



Any questions?

- Practicalities?
- Any open issues from yesterday?

Lecture 5: *audiovisual prosody*

Bosker, H. R., & Peeters, D. (2021). Beat gestures influence which speech sounds you hear. *Proceedings of the Royal Society B*, 288(1943), 1–9.
doi:[10.1098/rspb.2020.2419](https://doi.org/10.1098/rspb.2020.2419).

Hans Rutger Bosker

Speech Perception in Audiovisual Communication [SPEAC] lab

Donders Institute, Radboud University, Nijmegen, The Netherlands

<https://hrbosker.github.io>

hansrutger.bosker@donders.ru.nl





Listening with your eyes!

- We can ‘see’ speech from articulatory movements alone
 - Lip-reading / speech-reading
- We use visual articulatory signals to inform speech perception
 - Audiovisual enhancement of speech-in-noise
 - Classic McGurk effect

Listening with your eyes!

- Classic McGurk effect



Listening with your eyes!

- Classic McGurk effect



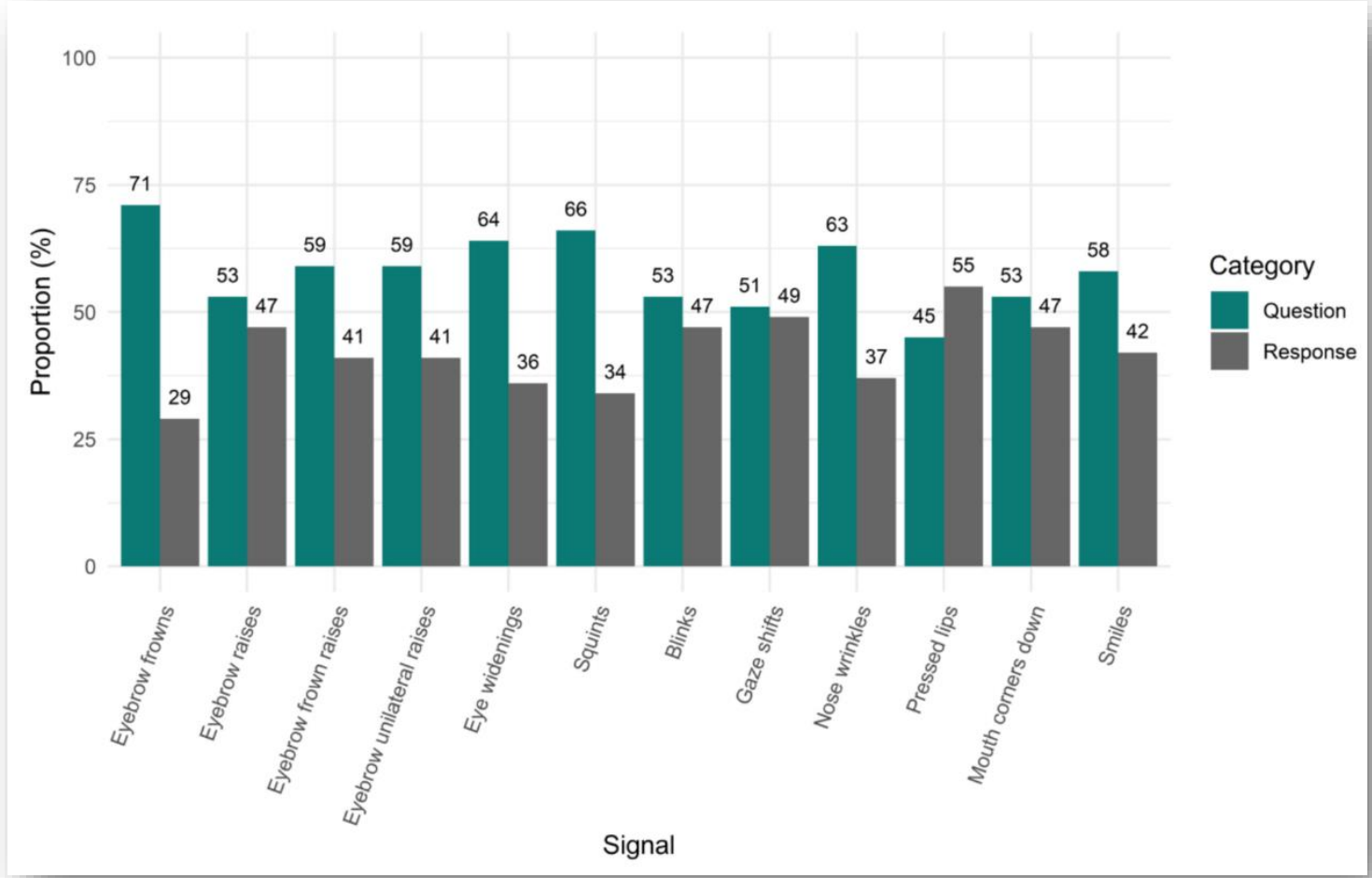
Listening with your eyes!

- Classic McGurk effect



Can we *see* prosody?

McGurk & MacDonald, 1976





Can we *lip-read* prosody?

- Speech rate estimates are equally accurate in V-only vs. A-only (Green, 1987)
- ??? intensity ???
- Phonological contrasts that primarily rely on f_0 are much harder to lip-read:
 - Lexical tone
 - Video-only accuracy is at chance (Burnham et al., 2001)
 - ...but V-only training can help (Chen & Massaro, 2008)
 - Lexical stress
 - Video-only accuracy above chance
 - ...and Dutch (less segmental r)



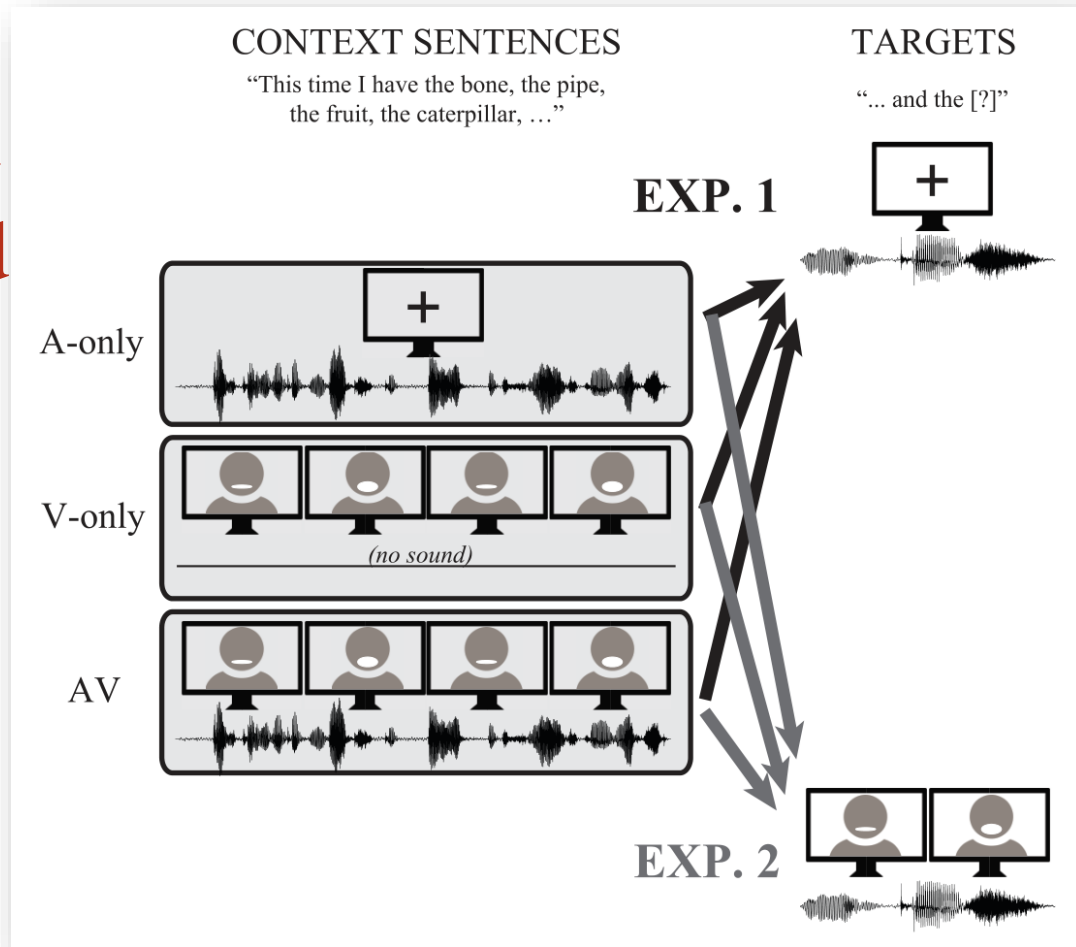
Do we *use* visual articulatory cues to prosody in audiovisual speech perception?

- Speech rate
 - A video of a fast /pi/ + /bi-pi/ auditory continuum: +/p/ responses
Green & Miller, 1985
 - A video of a fast /ti/ + /bi-pi/ auditory continuum: +/p/ responses
Brancazio & Miller, 2005
 - Rate normalization by visual speech rate?



Do we
in aud

prosody





Do we
in aud

prosody

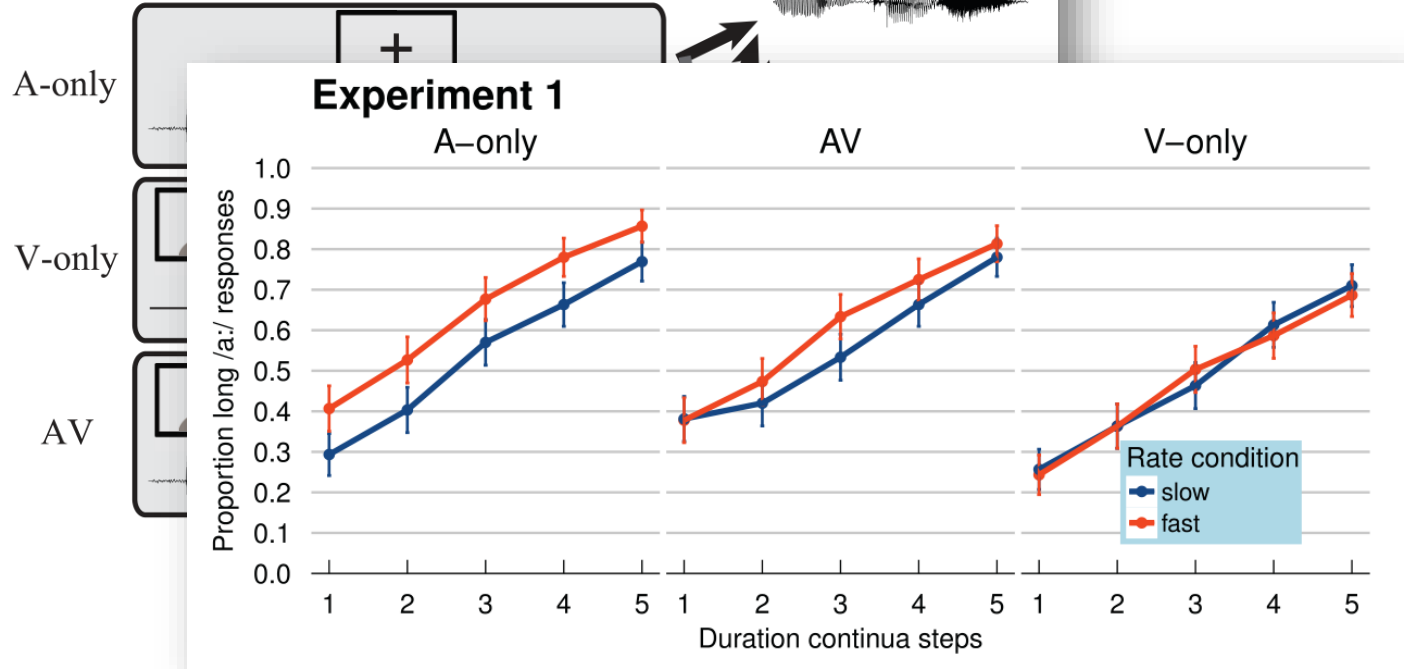
CONTEXT SENTENCES

“This time I have the bone, the pipe,
the fruit, the caterpillar, ...”

TARGETS

“... and the [?]”

EXP. 1





Prosody

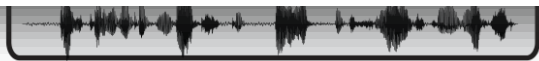
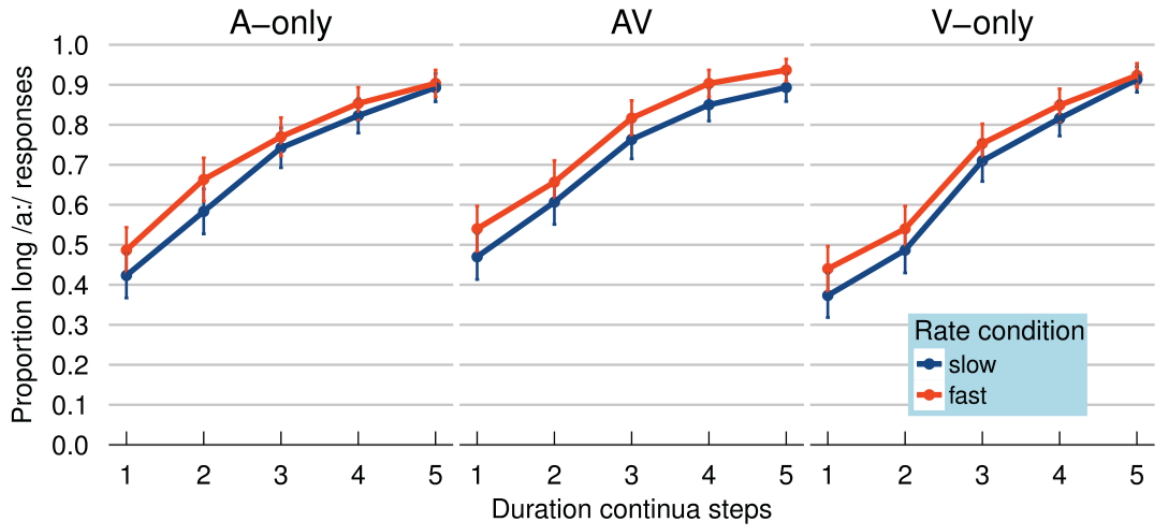
CONTEXT SENTENCES

"This time I have the bone, the pipe,

TARGETS

"... and the [?]"

