

Toward a grammar of timing in speech production

Christopher Geissler

Department of Linguistics

Carleton College

September 21, 2023

Slides available on cageissler.github.io/resources

Roadmap

- **Phonology and articulatory gestures**
- **Coordinating gestures: the Coupled Oscillator Model**
- **Problems**
 - **Tibetan tone study**
- **Toward solutions: Analysis-by-synthesis**
- **Conclusion**

LING 110-level

Categorical behavior

- In German, voiced consonants are voiceless when they occur at the end of words (but not elsewhere):
 - *Maus* ‘mouse’ [maʊs̥], but plural *Mäuse* [mɔʏzə]
 - *Rad* ‘wheel’ [ʀat̥], but plural *Räder* [ʀɛdɐ]
 - compare:
Rat ‘council’ [ʀat̥], but plural *Räte* [ʀɛtə]

LING 110-level

Categorical behavior

- In German, voiced consonants are voiceless when they occur at the end of words (but not elsewhere):
 - *Maus* ‘mouse’ [maʊs̥], but plural *Mäuse* [mɔʏzə]
 - *Rad* ‘wheel’ [ʀat̥], but plural *Räder* [ʀɛdɐ]
 - compare:
Rat ‘council’ [ʀat̥], but plural *Räte* [ʀɛtə]

Linguists are really good at this

LING 217-level

Probabilistic behavior

- In English, t/d at the end of a word sometimes isn't there
 - *rift* = [ɹɪft̚] or [ɹɪf_]; *build* = [bɪɫd] or [bɪɫ]
 - More likely among some groups
 - More likely in some social contexts
 - More likely around some sounds
 - More likely in *mist* than in *missed*

LING 217-level

Probabilistic behavior

- In English, t/d at the end of a word sometimes isn't there
 - *rif*t = [ɹɪft̚] or [ɹɪf_]; *bi*ld = [bɪɫd] or [bɪɫ]
 - More likely among some groups
 - More likely in some social contexts
 - More likely around some sounds
 - More likely in *mist* than in *missed*

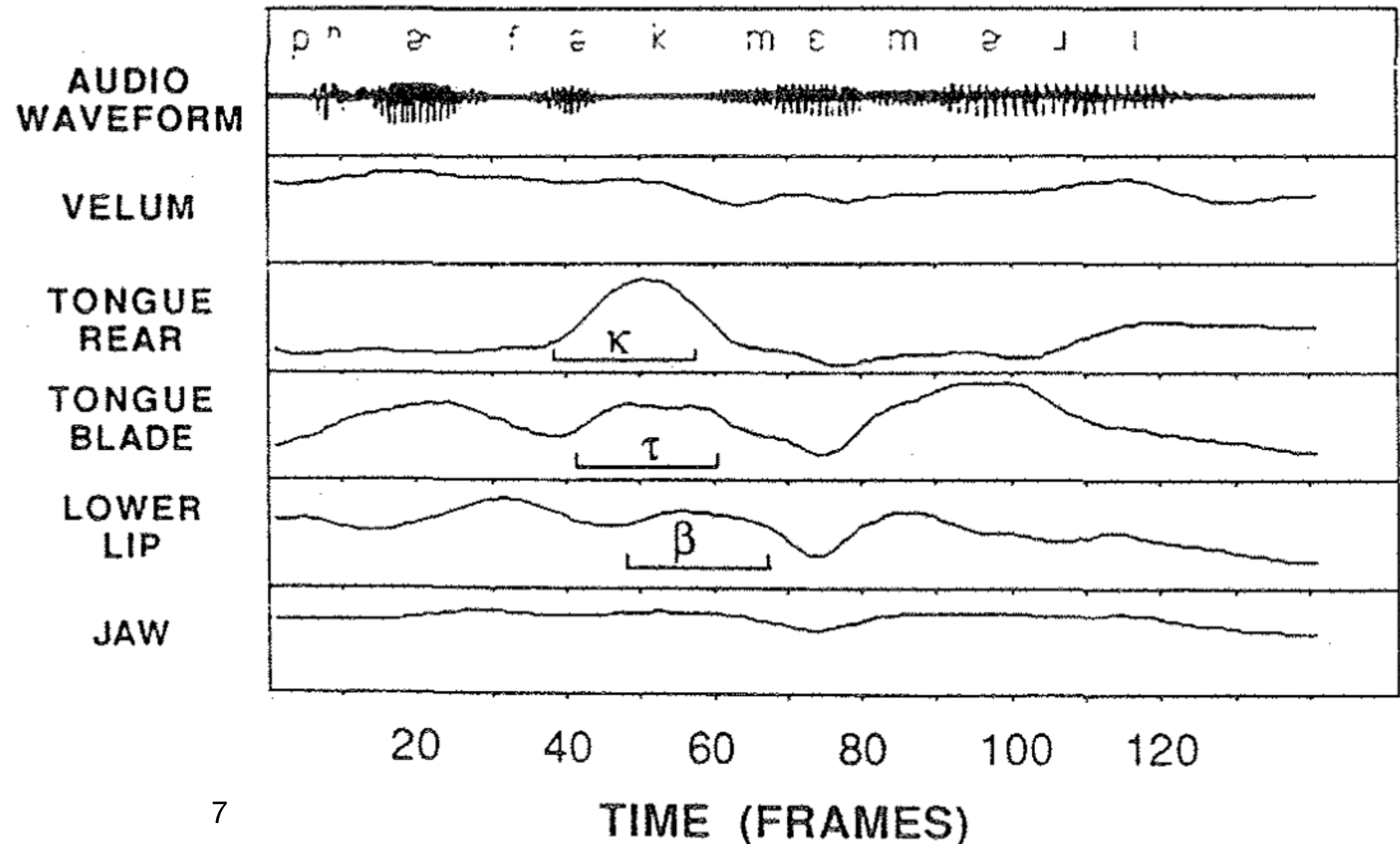
Linguists get excited about this

...uh-oh

- *Perfect memory*
- At least some “deleted” t’s/d’s are visible in articulation, but not in acoustics
- (Actually it’s most)

Midsagittal sections

(Browman & Goldstein 1988, Purse 2019)



