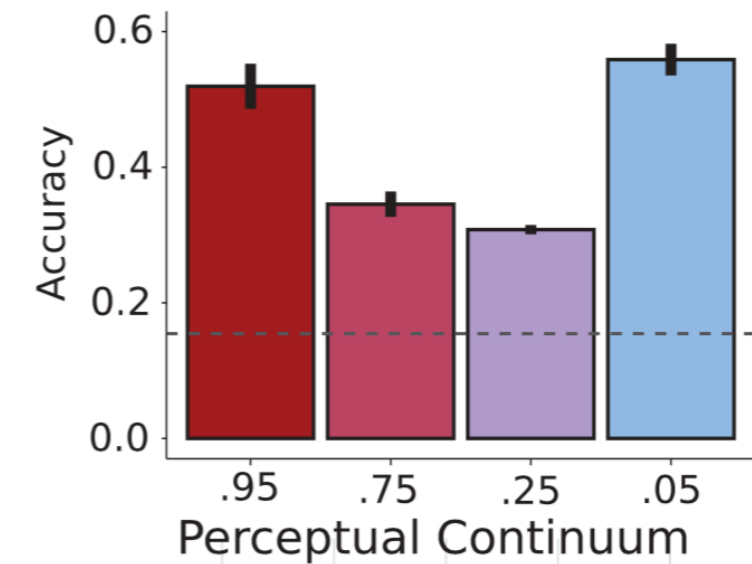
Thank you!



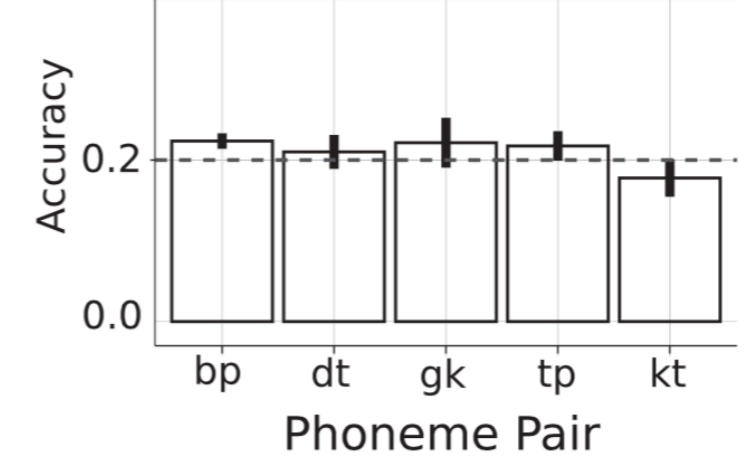
a Power Spectrum of Noise Burst (first 20 ms)