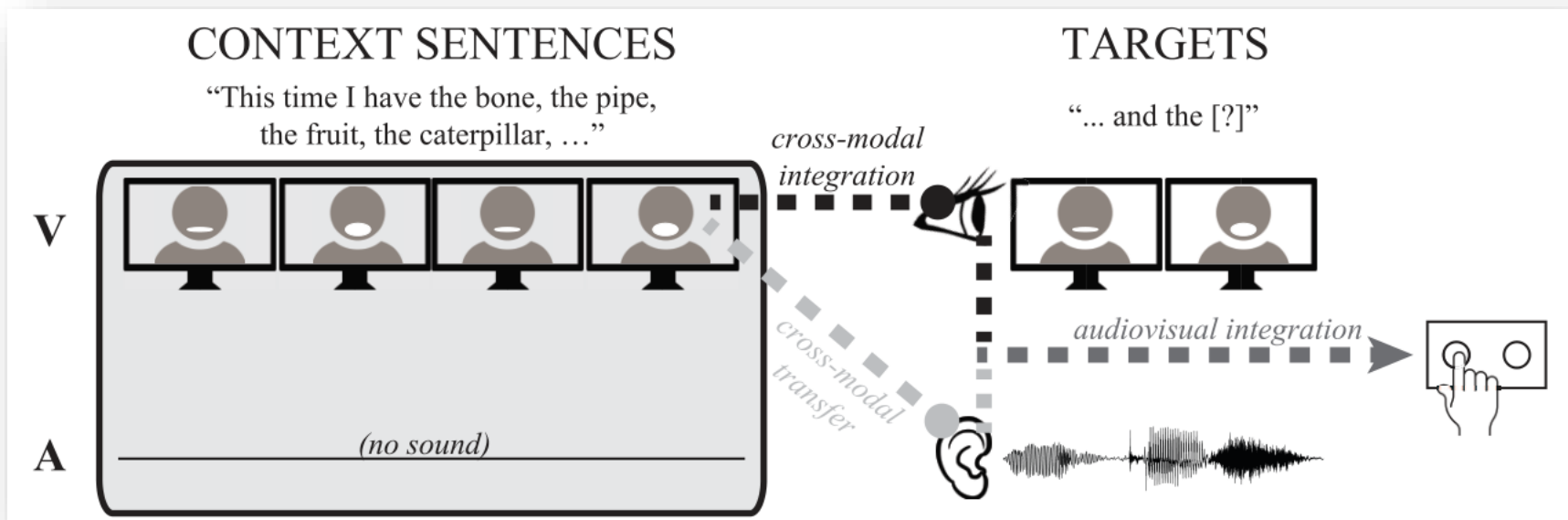
Experiment 2

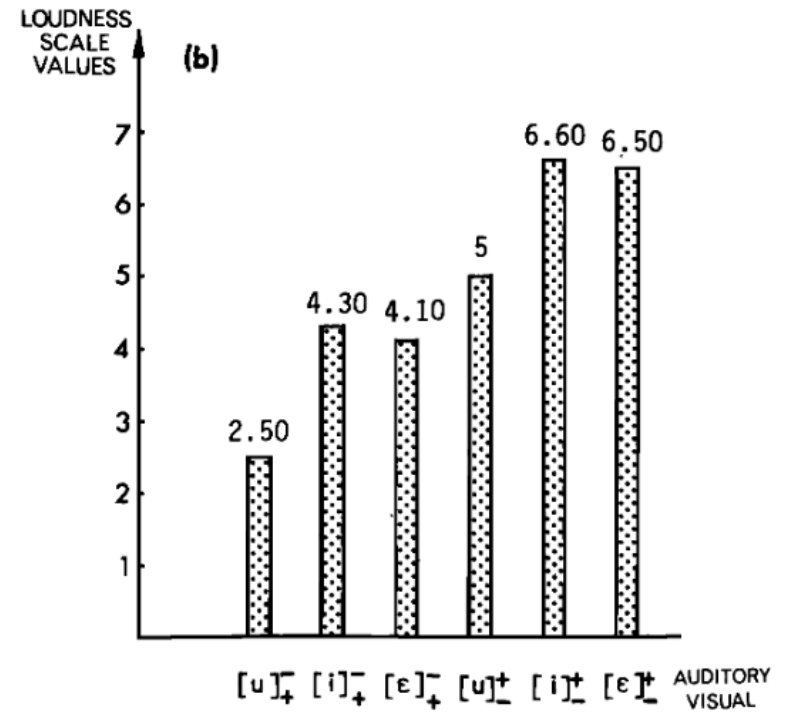
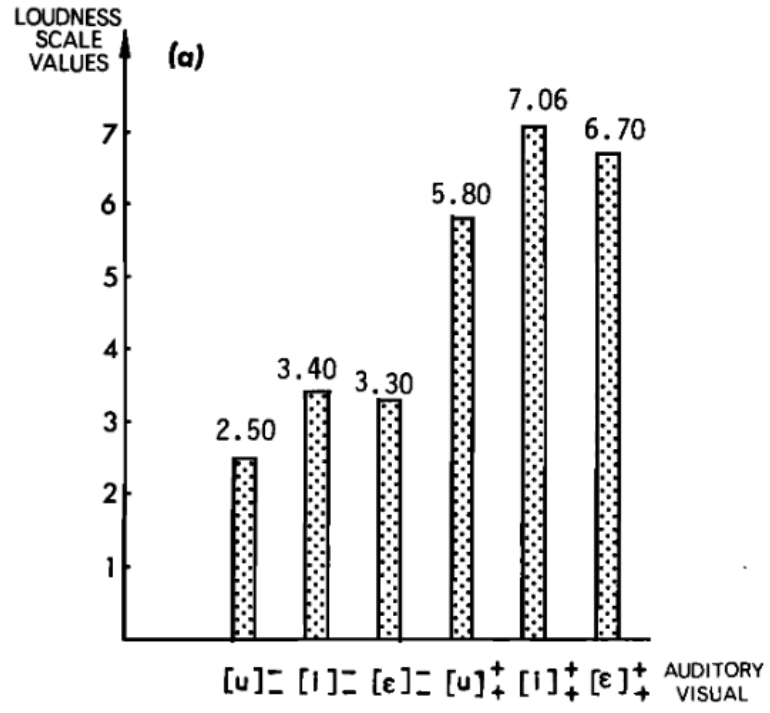


EXP. 2



Do we *use* visual articulatory cues to prosody in audiovisual speech perception?

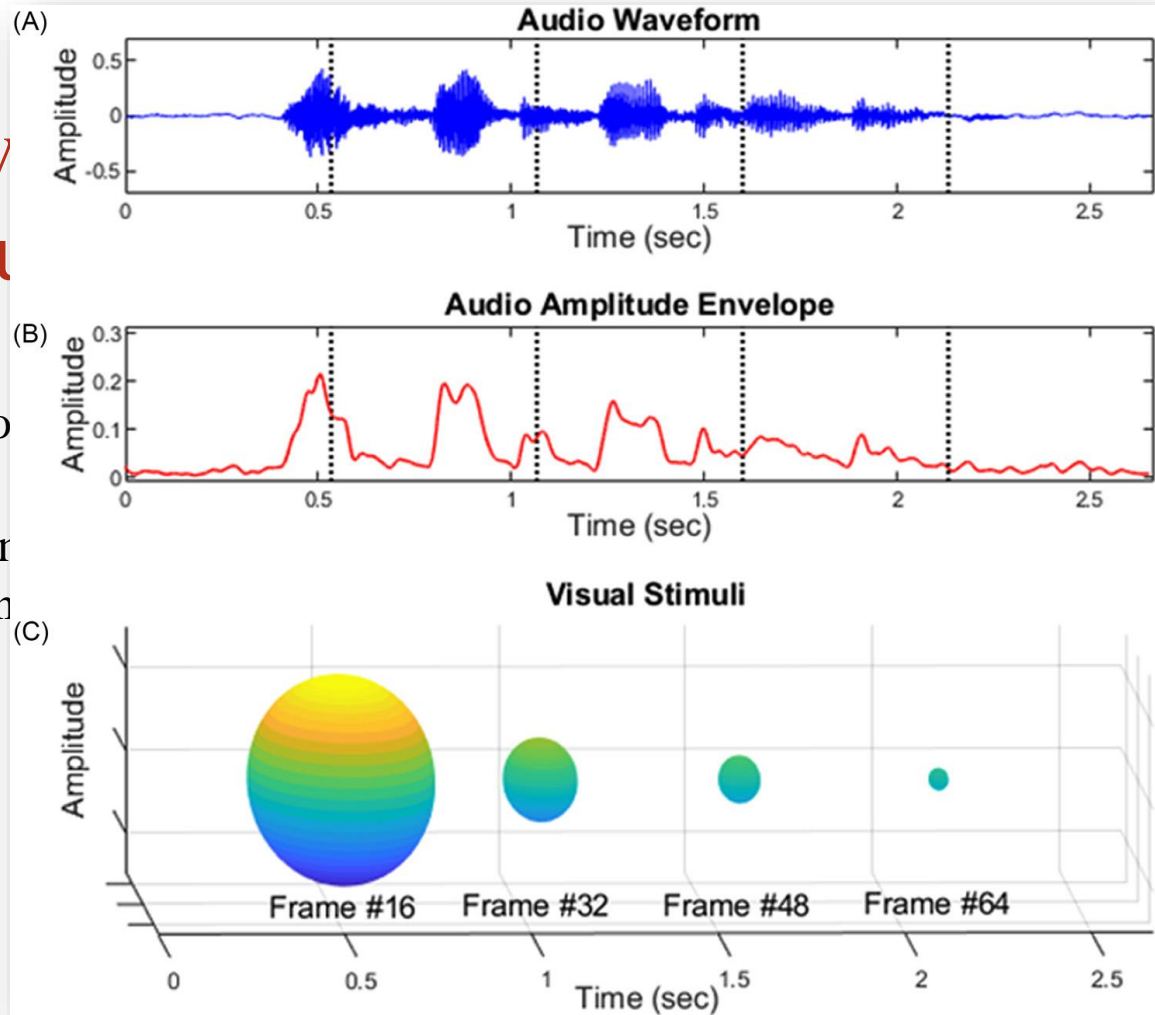






Do we use volume in audiovisual speech processing?

- Intensity
 - Seeing volume
 - Seeing a range of speech



, 1991
(m)
g'?) ...



Do we *use* visual articulatory cues to prosody in audiovisual speech perception?

- Intensity

- Seeing vocal effort boosts perceived loudness

Glave & Rietveld, 1979; Rosenblum & Fowler, 1991

- Seeing a modulating circle in sync with the amplitude envelope (rhythm) of speech boosts its intelligibility... (...link to neural ‘speech tracking’?)

Yuan et al., 2020

- ...or does it?

Strand et al., 2020



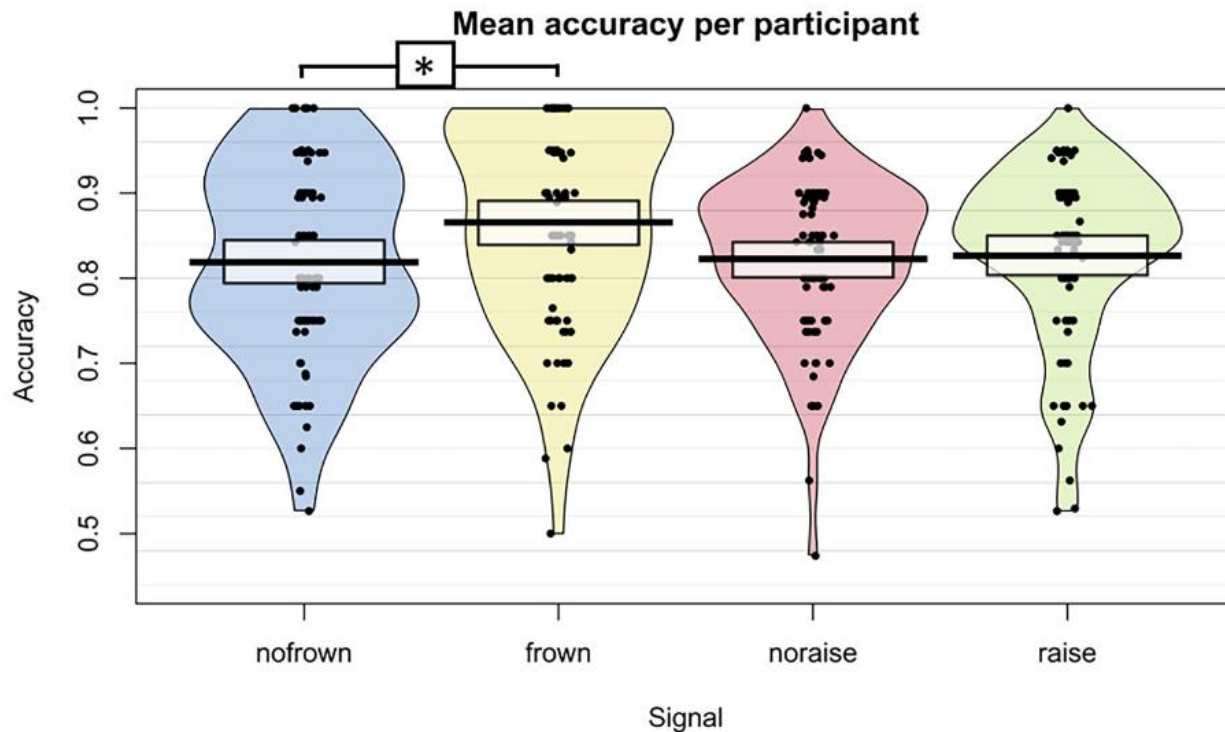
Do we *use* visual articulatory cues to prosody in audiovisual speech perception?