...uh-oh

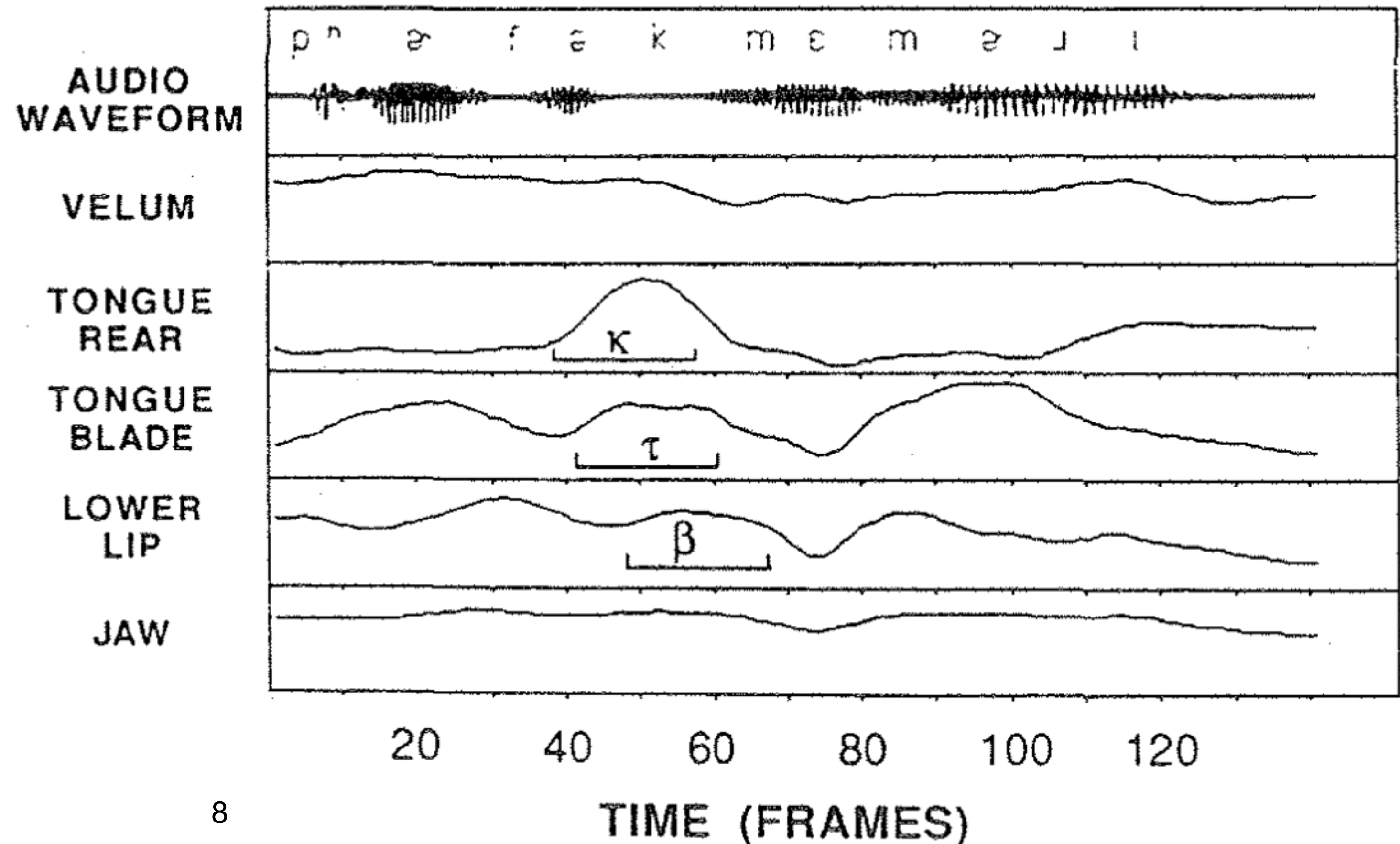
- *Perfect memory*
- At least some “deleted” t’s/d’s are visible in articulation, but not in acoustics
- (Actually it’s most)

Midsagittal sections

(Browman & Goldstein 1988, Purse 2019)

Gestures!

... but how are they coordinated?



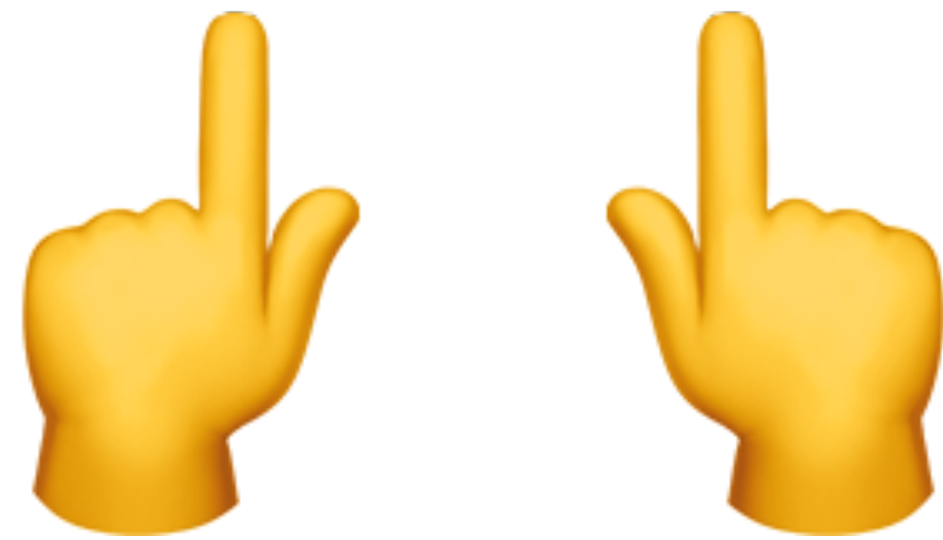
Gestures

- Definition 1: controlled movement of the vocal tract in speech
- Definition 2: abstract, hierarchical control unit for linguistically-defined goal-directed movement *(Pouplier 2020)*
- Motor equivalence, equifinality

Roadmap

- Phonology and articulatory gestures
- **Coordinating gestures: the Coupled Oscillator Model**
- Problems
 - Tibetan tone study
- Toward solutions: Analysis-by-synthesis
- Conclusion

Bimanual tapping interlude



Oscillators

- Synchronization in non-speech and speech movements:
 - “pa... pa... pa... pa.pa[...]pa.pa.pa.pa”
 - “ap... ap... ap... ap.ap[...]pa.pa.pa.pa”
- Tapping: “in-phase” more stable than “anti-phase”
(both more stable than any other phasing)
... in speech too?

CV vs. VC syllables

in-phase

[pa]	
LIPS	Labial closure
TONGUE TIP	
TONGUE BODY	pharyngeal wide

anti-phase

[ap]	
LIPS	labial closure
TONGUE TIP	
TONGUE BODY	pharyngeal wide

CV vs. VC syllables

in-phase

[pa]	
LIPS	Labial closure
TONGUE TIP	
TONGUE BODY	pharyngeal wide

/ti/ 'tea'	
LIPS	
TONGUE TIP	alveolar closure
TONGUE BODY	palatal narrow

anti-phase

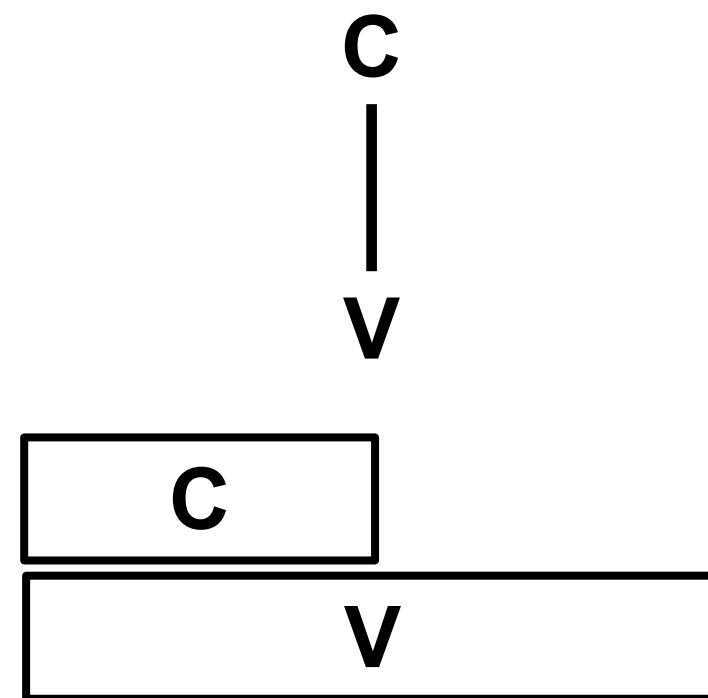
[ap]	
LIPS	labial closure
TONGUE TIP	
TONGUE BODY	pharyngeal wide