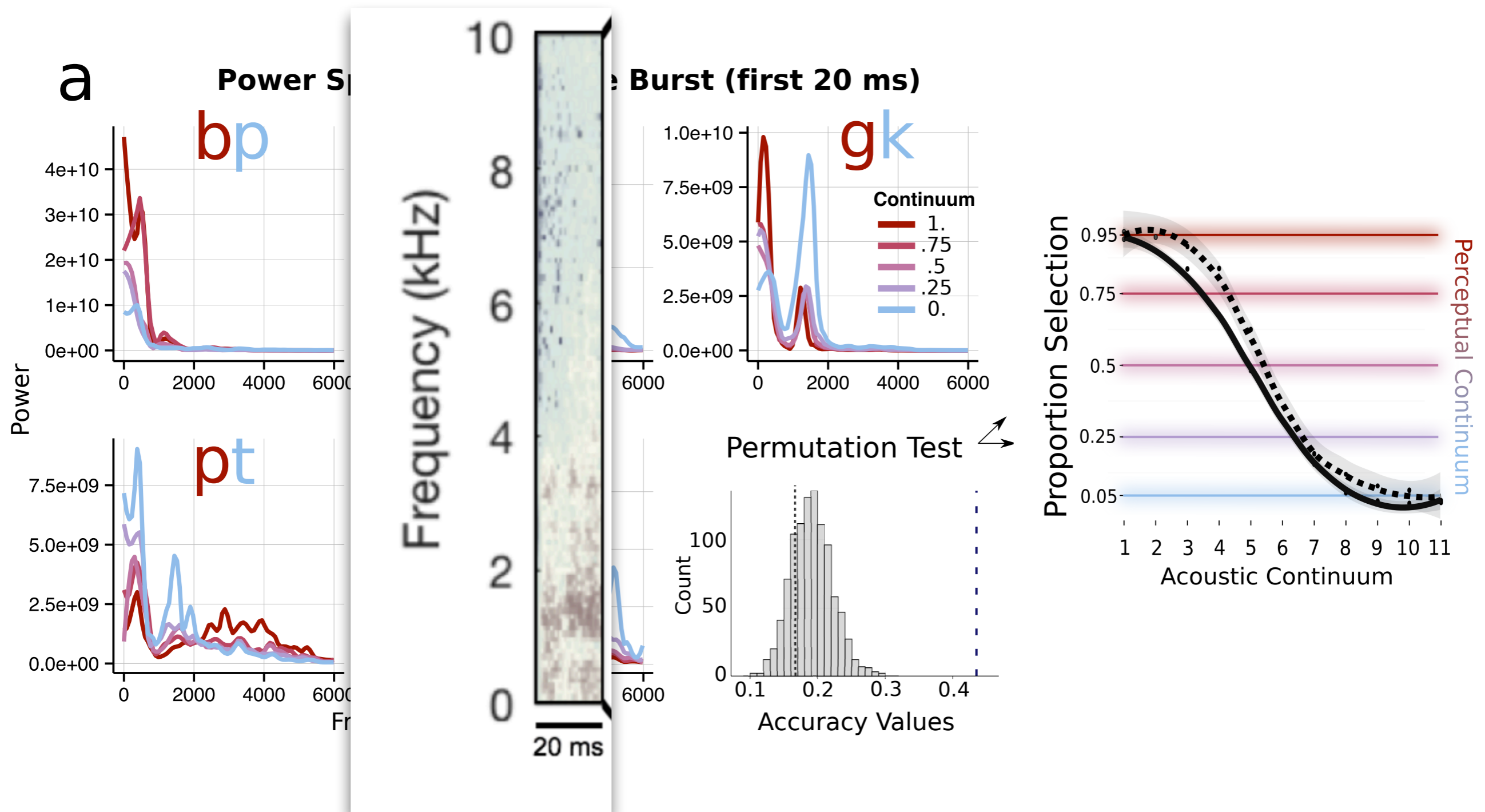
b Decode Phoneme



c Decode Step

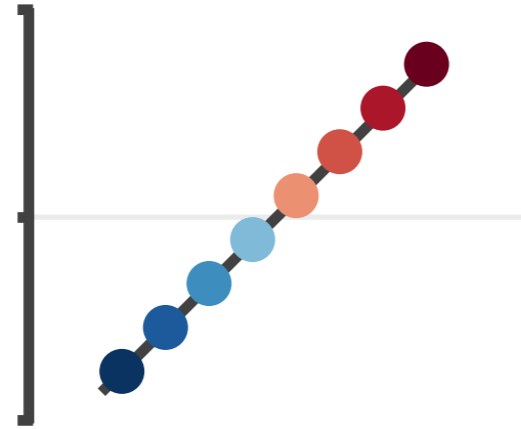


Is ambiguity correlated with acoustic properties?