- f_0
 - Sentence accent
 - Which word in this sentence is the most prominent?
 - Slow down in RTs when V is incongruent with A

Swerts & Kraemer, 2008



prosody



verts & Kraemer, 2008

Nota et al., 2022

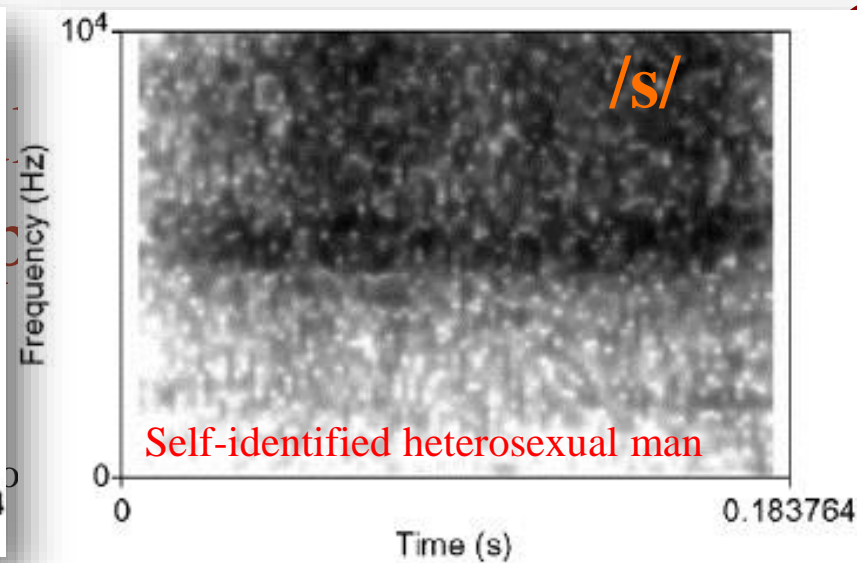
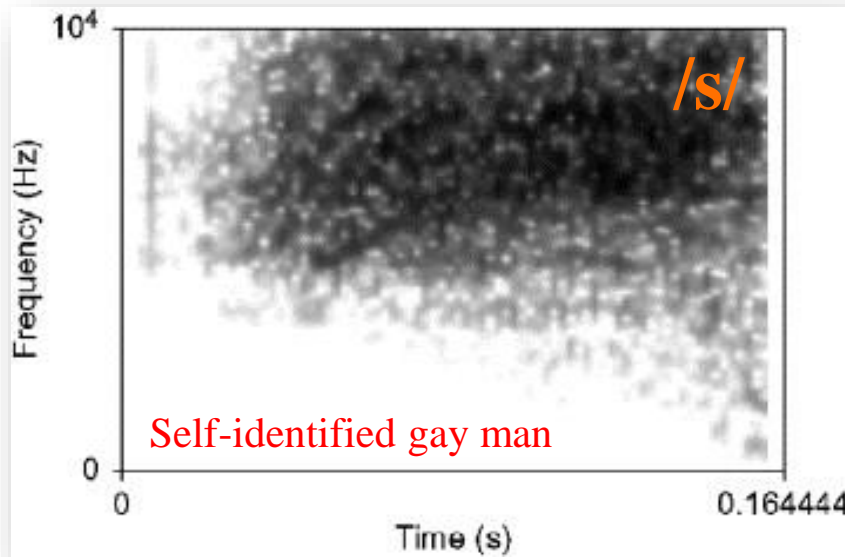


Do we *use* visual articulatory cues to prosody in audiovisual speech perception?

- f_0

Other visual cues to f_0 ?

- f_0 spectral energy in female talkers
Jongman et al., 2000
- Same fricative on /ʃ-s/ continuum more /ʃ/-like if female speech
Mann & Repp, 1980
- Just **seeing** a female face also leads to more /ʃ/-responses
Strand et al., 1996
- Even **imagining** listening to a female talker induces normalization
~ prior knowledge?
Johnson et al., 1999



- Common societal impression: higher-frequency spectral energy in gay talkers + lower-frequency spectral energy in lesbian talkers

Munson et al., 2006, *JPhon*

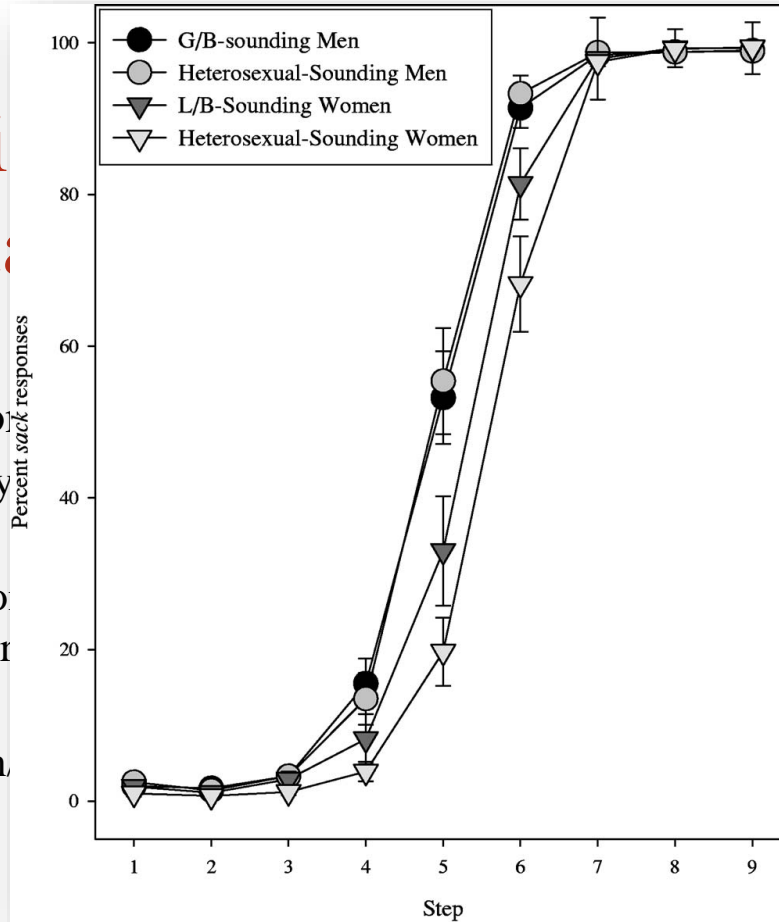
- Lesbian/bisexual-sounding talkers elicit more /s/-responses

Munson et al., 2006, *JASA*



Do we use visual cues in audiovisual speech processing?

- f_0
- Spectral noise
- Visually
- Common + lower
- Lesbian



o prosody

users

Winn et al., 2013

ral energy in gay talkers

Munson et al., 2006, *JPhon*

ponses

Munson et al., 2006, *JASA*



Do we *use* visual articulatory cues to prosody in audiovisual speech perception?

- f_0
 - Lexical tone
 - AV accuracy = A accuracy in Mandarin
Burnham et al., 2001
 - A ye1 + V ye2 = A ye1 (incongruent AV stimuli)
Hannah et al., 2017
 - AV advantage only surfaces when is degraded (noise/vocoded speech)
Mixdorff et al., 2005; Burnham et al., 2015
 - Perhaps no surprise because hardly any visual cues to lexical tone (V-only)
Burnham et al., 2001



Do we *use* visual articulatory cues to prosody in audiovisual speech perception?

- f_0
 - Lexical stress
 - Visual cues to stress are detectable in production...
Scarborough et al., 2009
 - ...and people can ‘lip-read’ stress from V-only videos...
Scarborough et al., 2009; Jesse & McQueen, 2014
 - ...so surely people should use these cues in AV speech perception, right?



Prosody on the face

Prosody on the face

CAnon
video

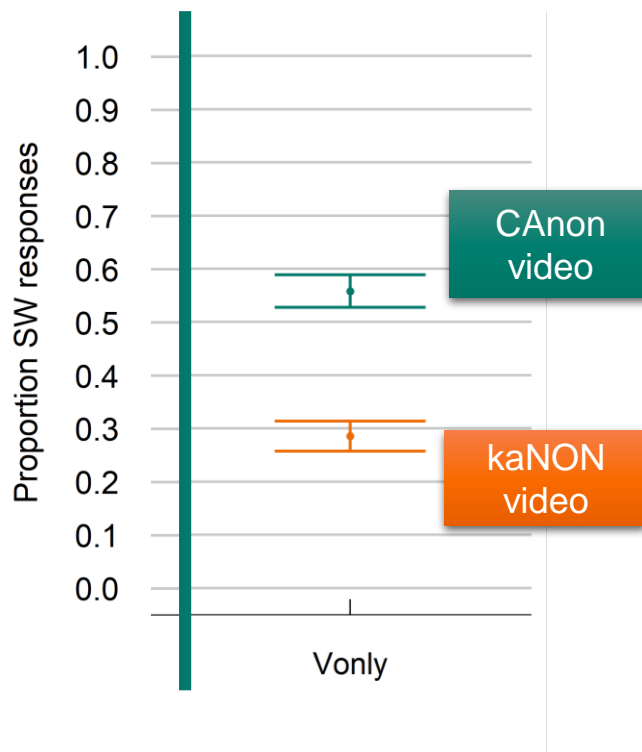


kaNON
video

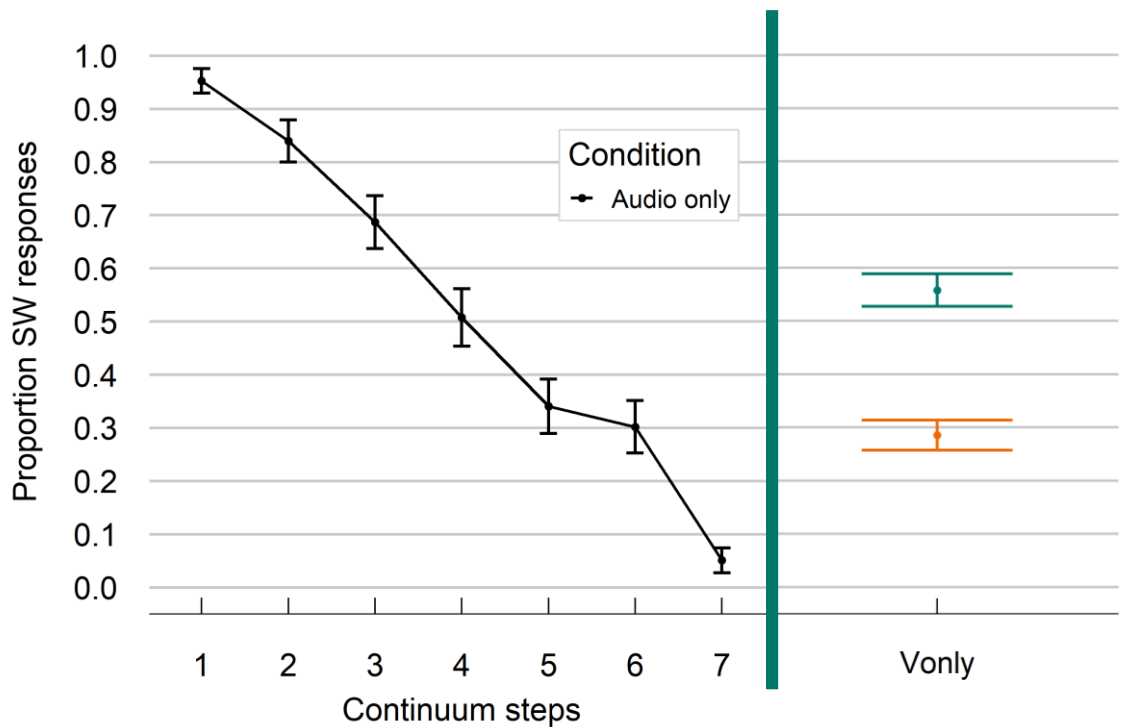




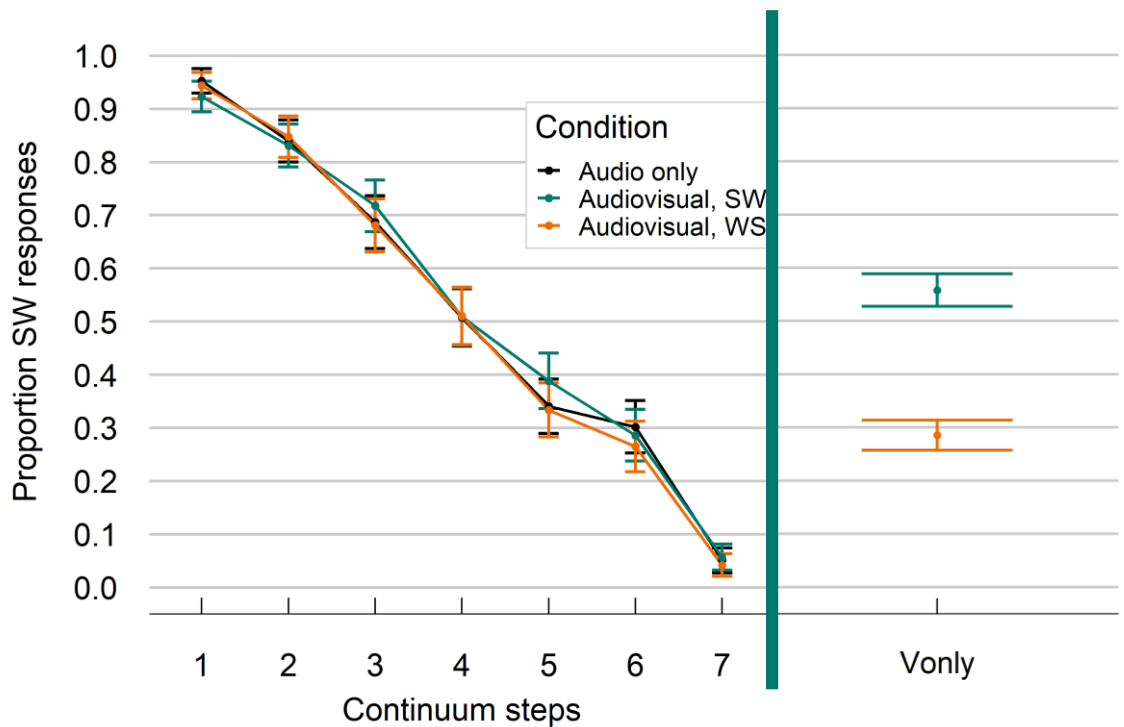
Prosody on the face



Prosody on the face



Prosody on the face





Prosody on the face

Proportion SW responses



Audio step



Do we *use* visual articulatory cues to prosody in audiovisual speech perception?

- f_0
 - Lexical stress
 - Visual articulatory cues on the face hardly contribute to AV speech perception
Bujok et al., 2024, L&S
 - This is in stark contrast to how people use visual segmental cues (McGurk)
 - Low reliability of visual prosody cues leads people to down-weight these cues in quiet...
 - ...but up-weight their contribution when listening is hard.

Prosody in the hands

- We move to the rhythm of our speech



“... and when and how
and on – what – ba sis...”