/it/ 'eat'	
LIPS	
TONGUE TIP	alveolar closure
TONGUE BODY	palatal narrow

CV vs. VC syllables

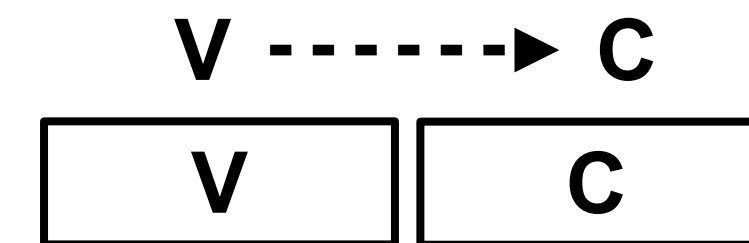
in-phase

[pa]	
LIPS	Labial closure
TONGUE TIP	
TONGUE BODY	pharyngeal wide



anti-phase

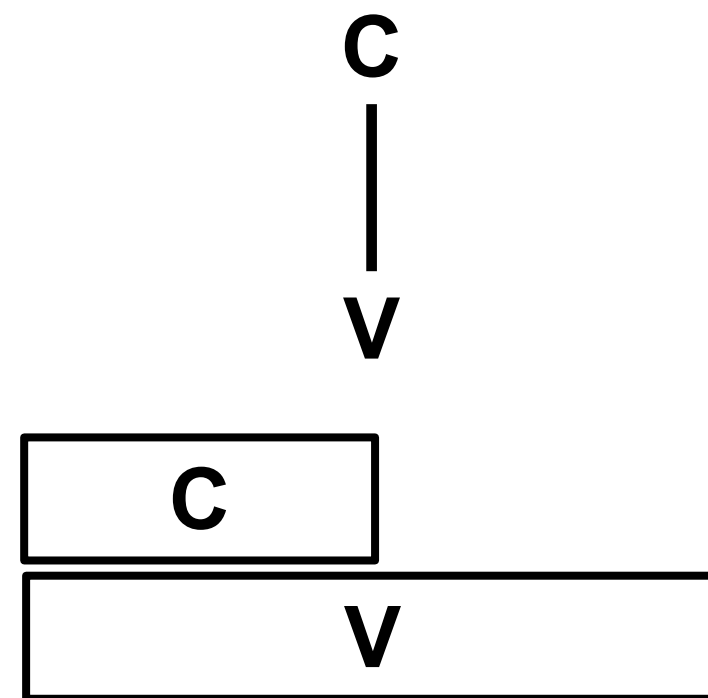
[ap]	
LIPS	labial closure
TONGUE TIP	
TONGUE BODY	pharyngeal wide



CV vs. VC syllables

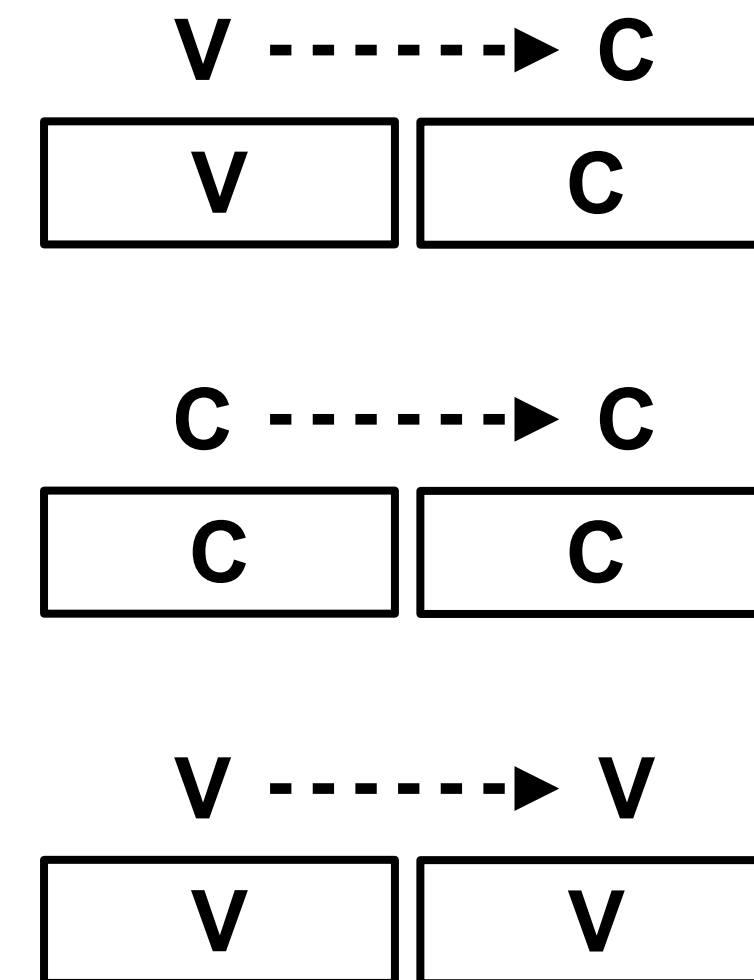
in-phase

[pa]	
LIPS	Labial closure
TONGUE TIP	
TONGUE BODY	pharyngeal wide



anti-phase

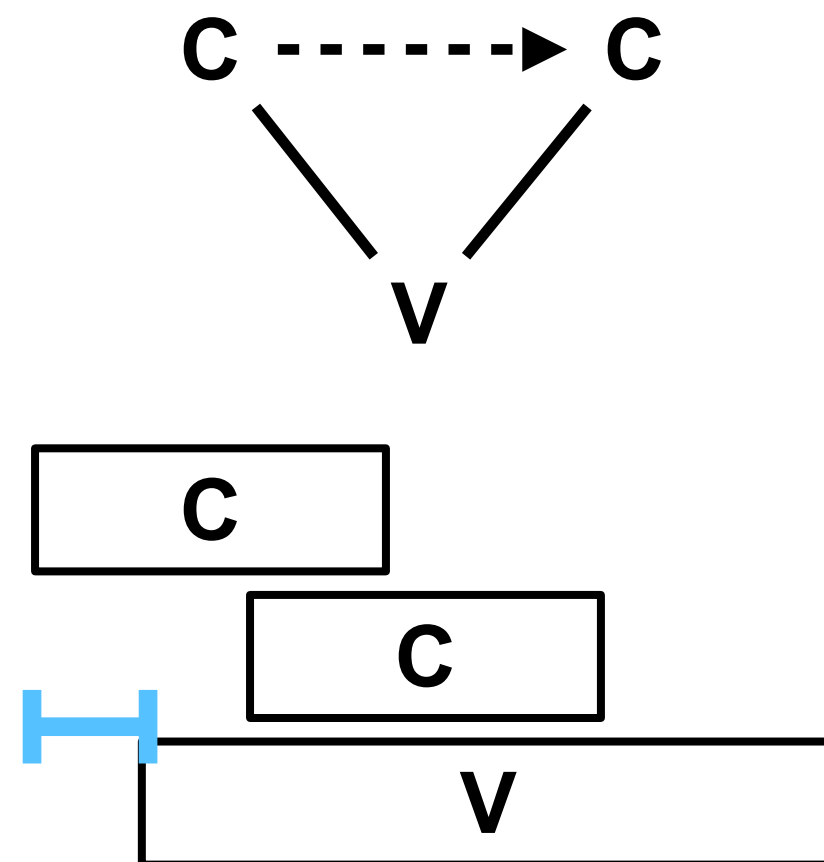
[ap]	
LIPS	labial closure
TONGUE TIP	
TONGUE BODY	pharyngeal wide



What about clusters?

- Empirically, onset clusters overlap

/spa/ 'spa'	
LIPS	labial closure
TONGUE TIP	alveolar critical
TONGUE BODY	pharyngeal wide



What about tone?

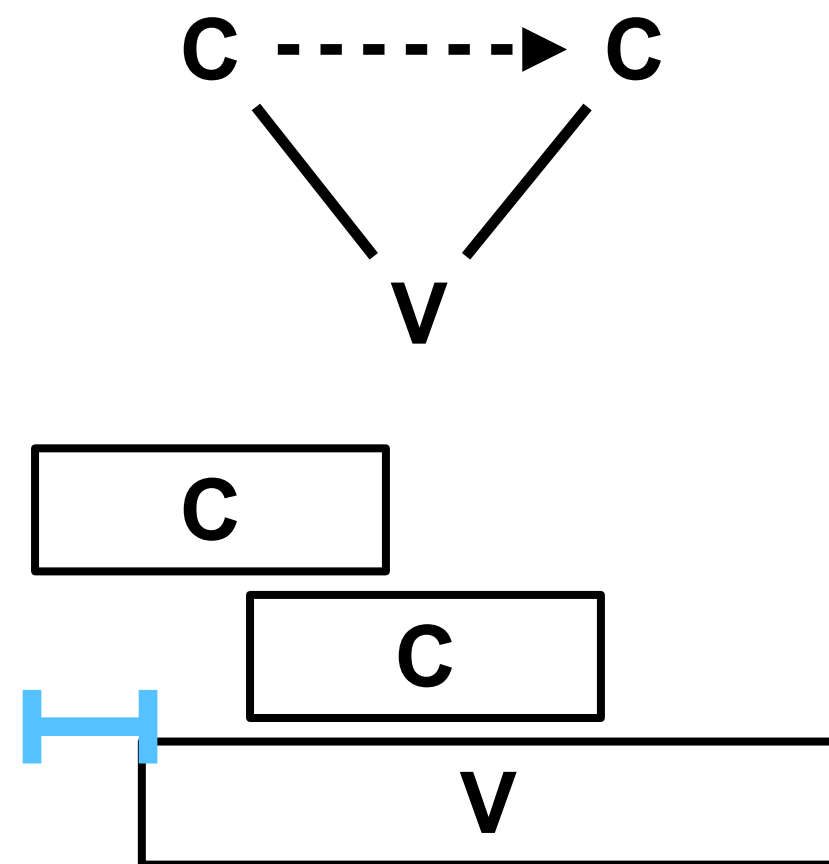
- Empirically, V lags following C
 - (In *lexical tone* languages only)