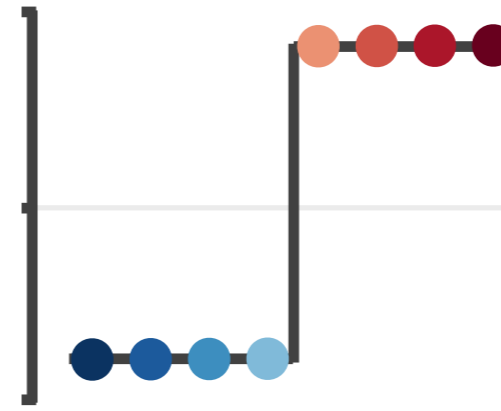


Predictive Coding

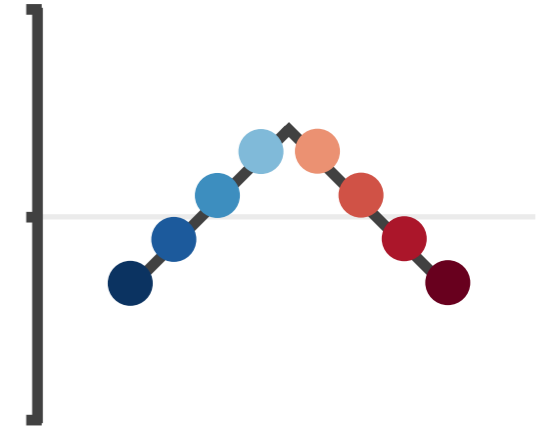
Linear Evidence



Categorical Percept

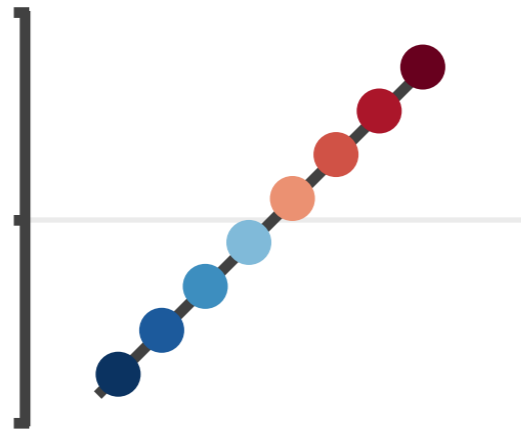


Ambiguity

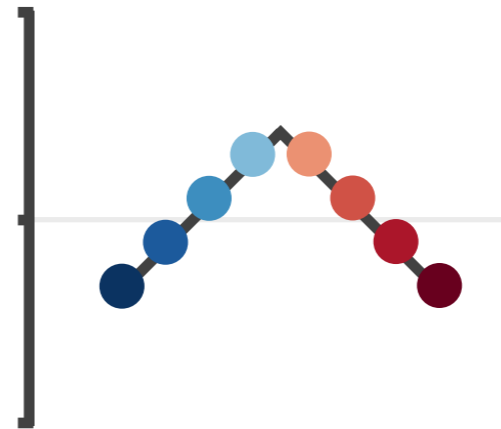


Neutralisation

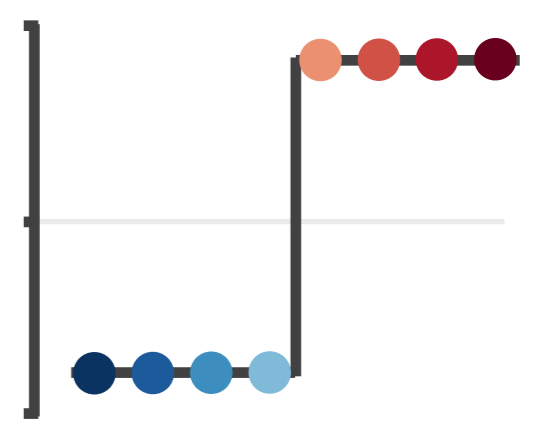
Linear Evidence