Prosody in the hands

- We move to the rhythm of our speech
- Relatively simple up-and-down hand movements make up over 90% of our everyday gestures

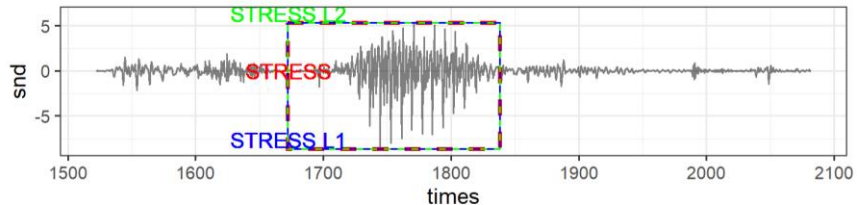
Shattuck-Hufnagel & Ren, 2018

- Closely coupled to speech prosody, typically falling on stressed syllables
- Gesture's apex time-locked to pitch peak in stressed syllable.

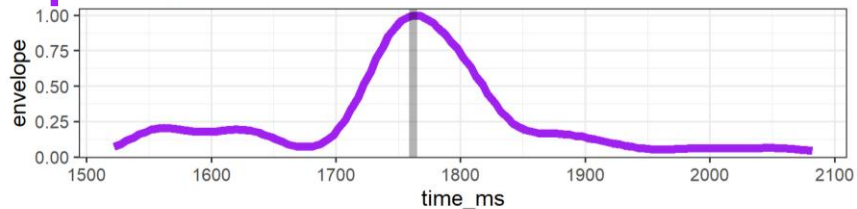


Prosody in the hands

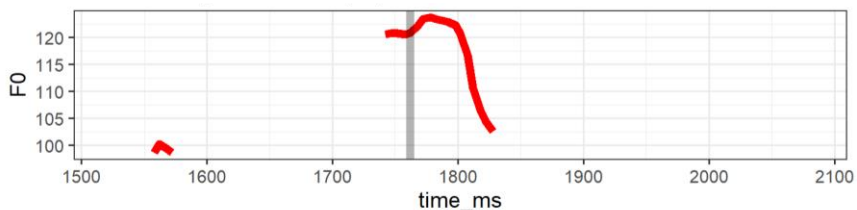
Spanish: “his.TÓ.ri.co”



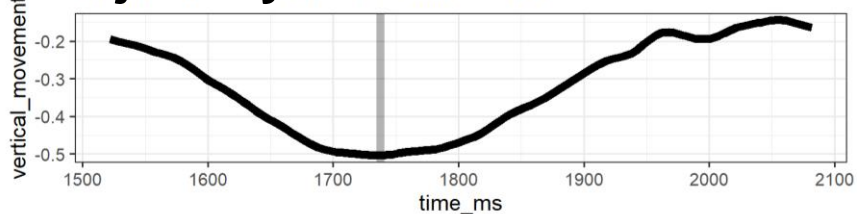
amplitude



f_0



hand trajectory





Prosody in the hands

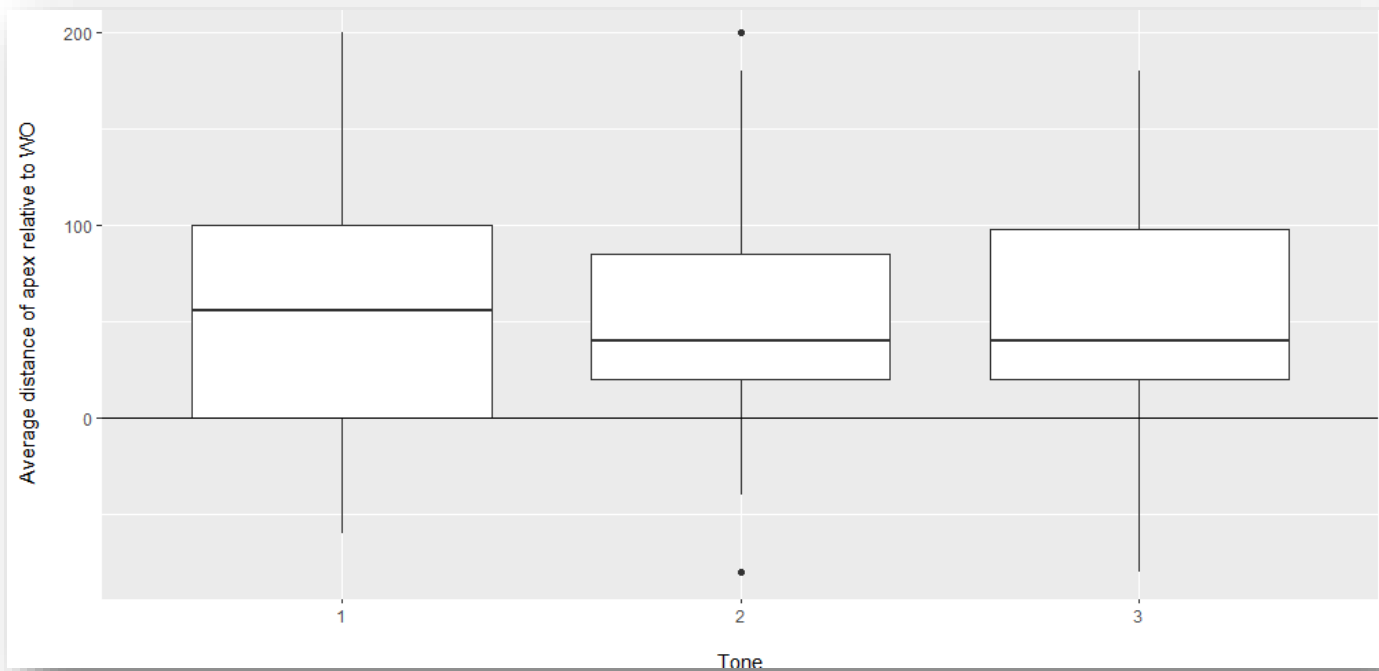
- Why do we beat on the stressed syllable?
 1. Universal biomechanics:
moving your limbs raises your voice (amplitude, f_0)
 2. Language-specific multimodal phonology:
no evidence for coupling to f_0 contour in tone languages

Pouw et al., 2020

Rohrer et al., 2024



Prosody in the hands



Pouw et al., 2020

Rohrer et al., 2024



Prosody in the hands

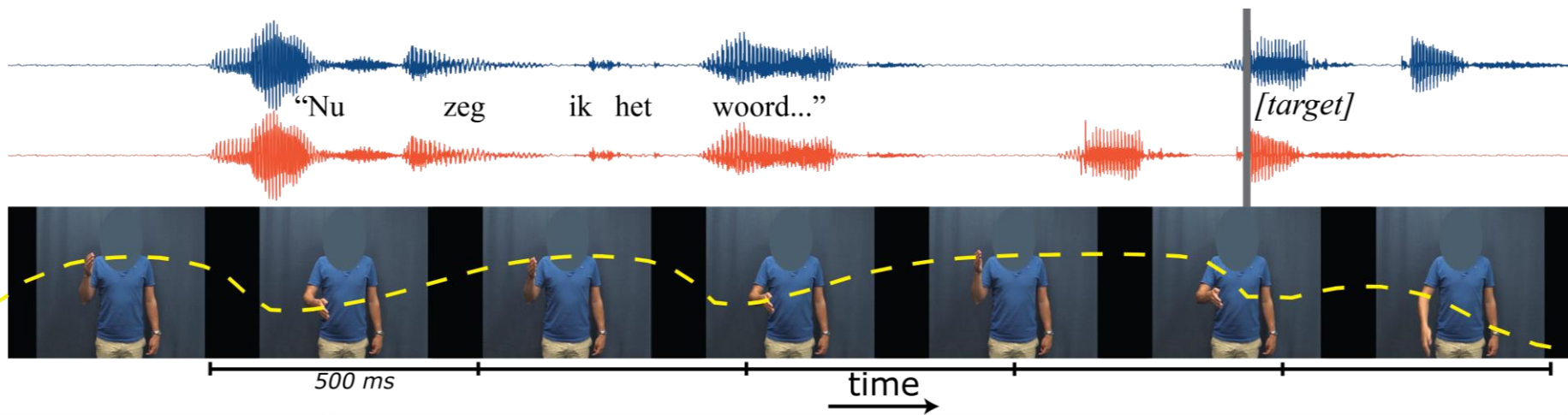
- If talkers tend to gesture on stressed syllables...
- ...can listeners use this visual temporal cue to perceive lexical stress?

Prosody in the hands

Nu zeg ik het woord... [PLAto/plaTEAU]



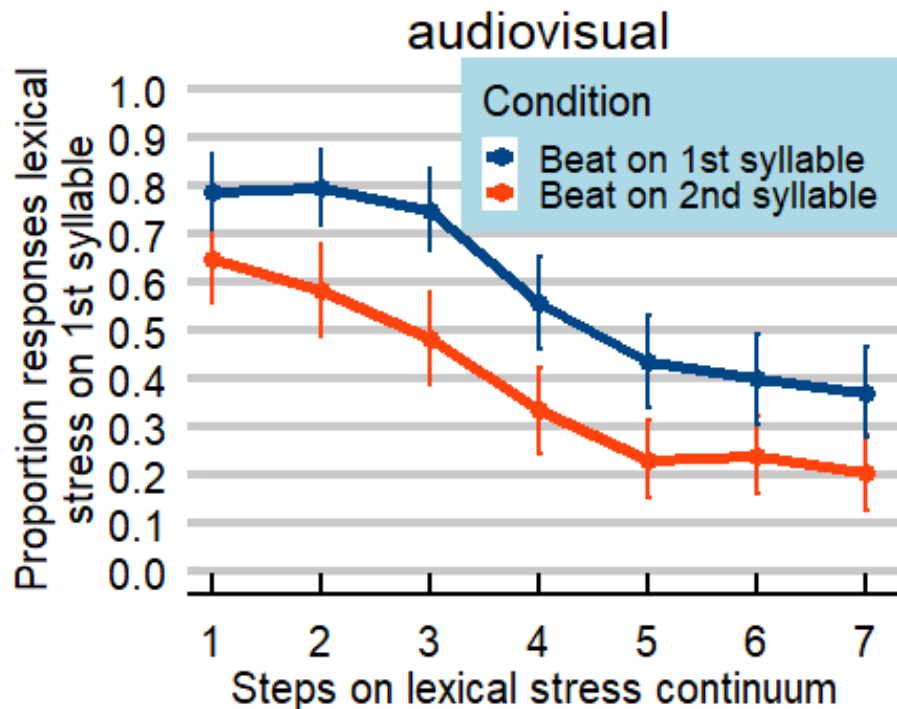
Prosody in the hands





Prosody in the hands

- ‘manual McGurk effect’



Prosody in the hands

nonwords

real words

shadowing

vowel length

Condition



Beat on 1st syllable

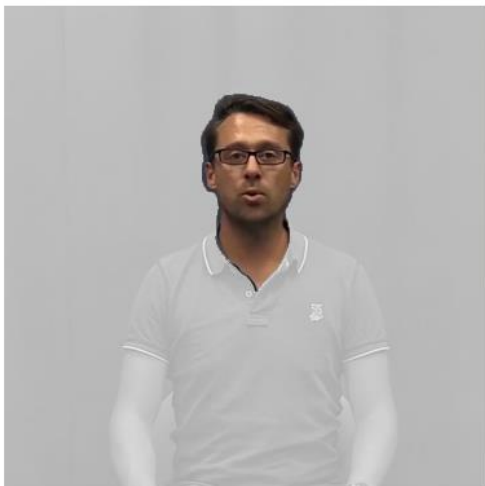
Beat on 2nd syllable



A beat gesture on the first syllable of /b?x.pif/ will bias towards: short /bax.pif/ or long /ba:x.pif/