/pá/	
LIPS	labial closure
TONGUE TIP	
TONGUE BODY	pharyngeal wide
pitch (?)	high

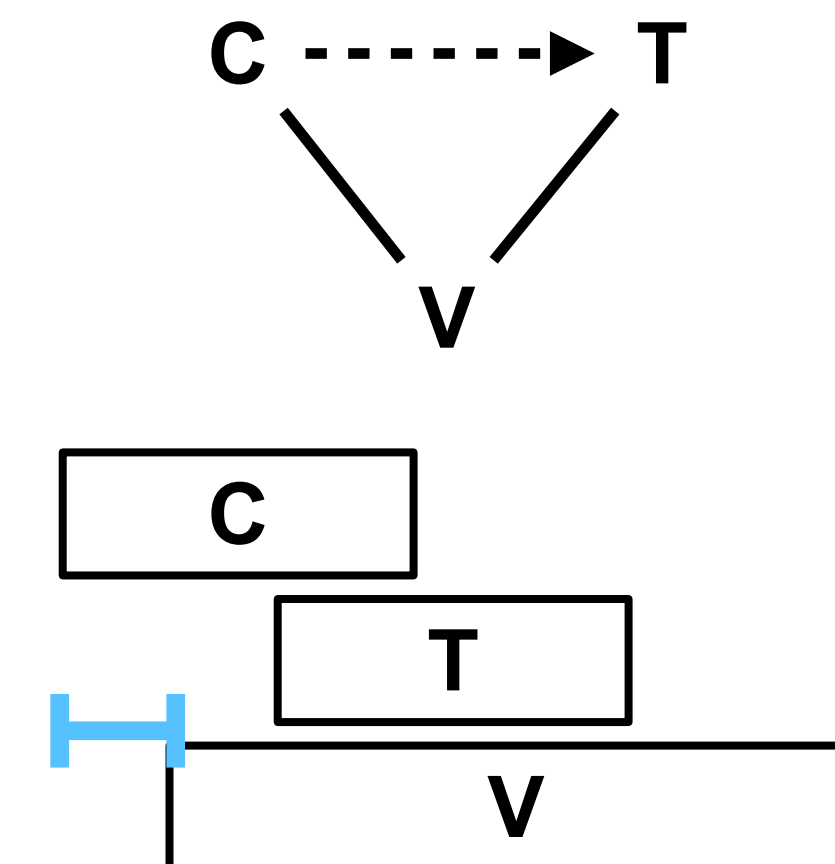
Competitive coupling account 🎉

- Unifies clusters and tone (neat for typology)
- Unifies syllables (and up?), contrast, and planning

/spa/ 'spa'	
LIPS	labial closure
TONGUE TIP	alveolar critical
TONGUE BODY	pharyngeal wide



/pá/	
LIPS	labial closure
TONGUE TIP	
TONGUE BODY	pharyngeal wide
pitch (?)	high



Roadmap

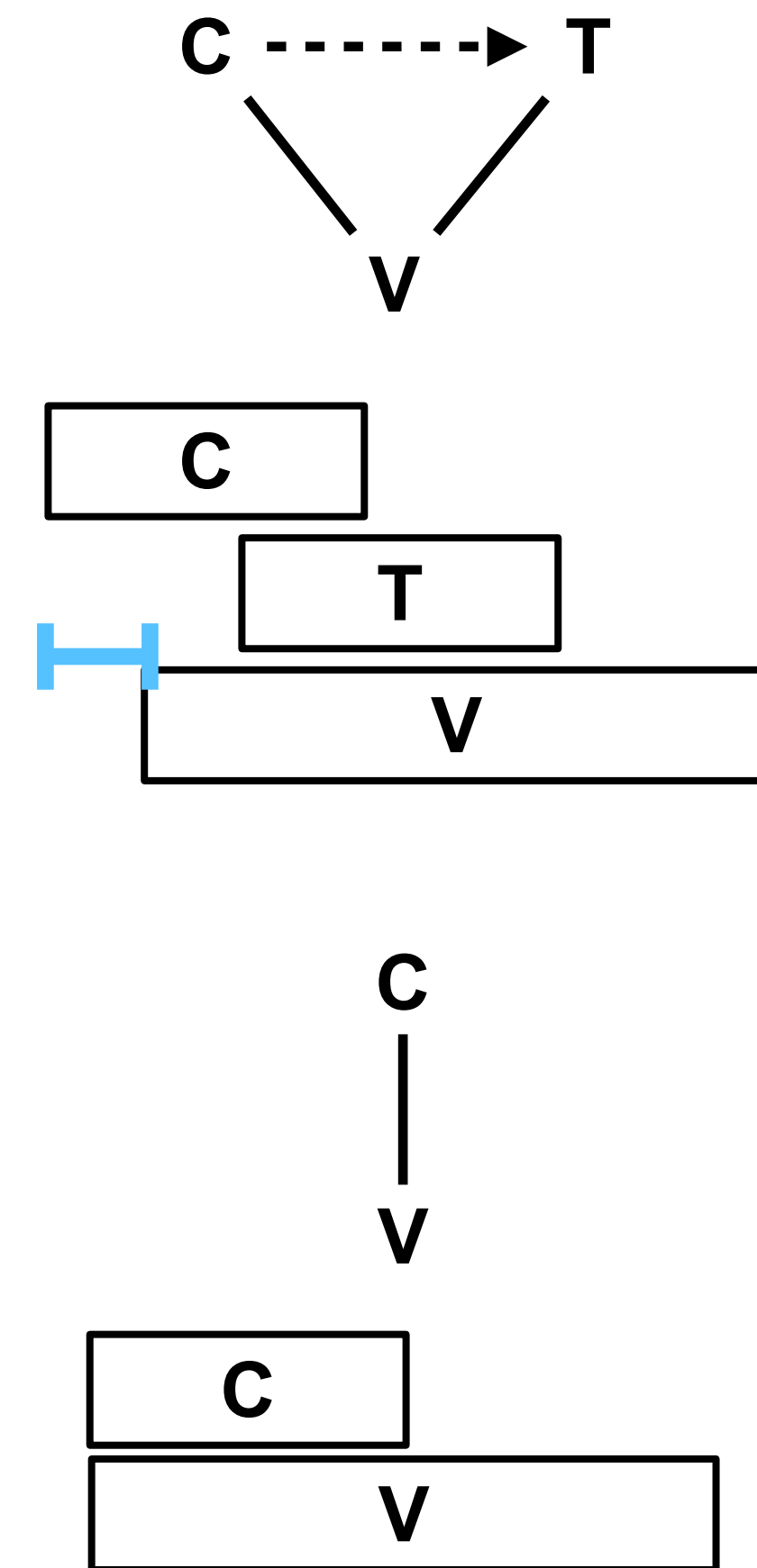
- Phonology and articulatory gestures
- Coordinating gestures: the Coupled Oscillator Model
- Problems
 - **Tibetan tone study**
- Toward solutions: Analysis-by-synthesis
- Conclusion

Doubts

- “Competitive coupling” ... is weird
 - Is it much better than just stipulating phasing?
 - Doesn't work for clusters of 3 +
- Can we really generalize from ‘papapapa’ to regular speech?
- Should we rely on *start* of a gesture or the *end* of a gesture?