Ambiguity

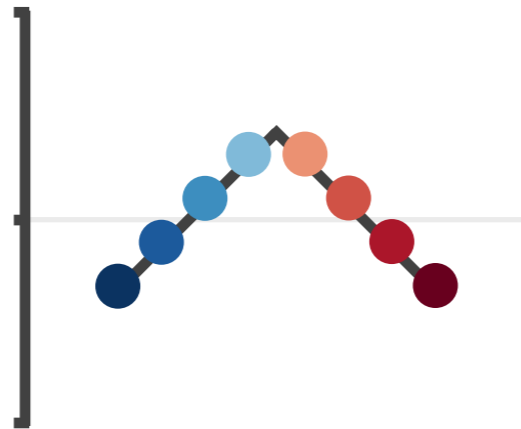


Categorical Percept

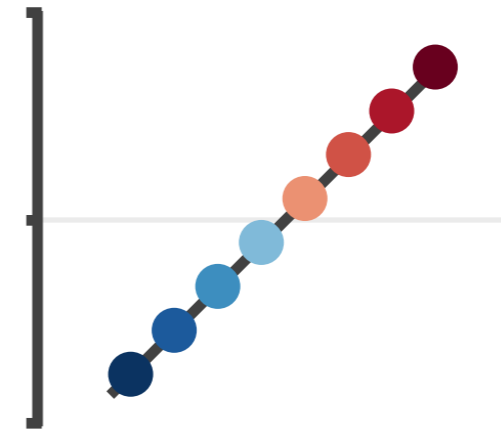


Cut-through connection

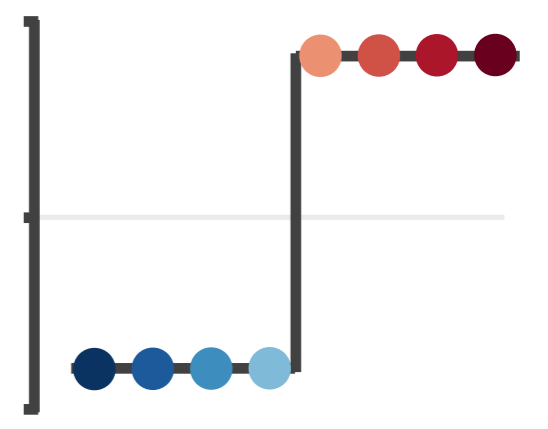
Ambiguity



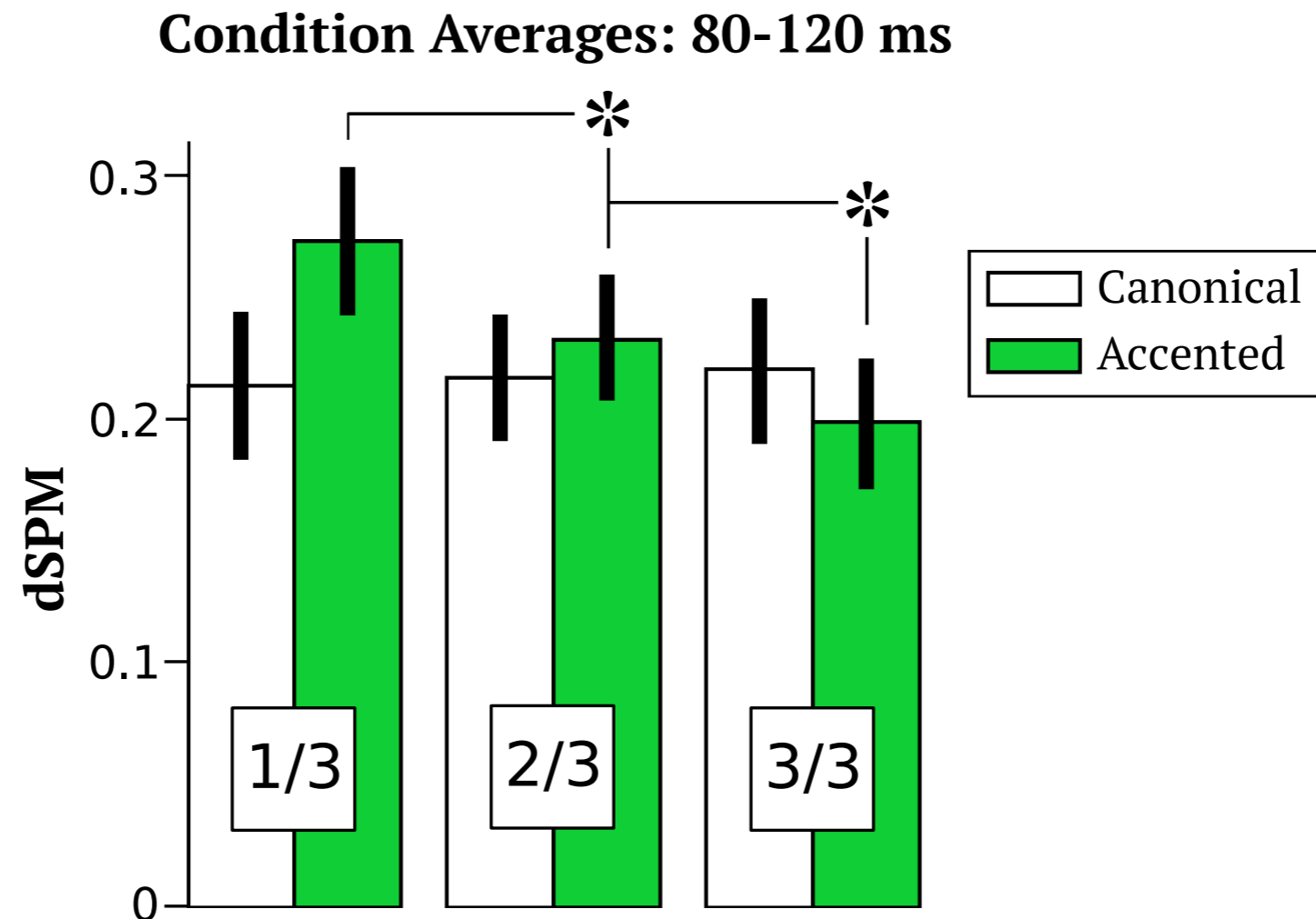
Linear Evidence



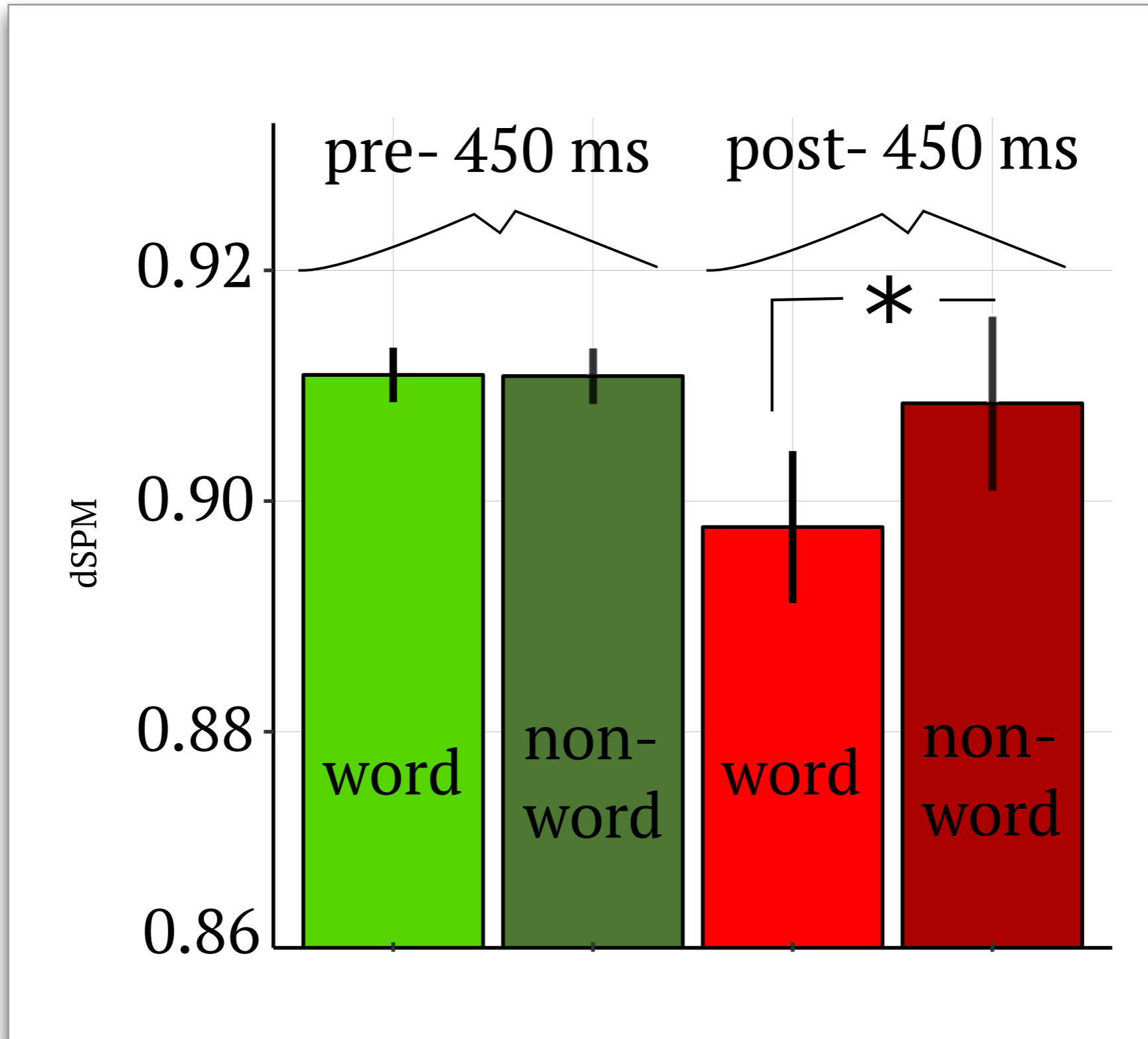
Categorical Percept



Interpretation

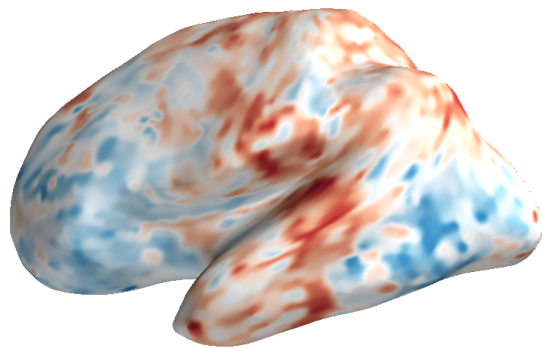


- Attunement is proposed to involve **re-tuning perceptual boundaries** between phonological categories (Norris et al., 2003; Kraljic and Samuel, 2005, 2006, 2007; Maye et al., 2008; see Samuel & Kraljic, 2009 for a review)