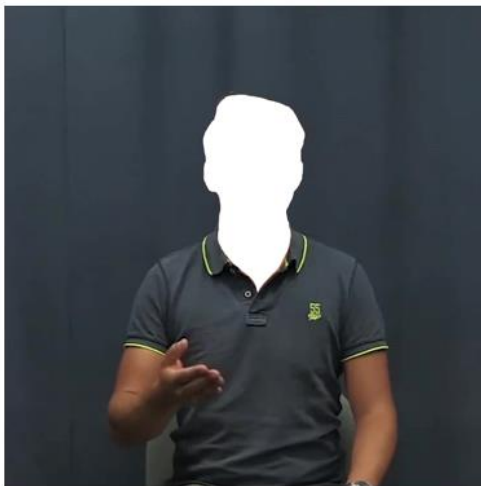
Prosody in the hands

*head from no-gesture
condition*



+

gesturing body



=

*new fully-crossed
stimuli*

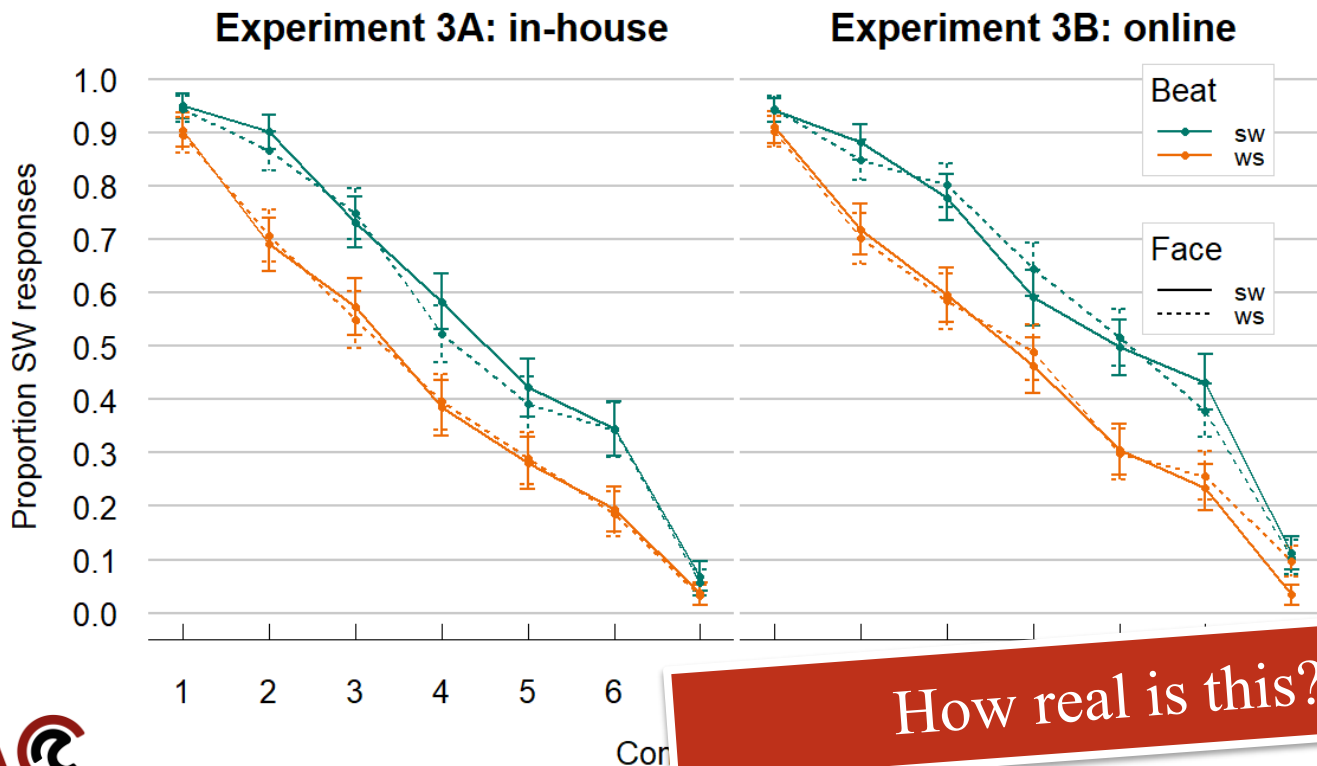




Prosody in the hands



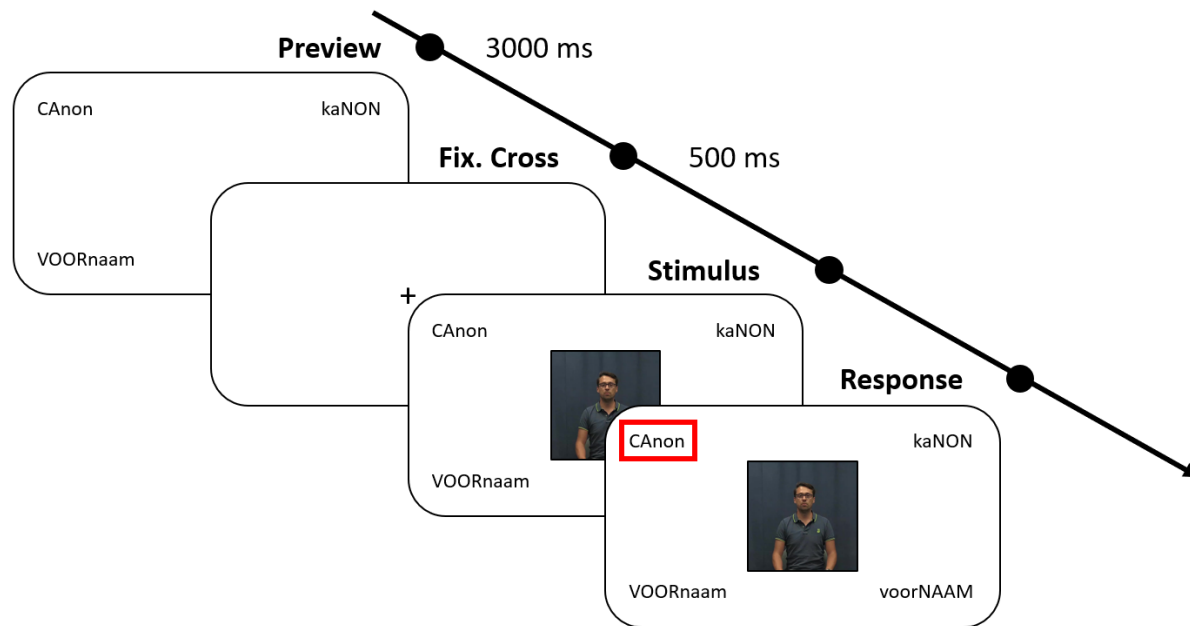
Prosody in the hands



How real is this?

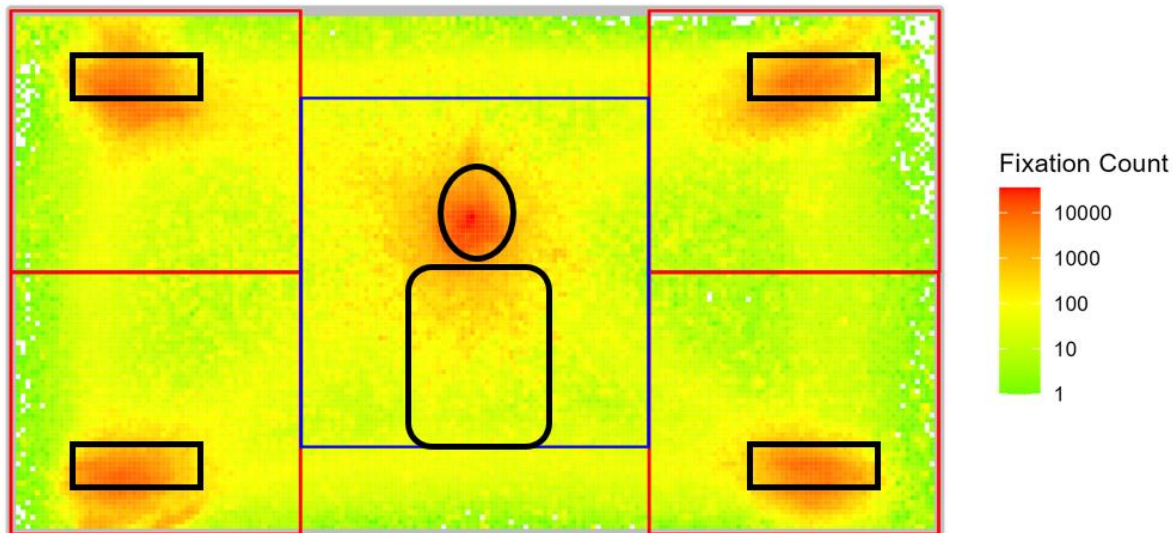
Prosody in the hands: immediate impact

- Eye-tracking: visual world paradigm



Prosody in the hands: immediate impact

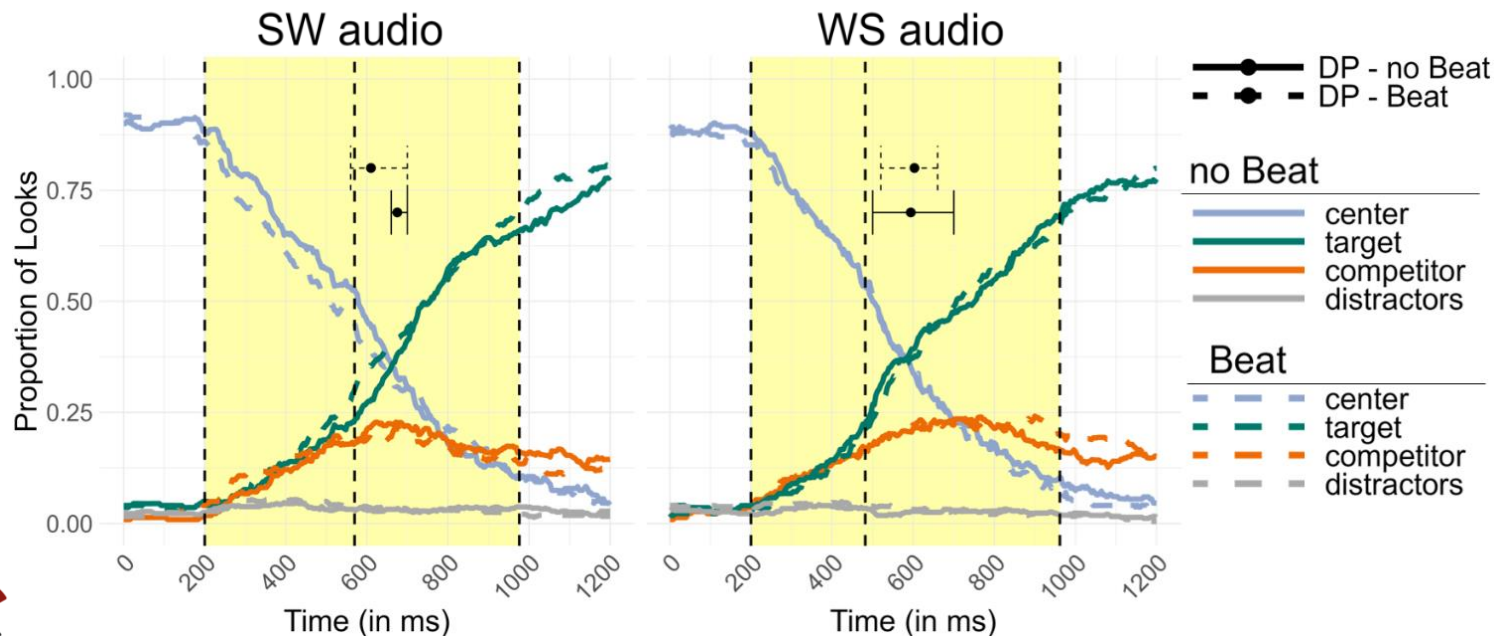
- Hardly any fixations on the gesture





Prosody in the hands: immediate impact

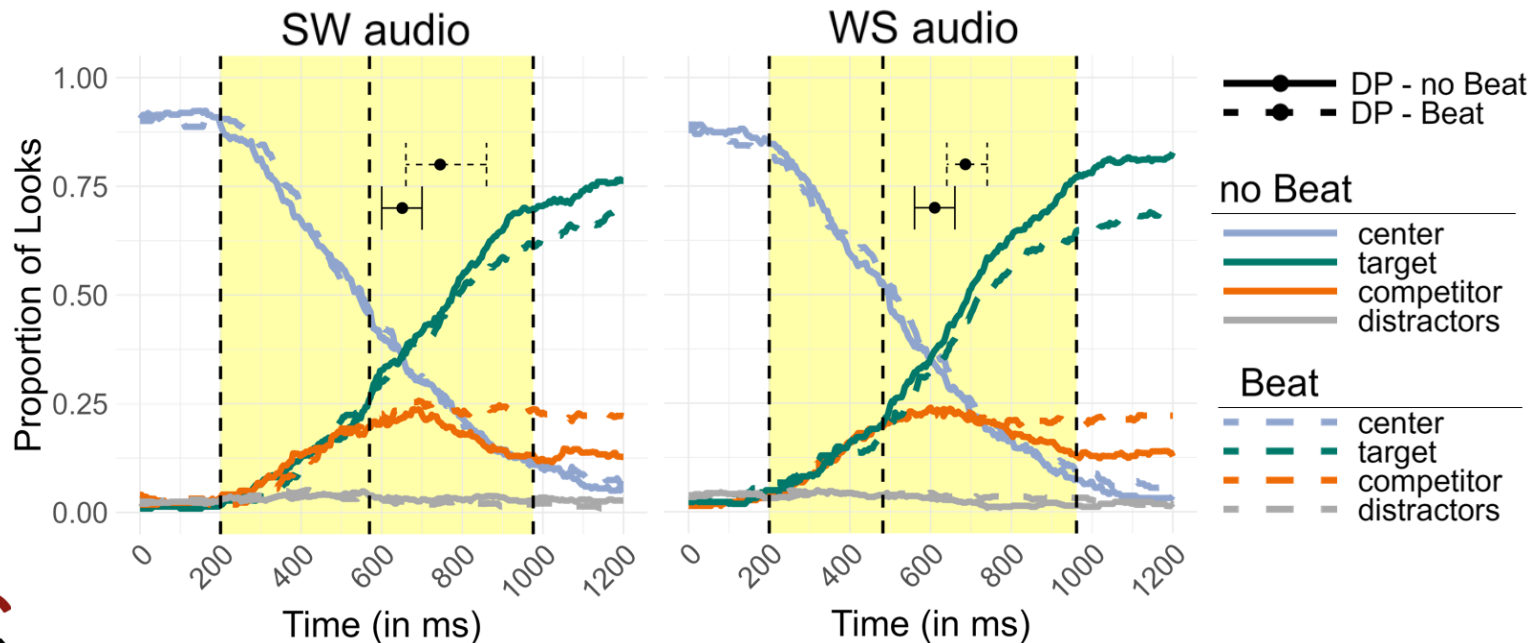
- When audio clear, no (reliable) beneficial effect of stress-congruent beat.





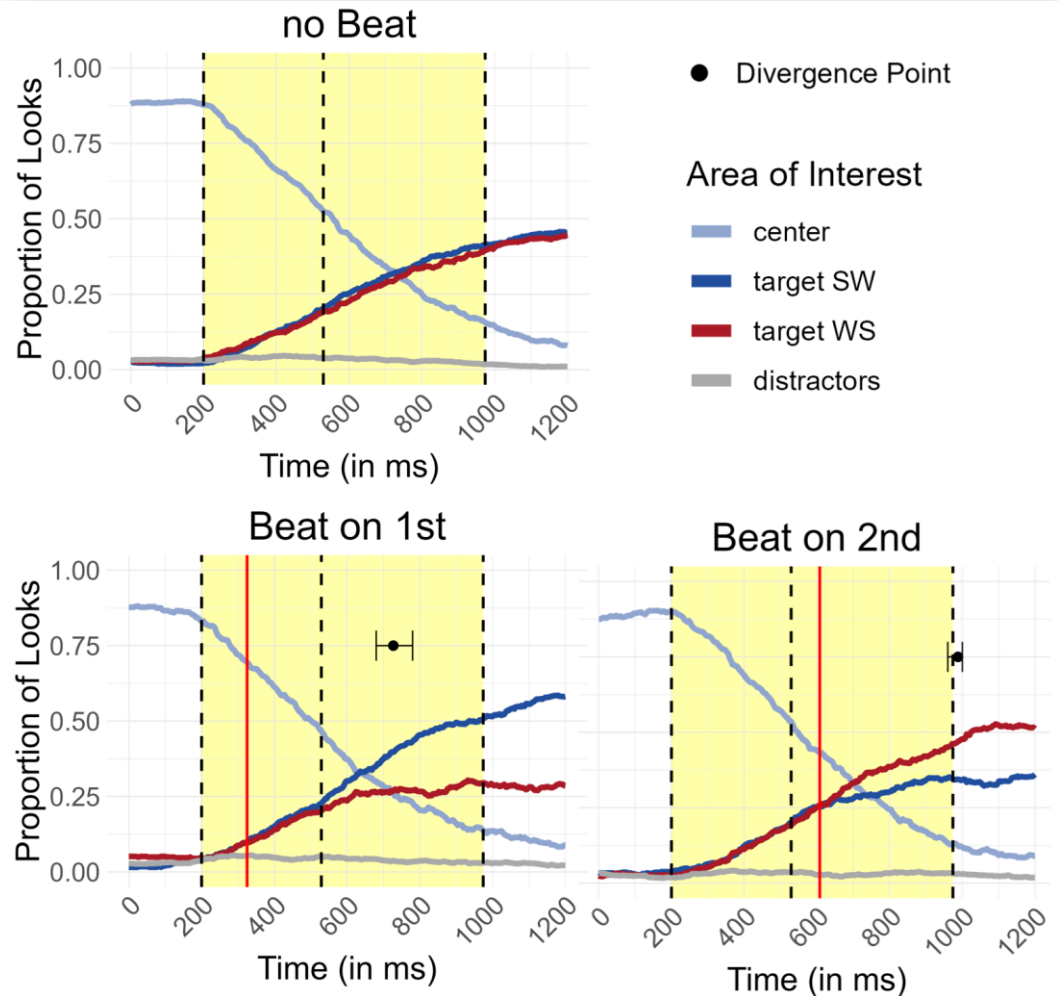
Prosody in the hands: immediate impact

- When audio clear, reliable interfering effect of stress-**incongruent** beat.