Predictions of Coupled Oscillator Model

- If there is a tone gesture in a syllable:
 - C-V timing like in clusters:
C-V lag positive, $\sim 50\text{ms}$
- If there is no tone in that syllable:
 - Simultaneous C & V:
C-V lag $\sim 0\text{ms}$



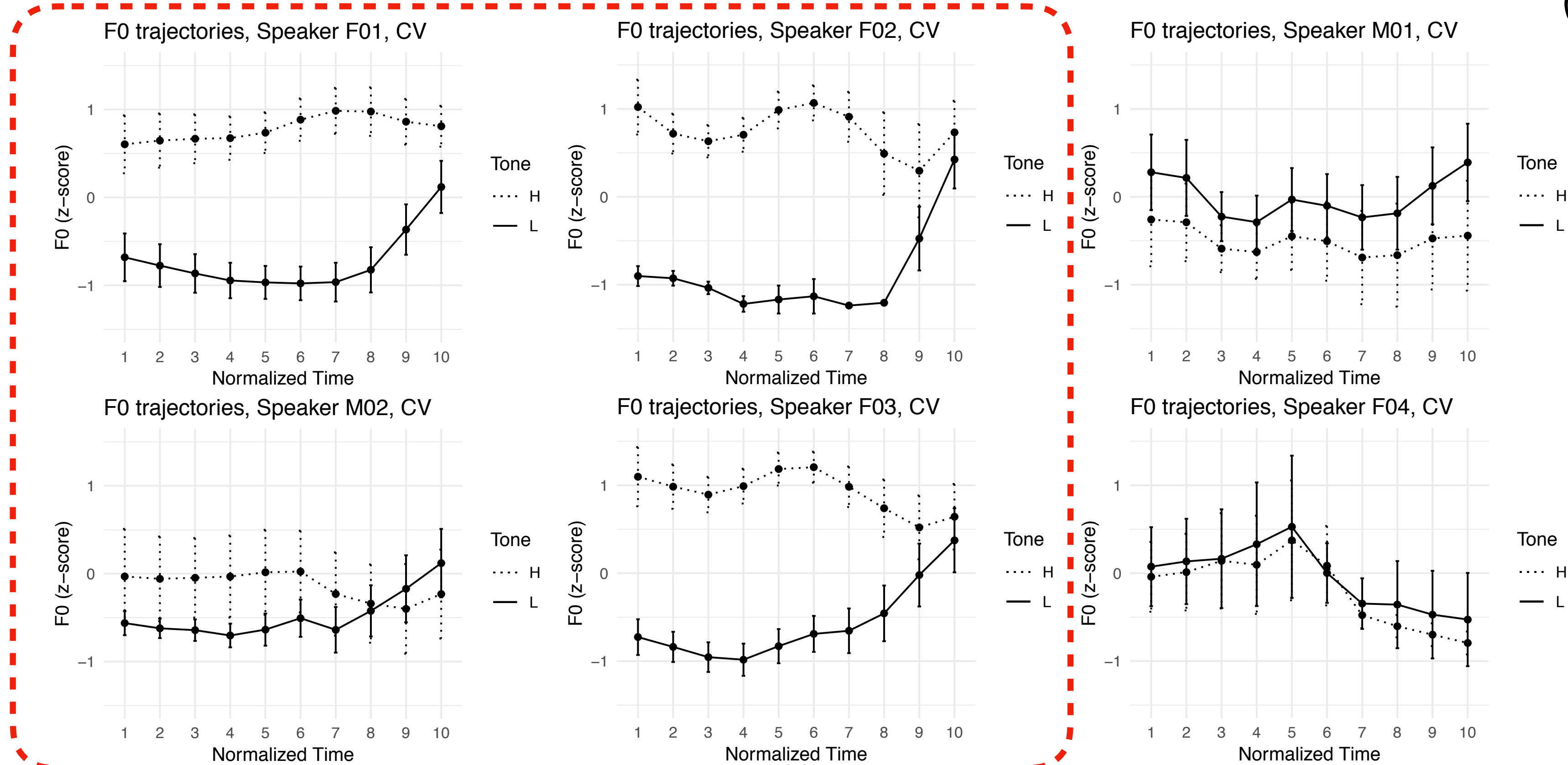
The perfect test case?

A language where some speakers produce tone and others don't

(Geissler 2019, 2021)

- 4 speakers produce a tone contrast, two do not (images: /mV/)

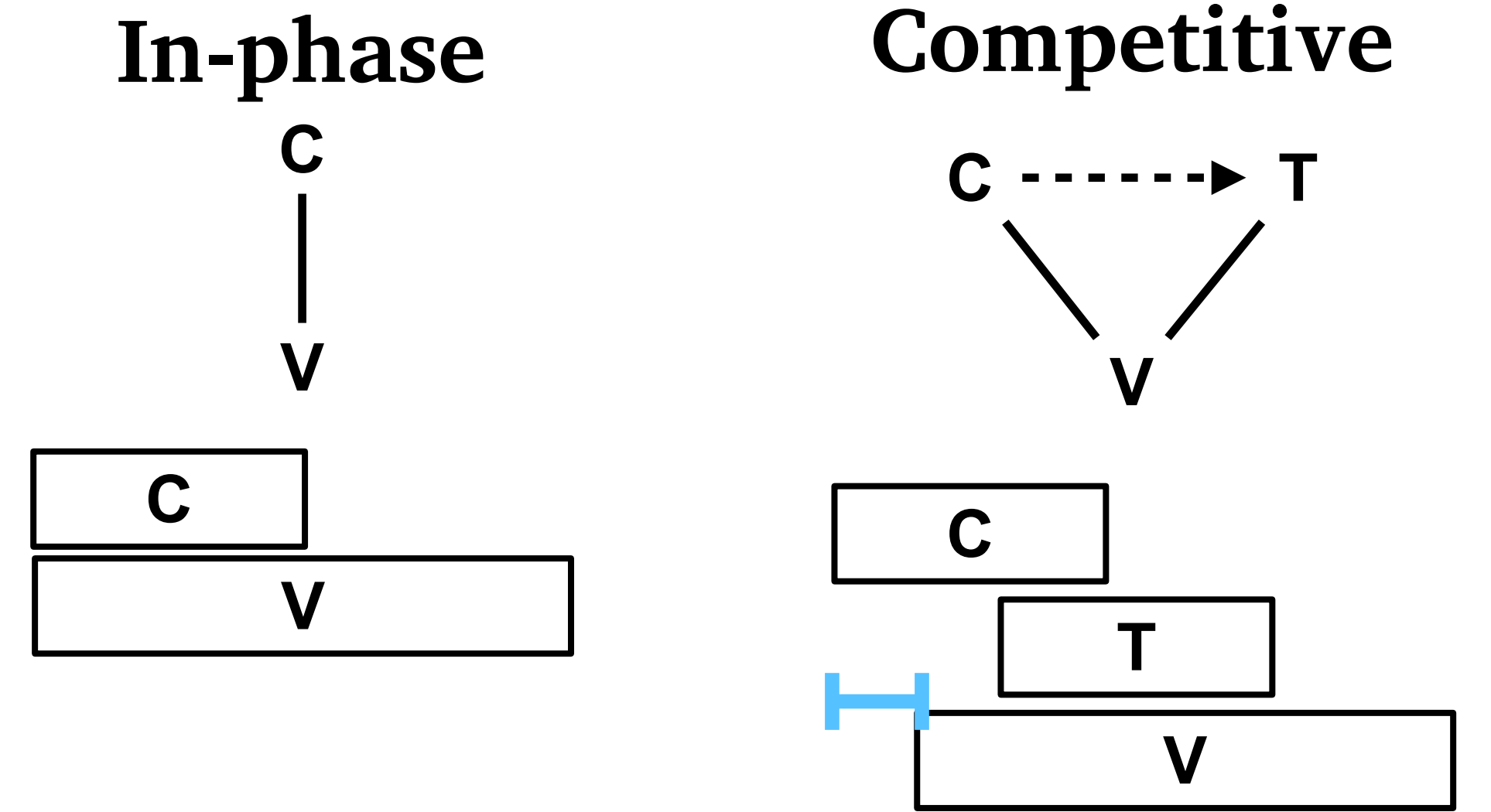
(Geissler et al. 2021)