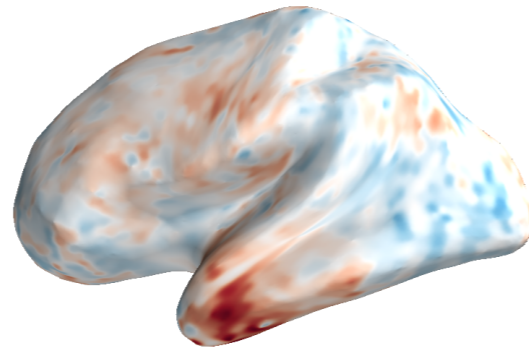
(Preliminary) Localisation of effects

#syllables



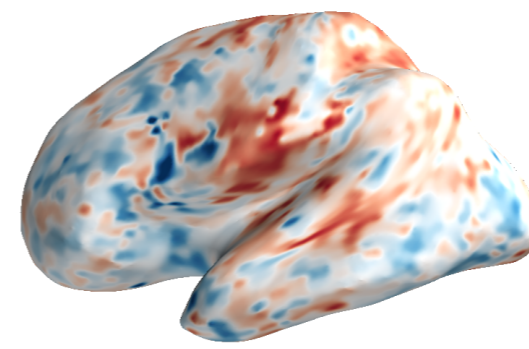
-0.0512 -0.0366 -0.0219 -0.00731 0.00731 0.0219 0.0366 0.0512

#morphemes



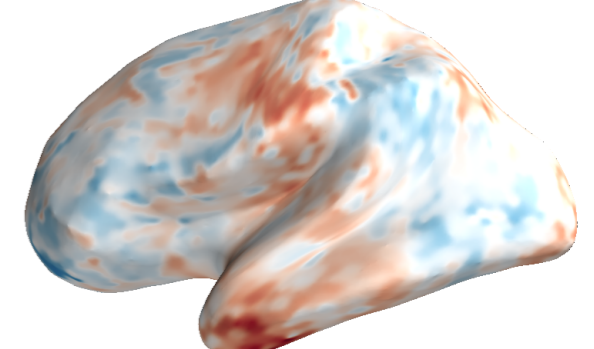
0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000

phon.
neighborhood



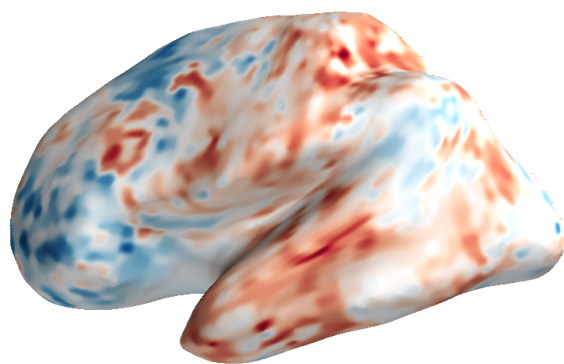
-0.0474 -0.0339 -0.0203 -0.00677 0.00677 0.0203 0.0339 0.0474

word
frequency



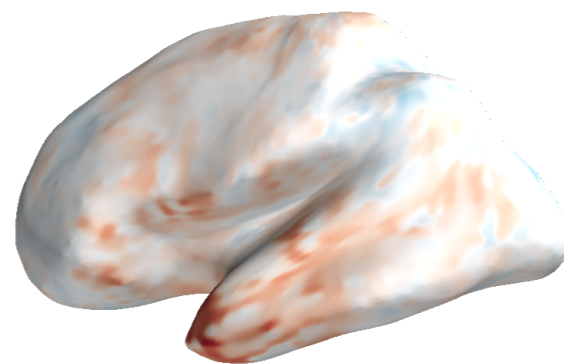
0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000