Prosody in the hard

- When audio ambiguous, beat-on-1 biases to SW, beat-on-2 biases to WS.
- Divergence point lies before word offset.
- Manual McGurk effect is an early perceptual bias, not only a decision-making effect.

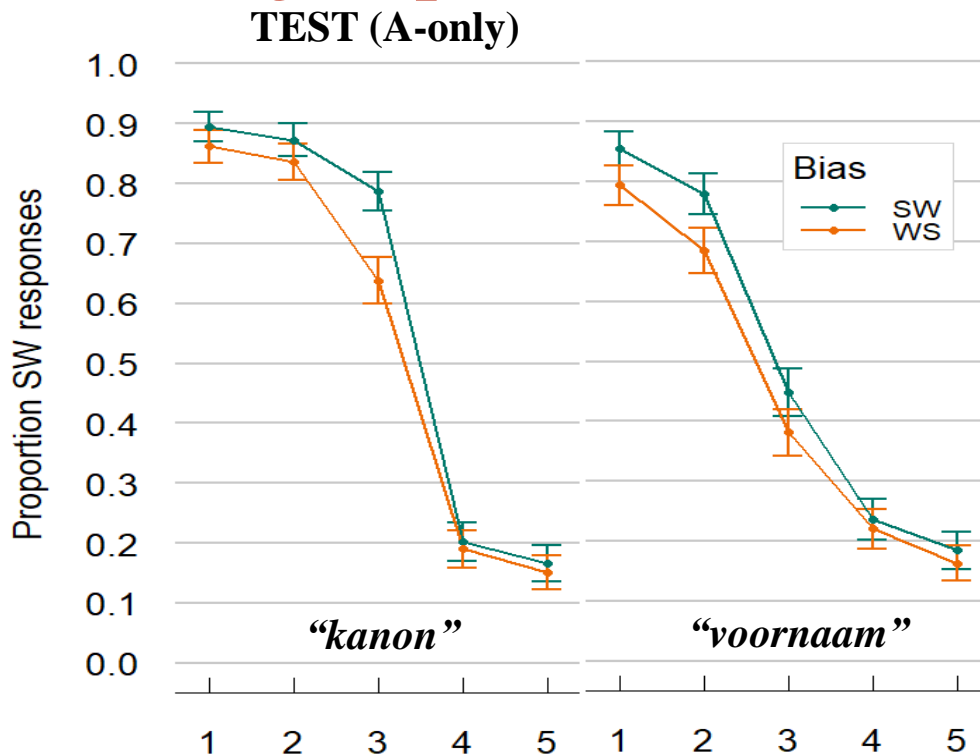


Prosody in the hands: lasting impact

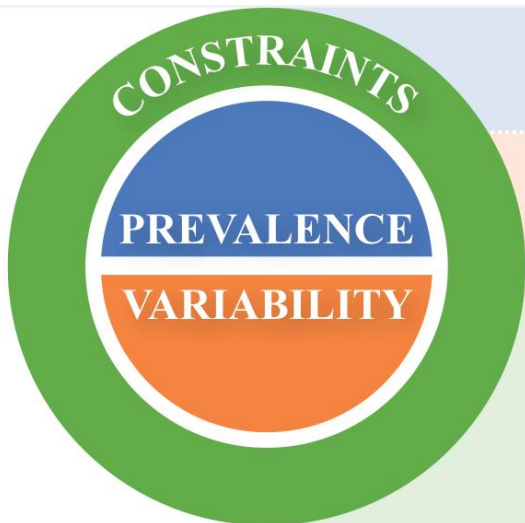
- Audiovisual recalibration of lexical stress

EXPOSURE (AV)

- Group 1: /ka.nɔn/?
+ beat on syllable 1
- Group 2: /ka.nɔn/?
+ beat on syllable 2



ERC StG HearingHands



WP1
(PD)

Gesture-speech synchrony in production & perception in free-stress, fixed stress, and tonal languages

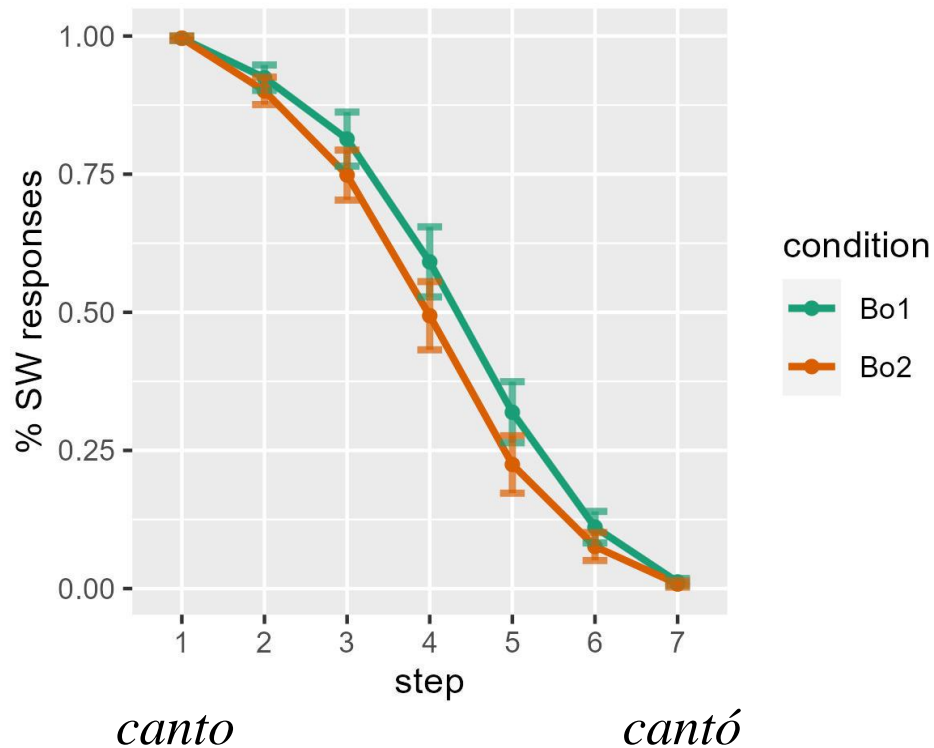
WP2
(PD)

Neurobiology of gesture-speech synchrony integration and individual differences in the neurotypical and autistic population

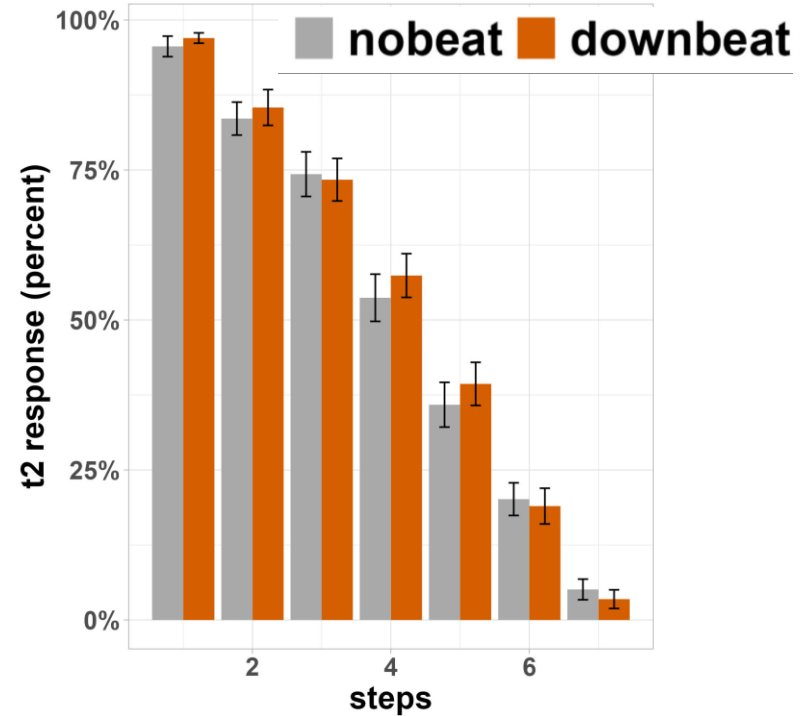
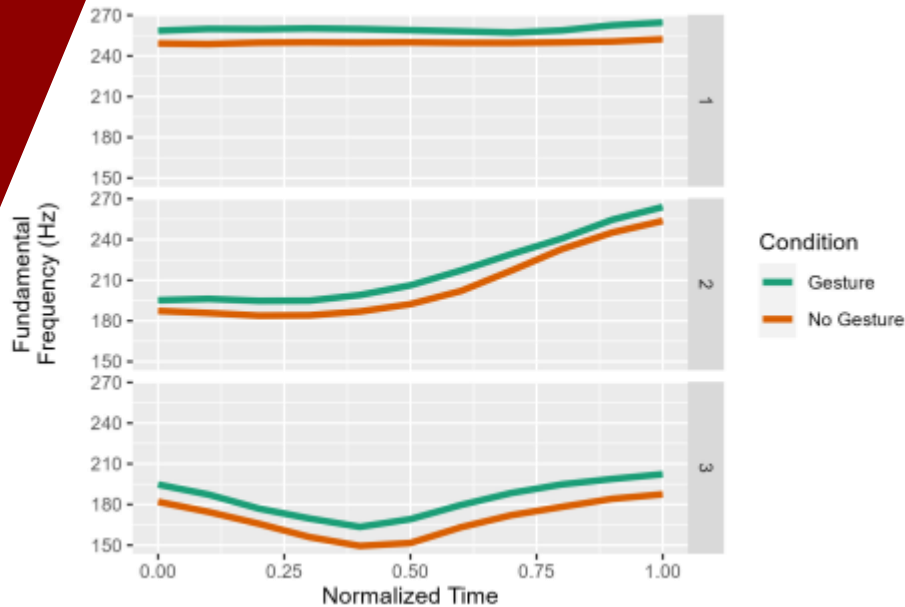
WP3
(PHD)

Communicative, situational, and cognitive constraints on the effects of gesture-speech synchrony on spoken word recognition in eye-tracking and VR

ONGOING: “Efecto McGurk manual”