EMA study

articulatory trajectories

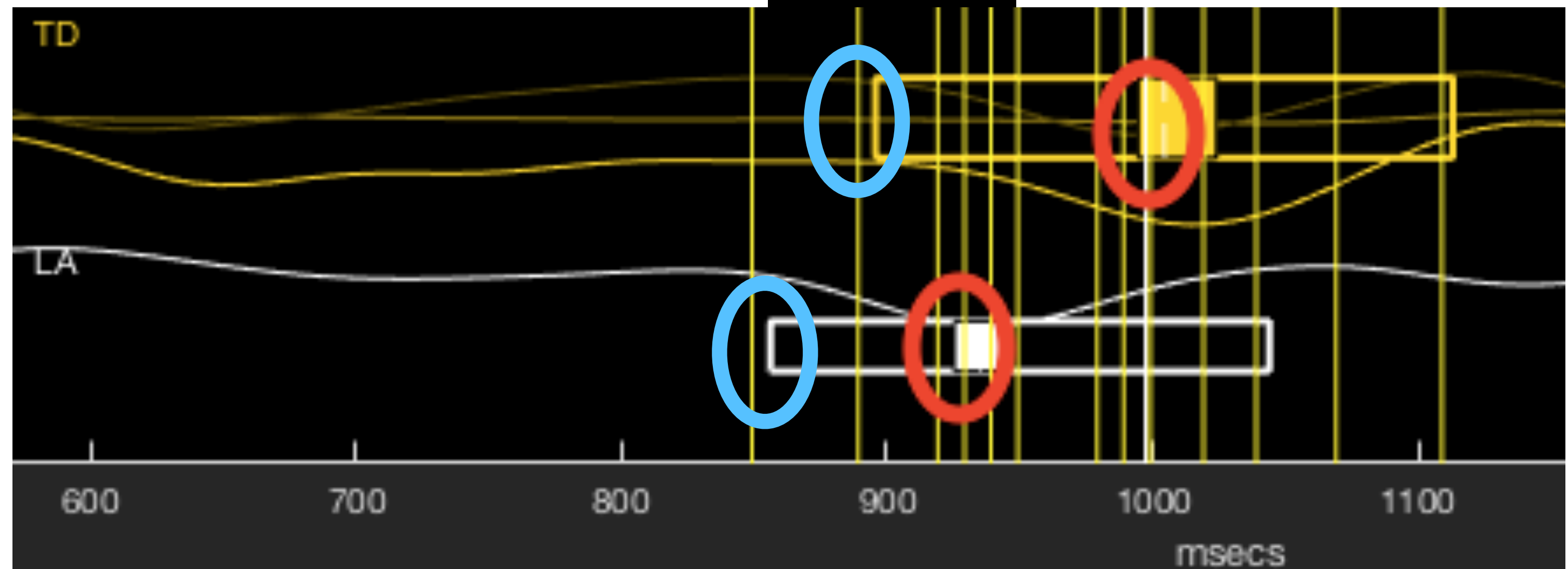
- [p p^h m]: distance between lip sensors
- [i] → [u o a]: tongue dorsum retraction
- H, L tones; 1- and 2-syllable words
- **C-V lag** as diagnostic of tone



[mu]



Tongue Dorsum front ↓ back
Lip Aperture open ↓ closed

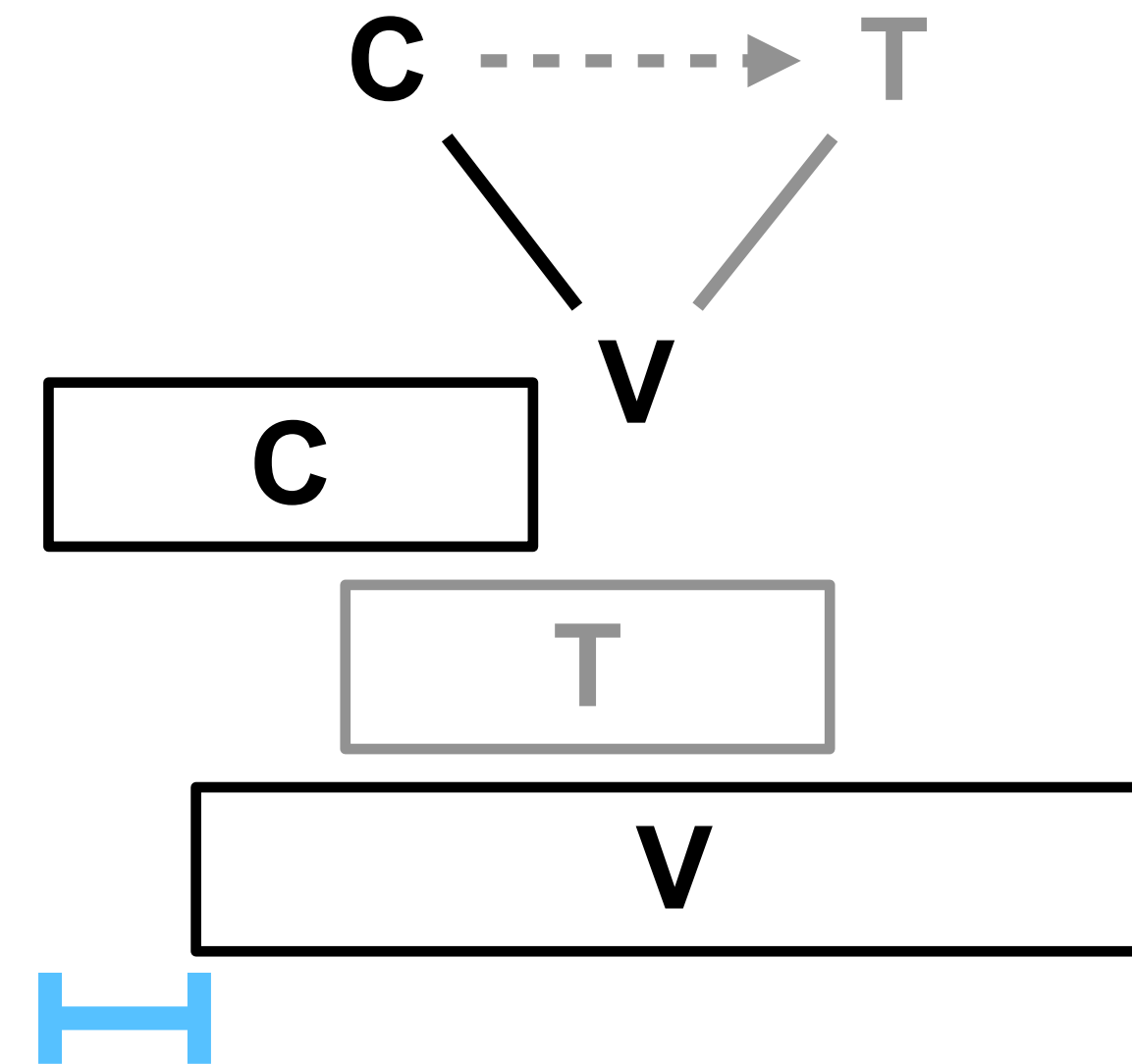
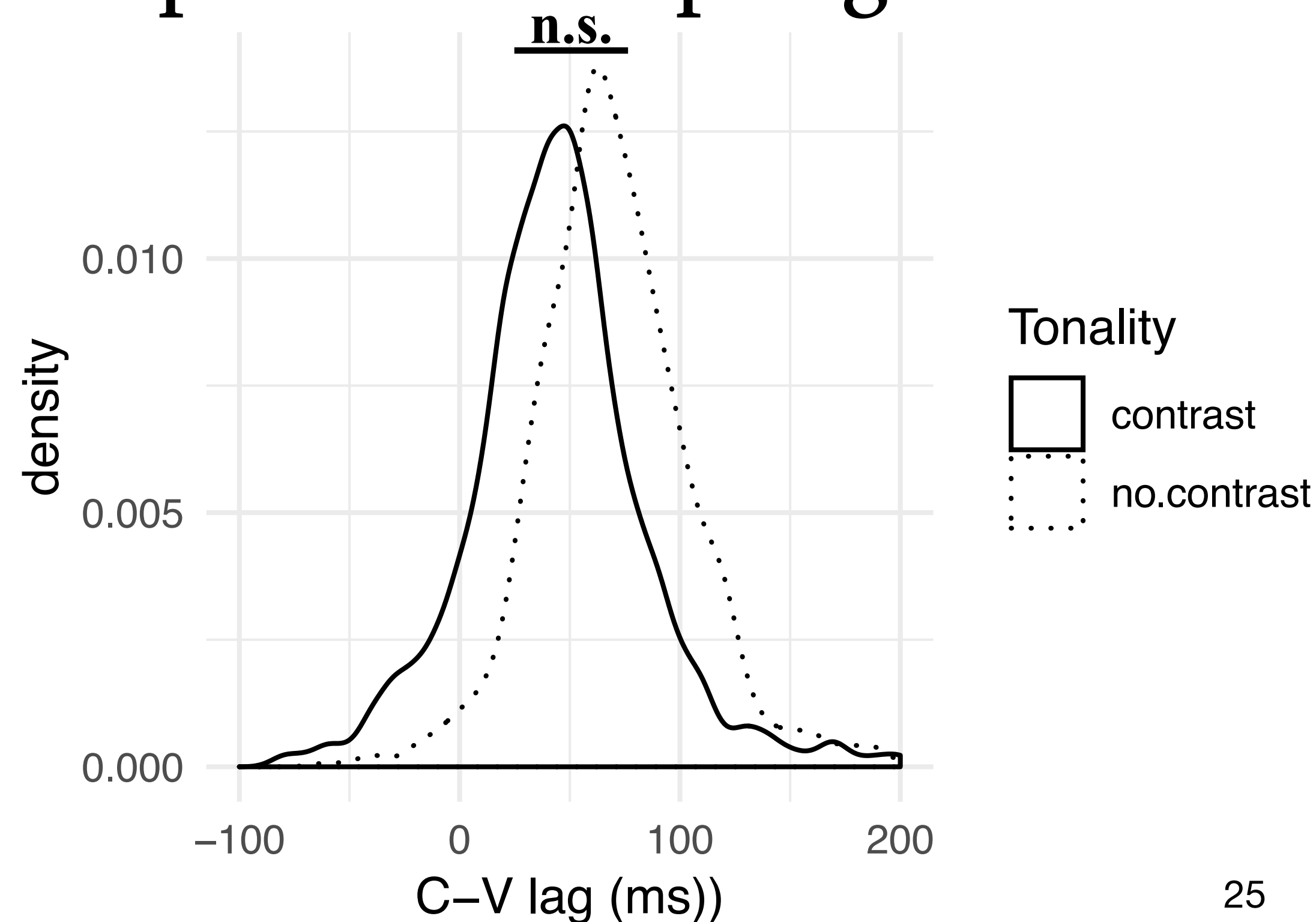