#opening nodes



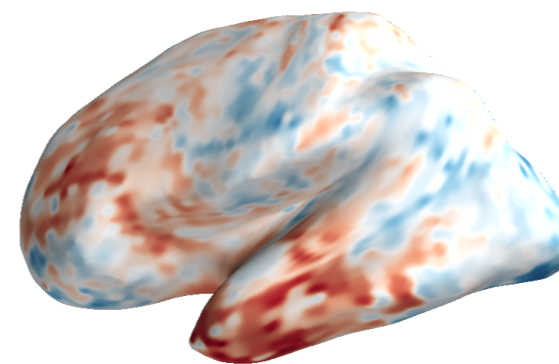
-0.0511 -0.0365 -0.0219 -0.00730 0.00730 0.0219 0.0365 0.0511

#merge



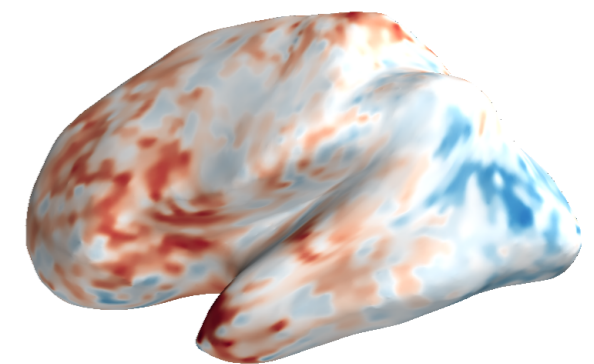
-0.115 -0.0818 -0.0491 -0.0164 0.0164 0.0491 0.0818 0.115