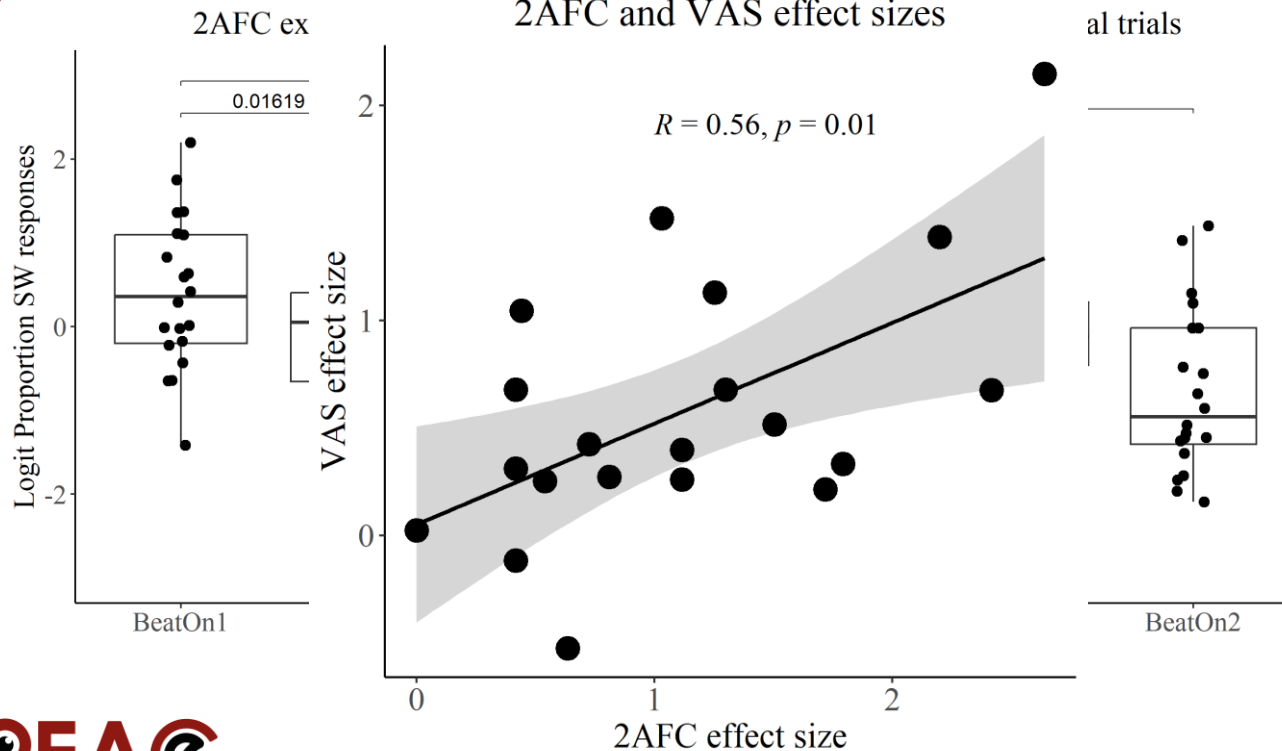


ONGOING: Gesture-speech perception in Mandarin



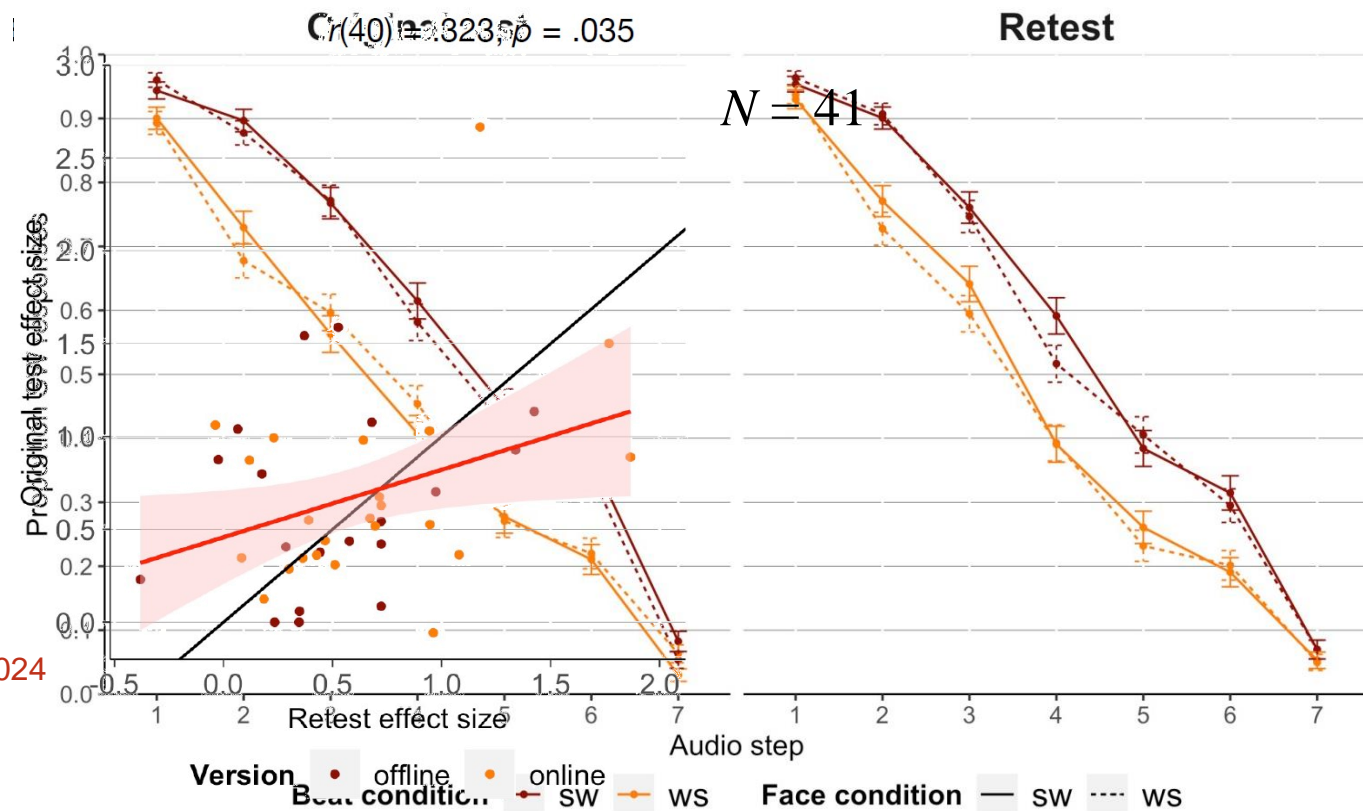
ONGOING:

Minitest for inclusion in test batteries



ONGOING:

Test-retest reliability after >1.5 years

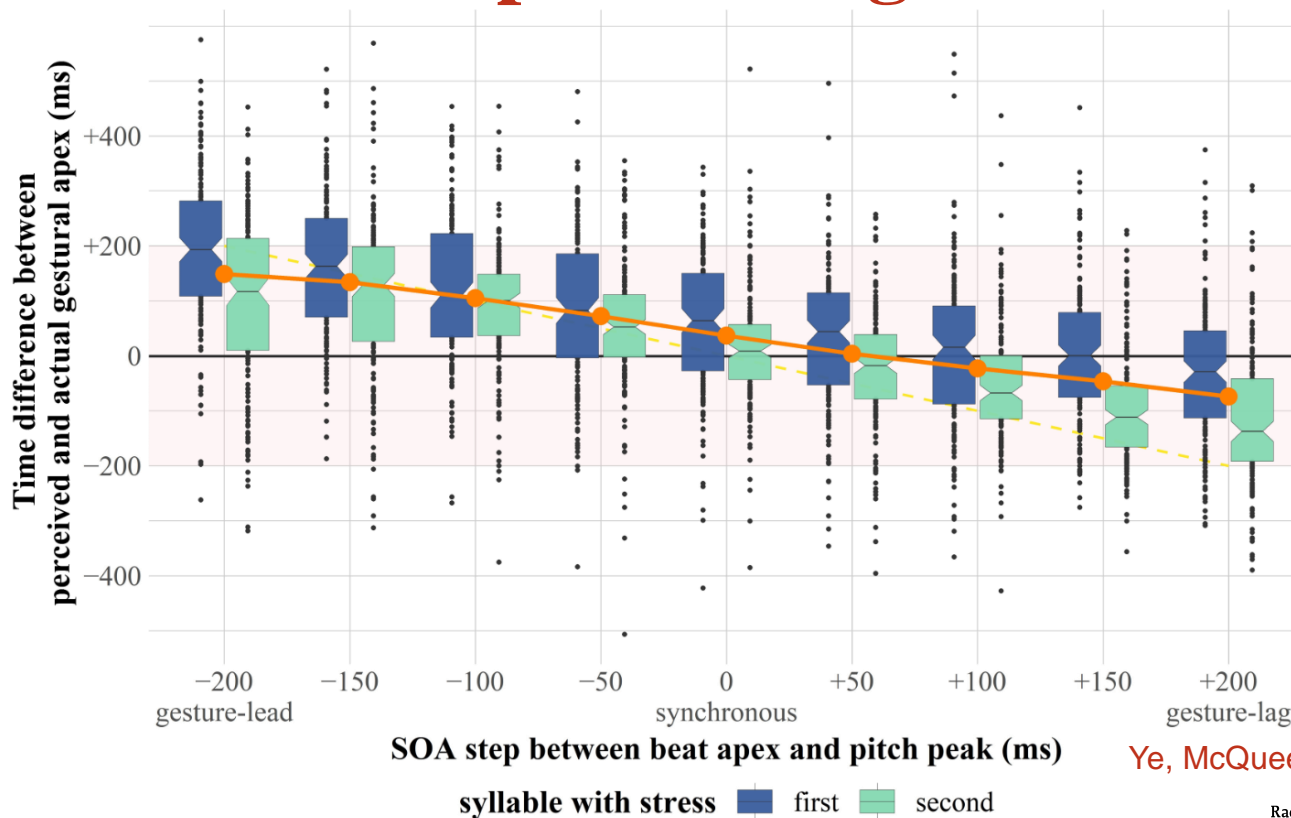


Cos, Bujok, & Bosker, 2024



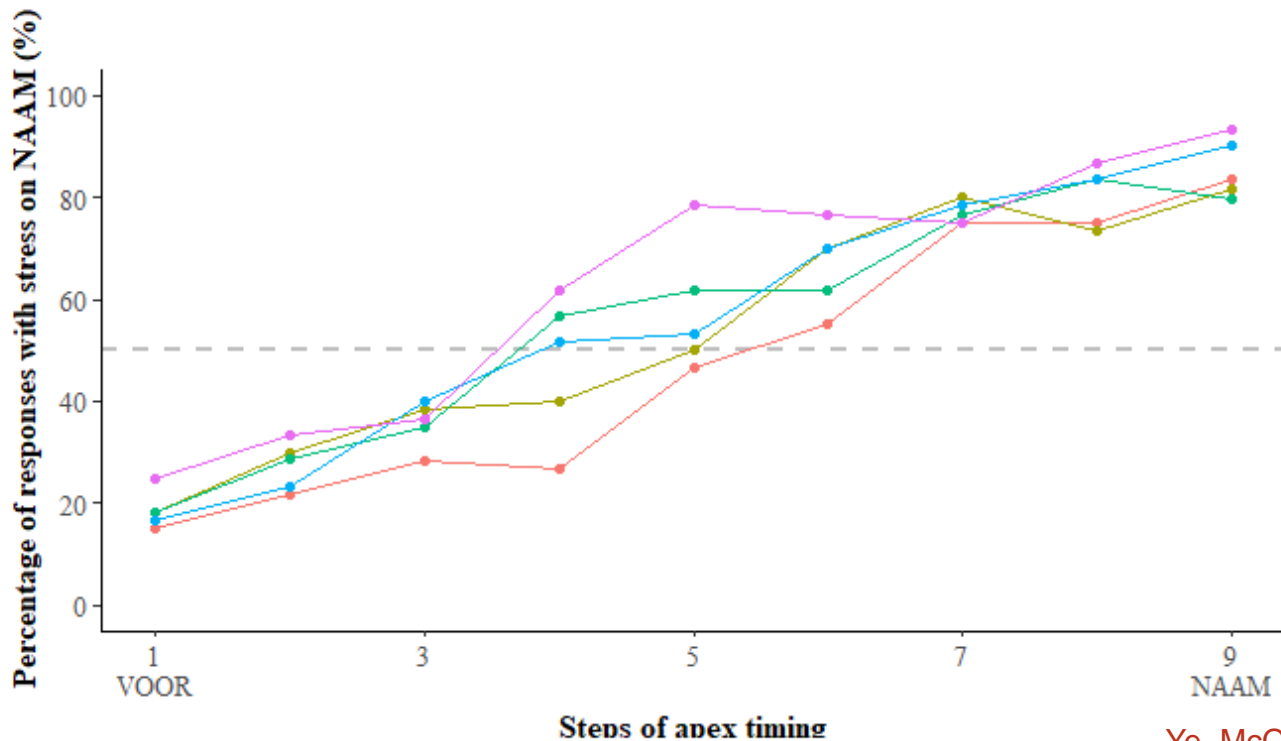
ONGOING:

stress attracts perceived gestural timing



Ye, McQueen, & Bosker, *in prep.*

stress attracts perceived gestural timing



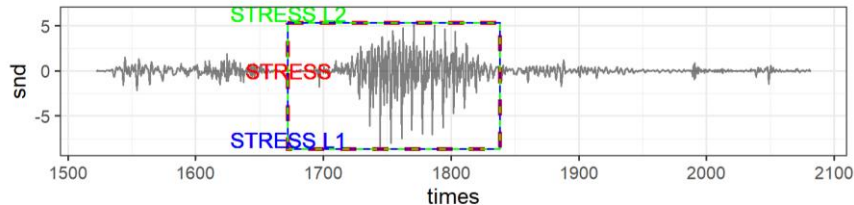
Ye, McQueen, & Bosker, *in prep.*

Steps of pitch and intensity cues to lexical stress — 1 — 3 — 4 — 5 — 7

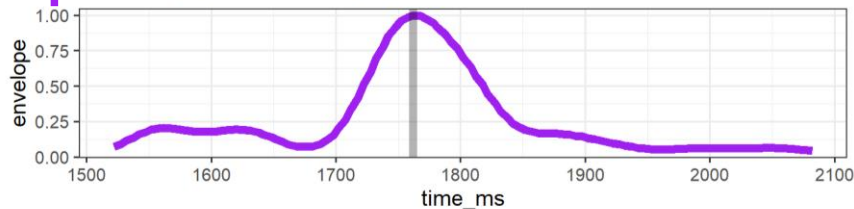
ONGOING: gesture-speech synchrony in L2

Spanish: “his.TÓ.ri.co”

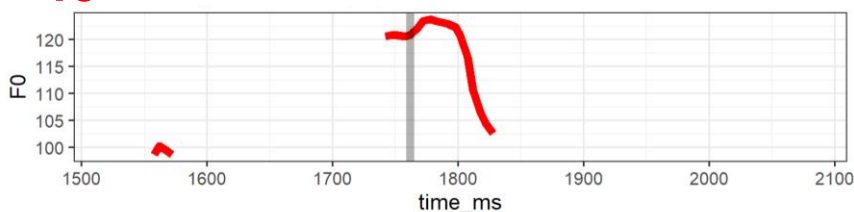
- How do L2 learners acquire L2 prosody in the auditory and gestural systems?