(Data: Zhang, Geissler, & Shaw 2019)

(Mview software: Tiede 2005)

Results: C-V lag

- There is a positive C-V lag... for speakers with *and* without the tone contrast (and in both tones)
- Competitive Coupling has no explanation for the 50ms lag



Cross-linguistic evidence (before)

No tone,
no C-V lag

Arabic

Catalan

English

German

Georgian

Italian

Romanian

Tone

Swedish

Serbian

C-V lag

Mandarin

Thai

Cross-linguistic evidence (after)

No tone,
no C-V lag

Arabic

Catalan

English

German

Georgian

Italian

Romanian

Tone

Swedish

Serbian

C-V lag

Mandarin

Thai

Tibetan

also Tibetan

Roadmap

- Phonology and articulatory gestures
- Coordinating gestures: the Coupled Oscillator Model
- Problems
 - Tibetan tone study
- **Toward solutions: Analysis-by-synthesis**
- Conclusion

There's another problem





There's another problem

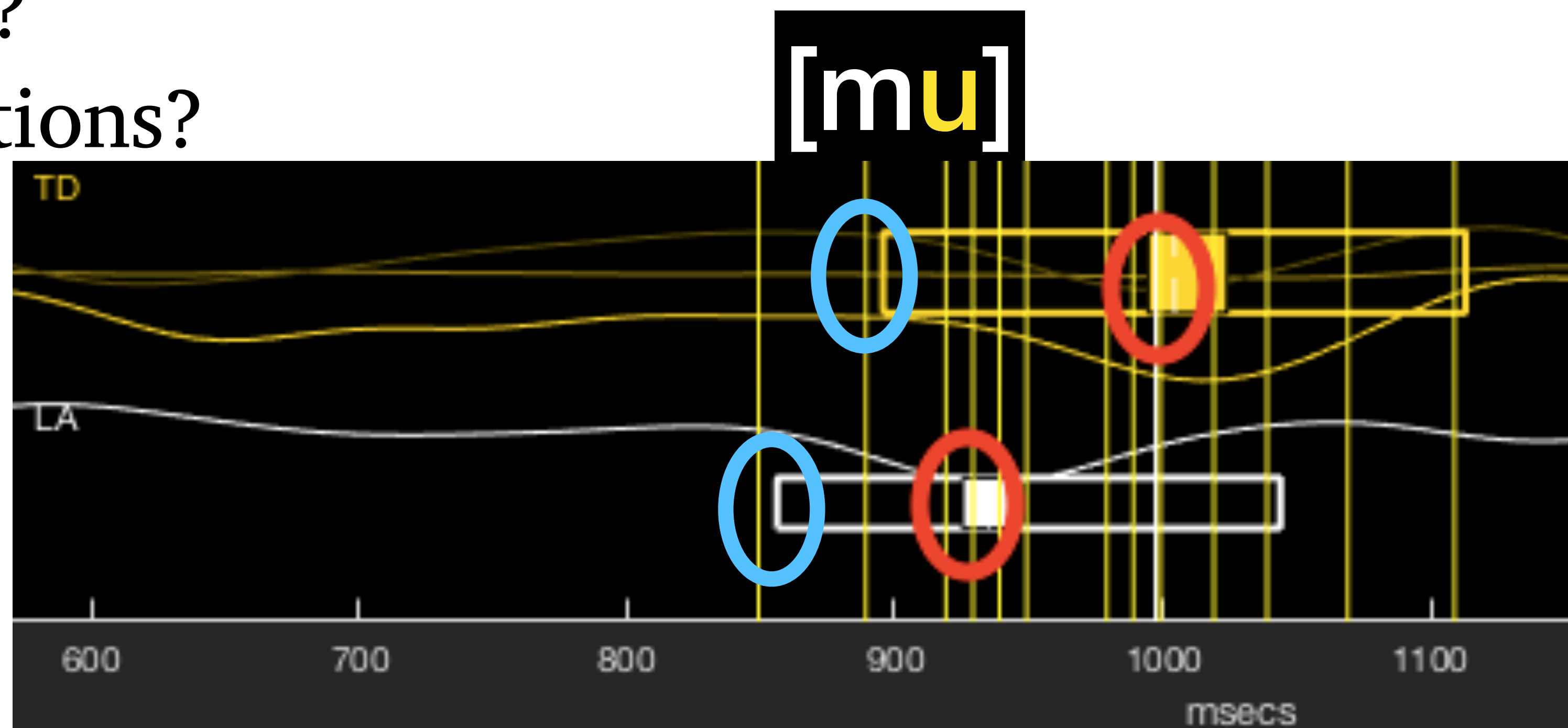
WHEN DOES A GESTURE START

- Velocity zero-crossing?
- Velocity 20% of peak?
- Acceleration maximum?
- Divergence from repetitions?