word number



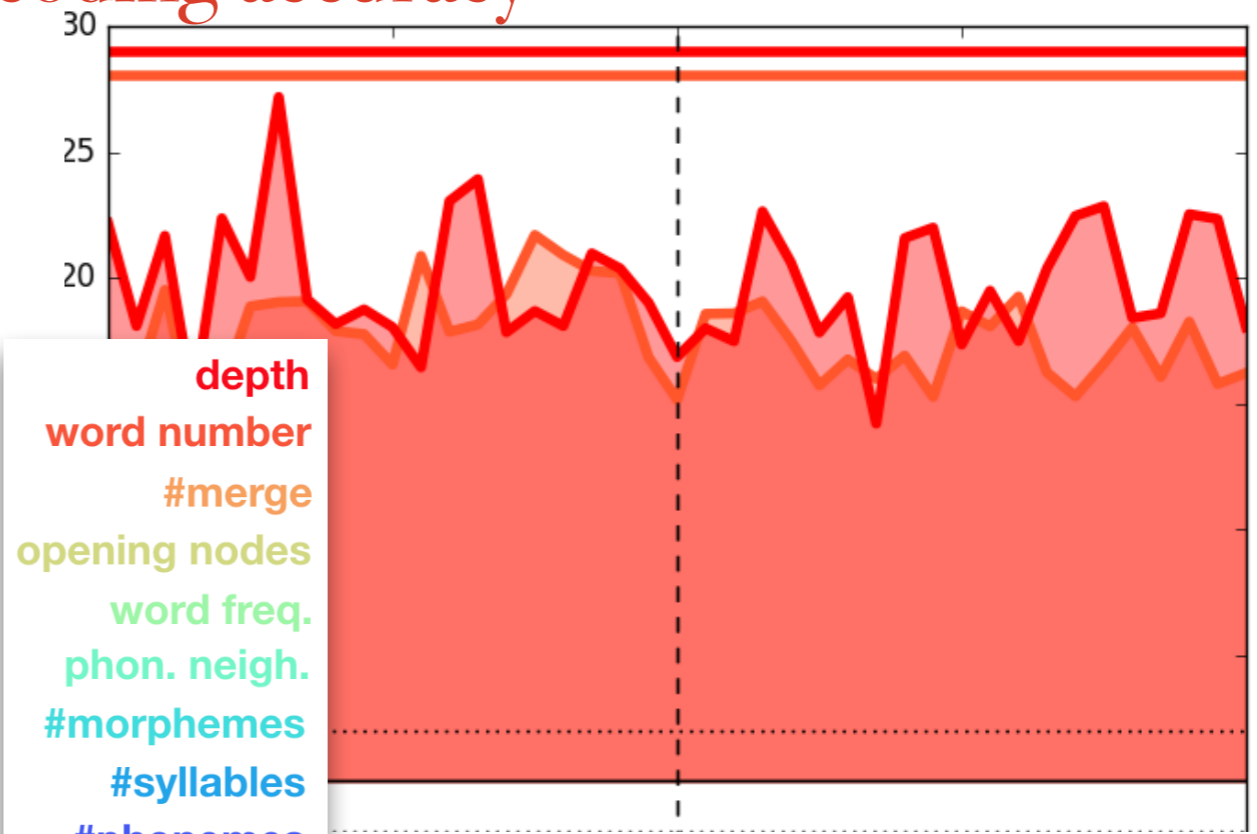
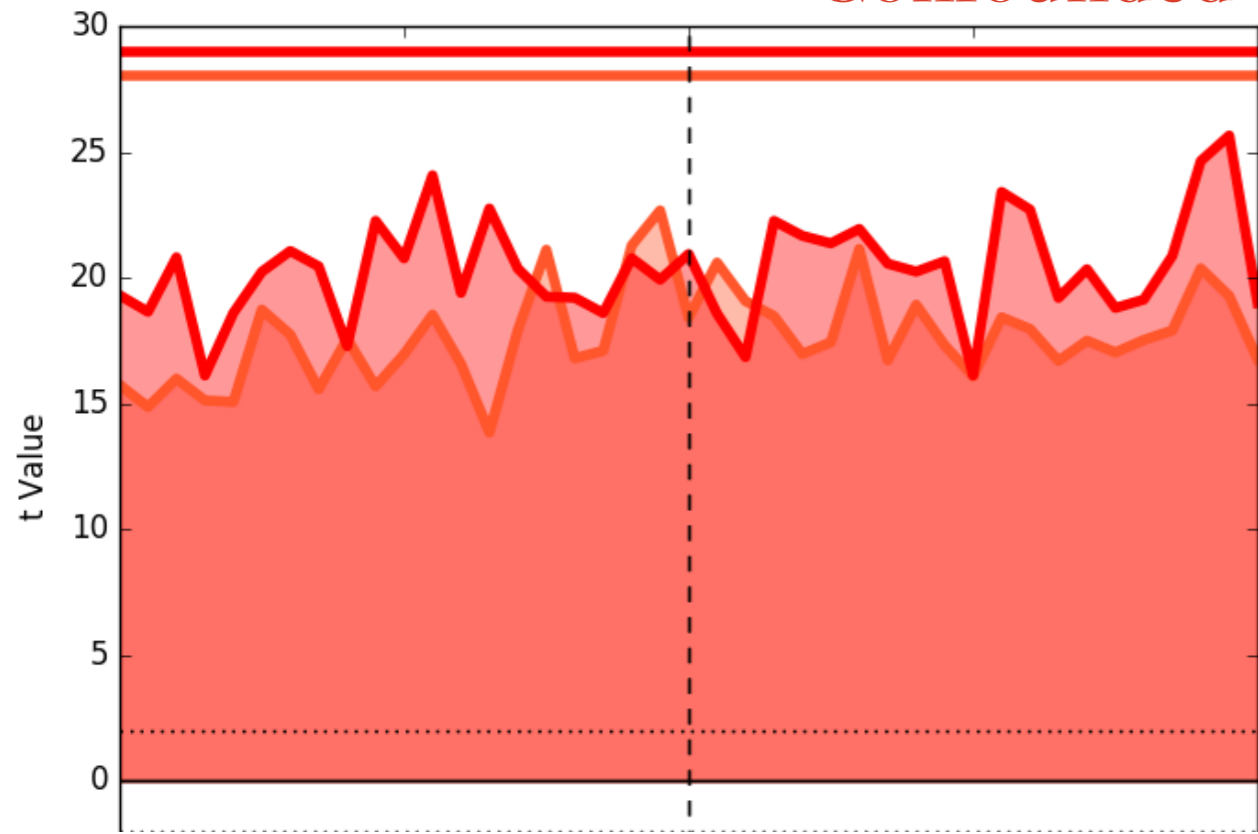
-0.0526 -0.0375 -0.0225 -0.00751 0.00751 0.0225 0.0375 0.0526

depth



-0.0728 -0.0520 -0.0312 -0.0104 0.0104 0.0312 0.0520 0.0728

Confounded decoding accuracy



depth
word number
#merge
opening nodes
word freq.
phon. neigh.
#morphemes
#syllables
#phonemes