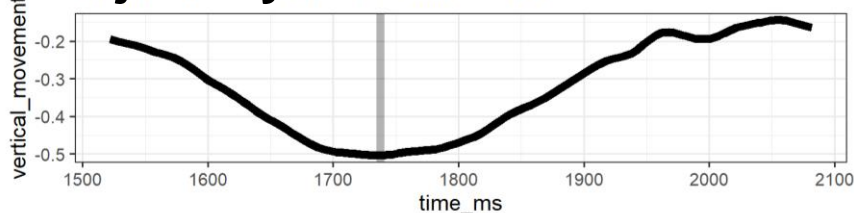
amplitude



f_0



hand trajectory



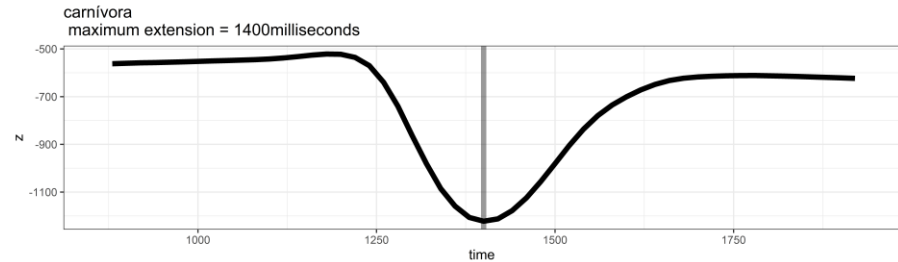
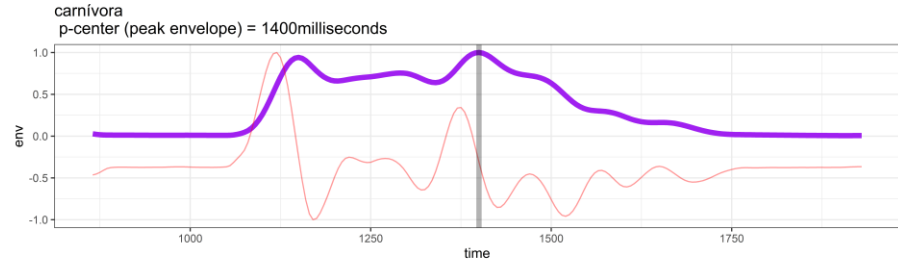
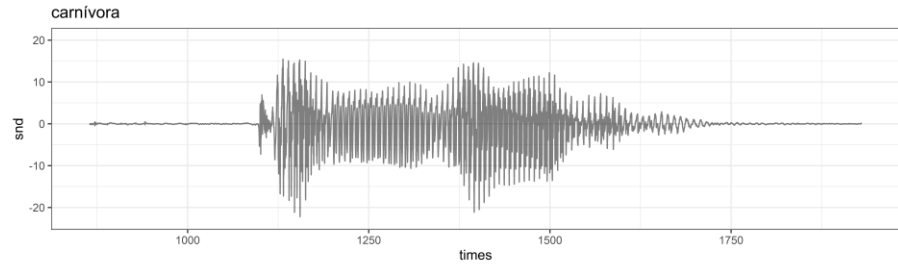


Example data

stress-mismatch

Dutch: *car.ni.VOOR*

Spanish: *car.NÍ.vo.ra*



TEAM SCIENCE:
Bosker, Hoetjes, Pouw, Van Maastricht, *subm.*
<https://osf.io/w2ezs>,

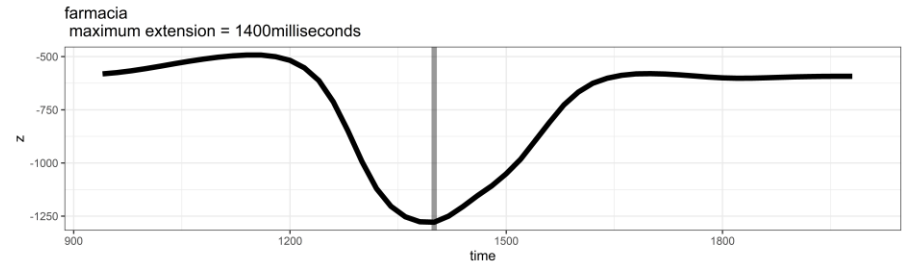
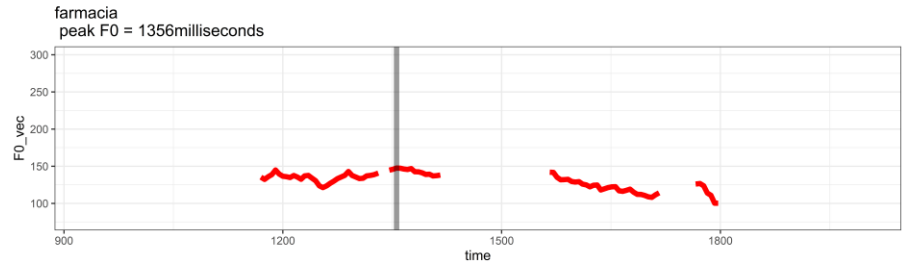
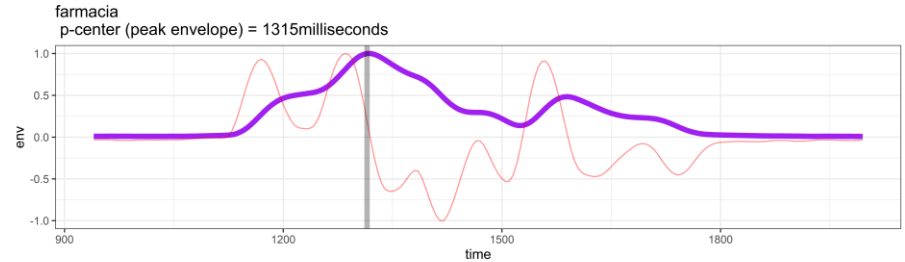
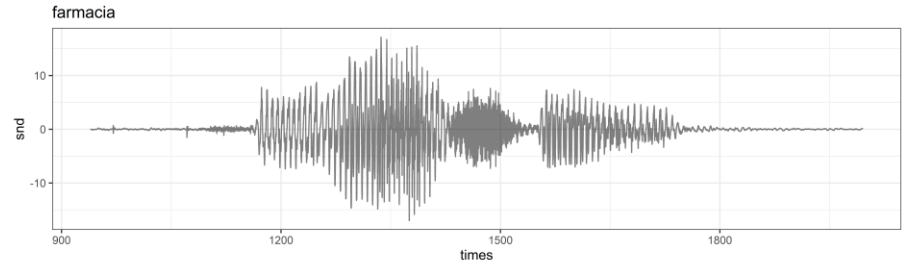


Example data

stress-mismatch

Dutch: *far.ma.CIE*

Spanish: *far.MA.cia*



TEAM SCIENCE:
Bosker, Hoetjes, Pouw, Van Maastricht, *subm.*
<https://osf.io/w2ezs>,

Results: gesture-speech synchrony



L1 attraction

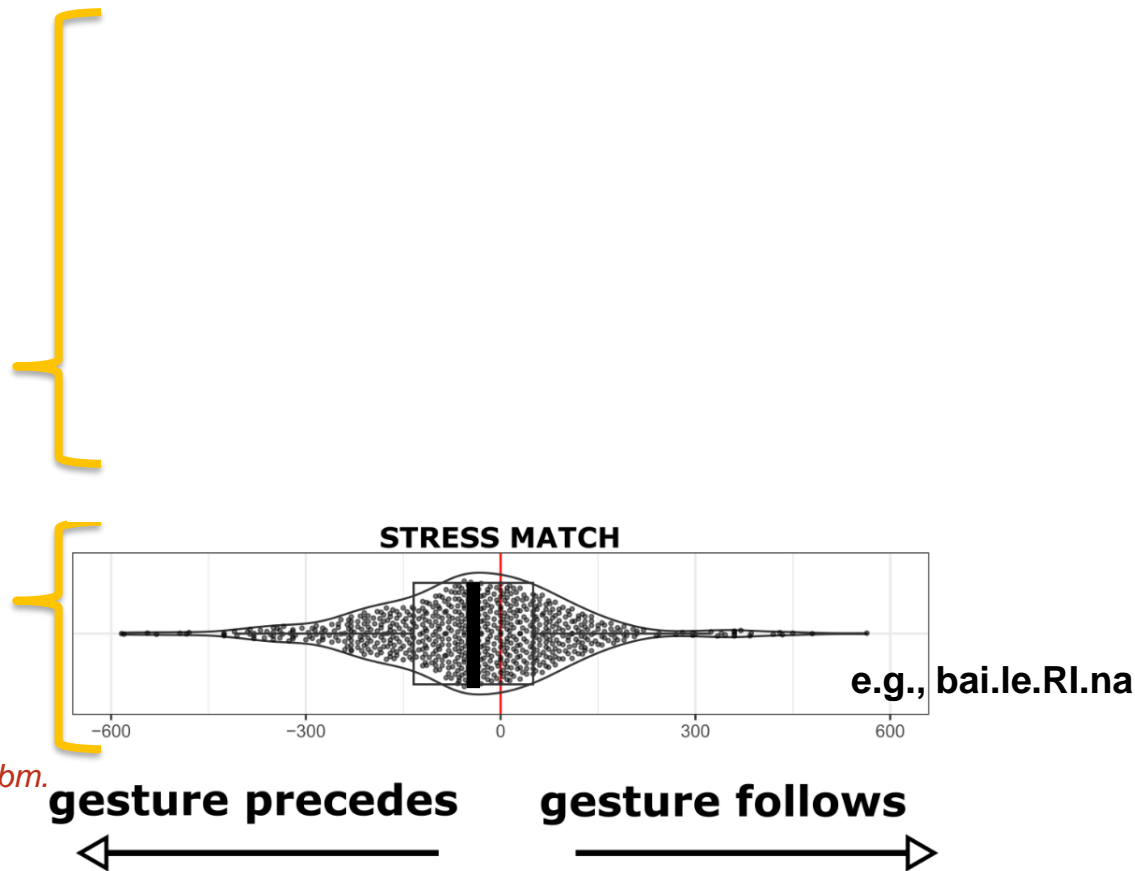
stress-mismatching cognates,
but correct L2 stress placement

stress-matching cognates
with correct L2 stress placement

TEAM SCIENCE:

Bosker, Hoetjes, Pouw, Van Maastricht, *subm.*

<https://osf.io/w2ezs>,



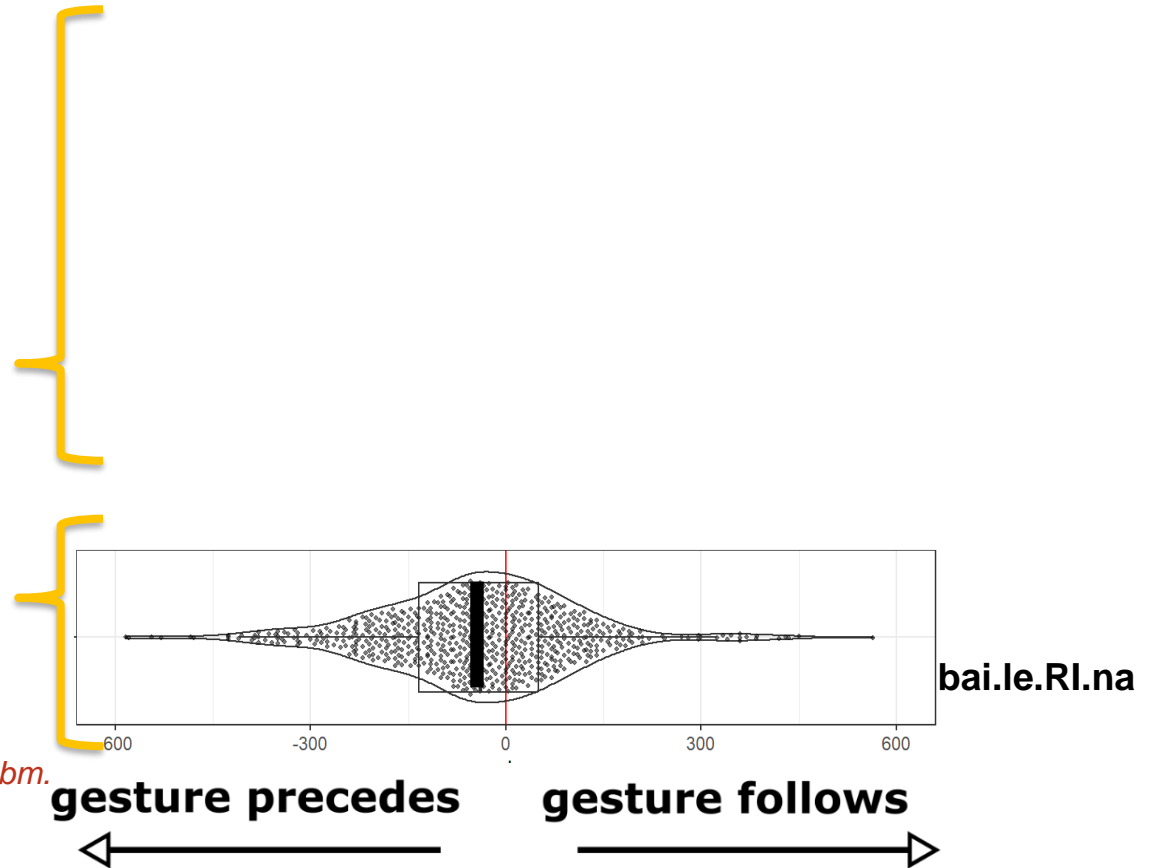
Results: gesture-speech synchrony



L2 attraction

stress-mismatching cognates,
but ***incorrect*** L2 stress placement
(i.e., on L1 target syllable)

stress-matching cognates
with correct L2 stress placement



TEAM SCIENCE:

Bosker, Hoetjes, Pouw, Van Maastricht, *subm.*

<https://osf.io/w2ezs>,



Prosody in the hands: L2

- When L2 stress placement was **correct**, the hands showed **L1 attraction**
 - ‘**kinematic accent**’
- When L2 stress placement was **incorrect**, erroneously stressing the L1 target, the hands showed **L2 attraction**

TEAM SCIENCE:

Bosker, Hoetjes, Pouw, Van Maastricht, *subm.*

<https://osf.io/w2ezs>,



Prosody in the hands: L2

- Snapshot of an unstable, developing, multimodal prosody system
- Overall strong evidence for multimodal gesture-speech coupling
- Yet separate unimodal L1/L2 attraction is reliably detectable
- Unique timing regimes for gestures vs. spoken prosody?
- Does a ‘kinematic accent’ of an L2 talker influence lexical stress perception in L1 listeners?

TEAM SCIENCE:

Bosker, Hoetjes, Pouw, Van Maastricht, *subm.*

<https://osf.io/w2ezs>,



Wrap-up of today

- Prosody is visible in articulatory movements, facial cues and expressions, and in the body.
- Multisensory cue weighting in audiovisual prosody perception, depending on the listening conditions
- Simple up-and-down hand gestures can influence what you hear.



Wrap-up of the course

- Prosody in Speech Perception
- Course aim: to reveal the central role that prosody plays in low-level speech perception and spoken word recognition.
- Course objectives:
 - to be familiar with key concepts in the area of speech prosody and speech perception
 - to be familiar with recent advances and paradigms in the speech perception literature
 - to understand how prosody influences the perception of vowels, consonants, and words
 - to understand the different processing mechanisms that underlie these influences
 - to understand the open issues and debates in the field of speech perception



Wrap-up of the course

- Normalization, neural speech tracking, prediction, talker-specific learning, audiovisual integration...
 - ...are tightly interconnected.
 - ...are not exclusive
 - ...showcase the ‘cognitive toolkit’ listeners have at their disposal

Prosody matters!

WANNA
JOIN US?

Hans Rutger Bosker

Speech Perception in Audiovisual Communication [SPEAC] lab

Donders Institute, Radboud University, Nijmegen, The Netherlands

<https://hrbosker.github.io>

hansrutger.bosker@donders.ru.nl