...

Tongue Dorsum front
↓
back

Lip Aperture open
↓
closed



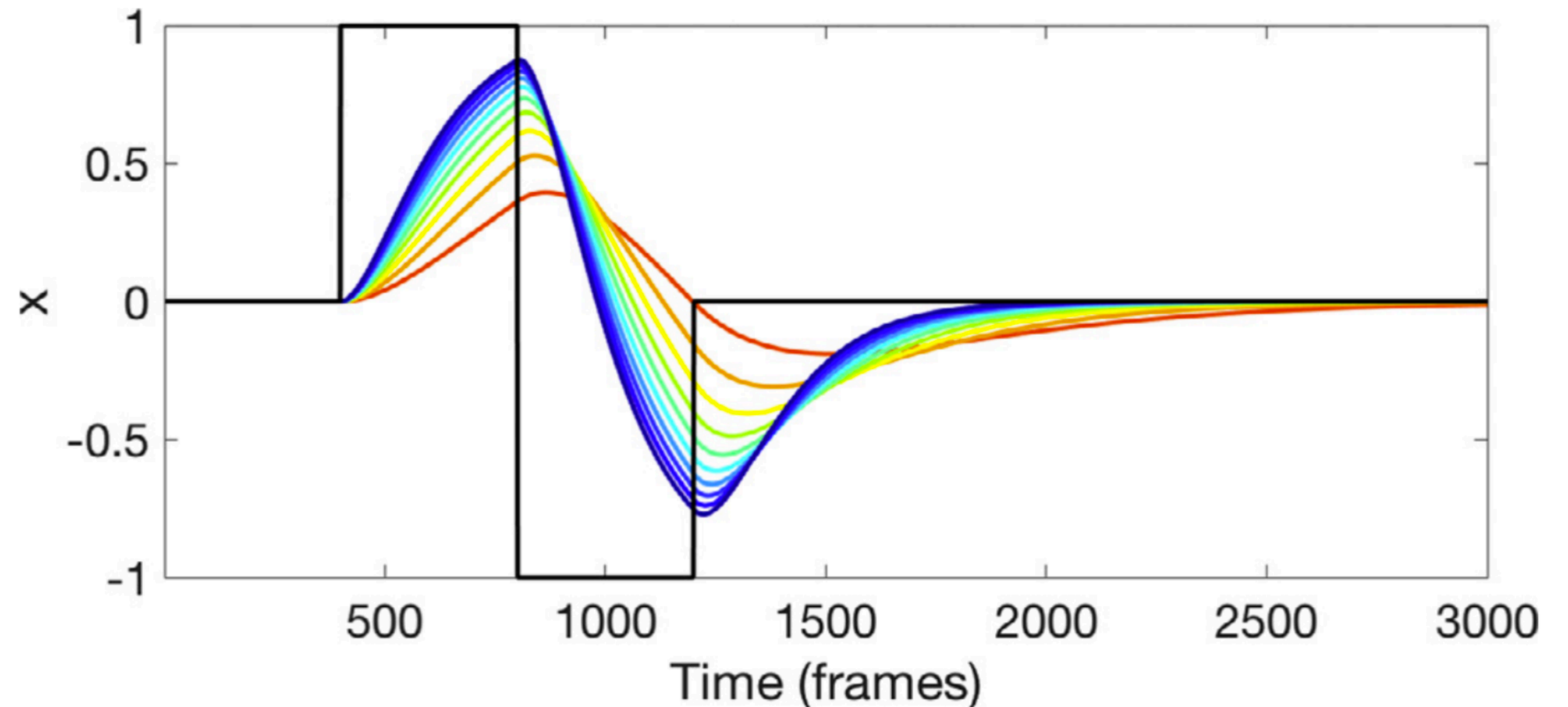
Measuring to modeling

(Haken et al. 1985, Saltzman & Munhall 1989, Nam & Saltzman 2003)

- Model kinematics as critically-damped mass-spring oscillator
- In-phase/anti-phase/etc. determine relative timing

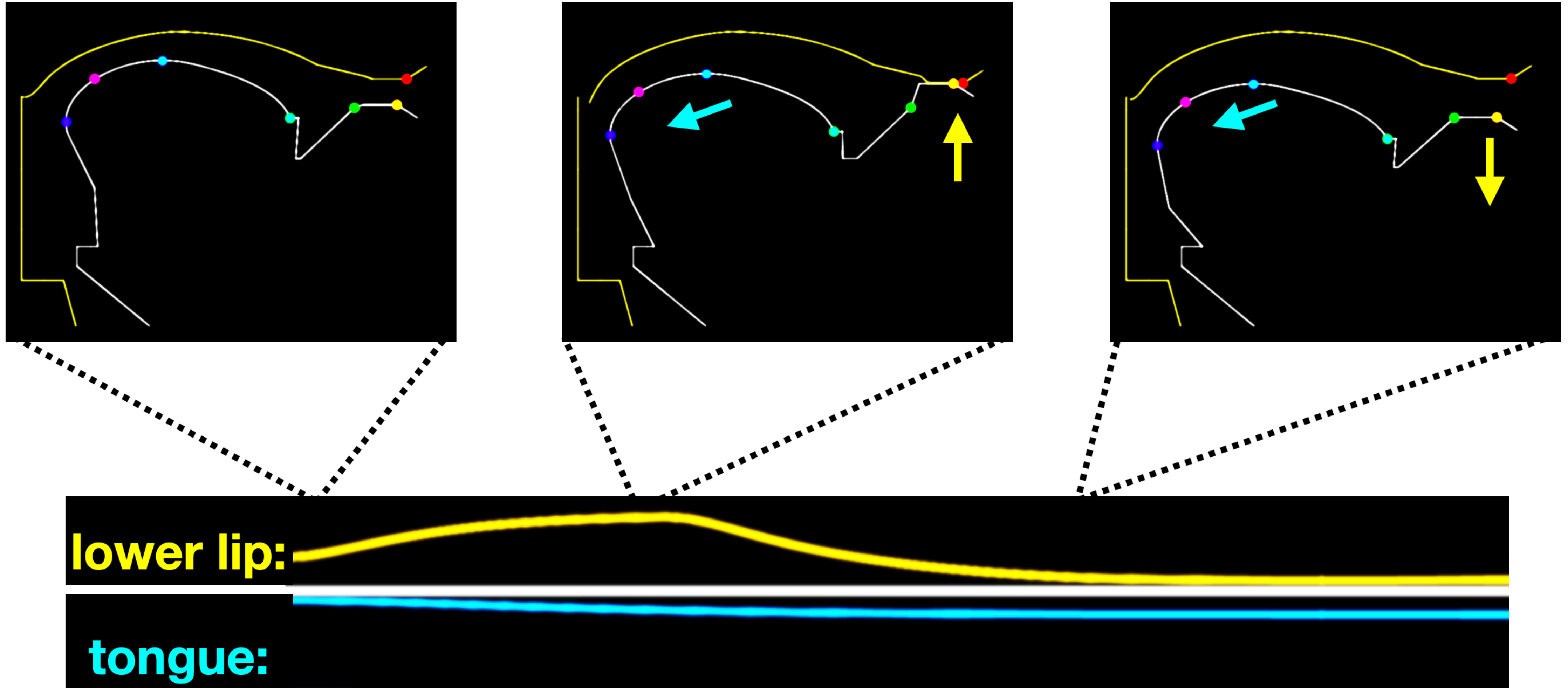
$$ma + bv + k(x - C) = 0$$

acceleration
velocity
position
stiffness
target



Articulatory simulation

TADA: Task Dynamics Application *(Nam et al. 2004)*

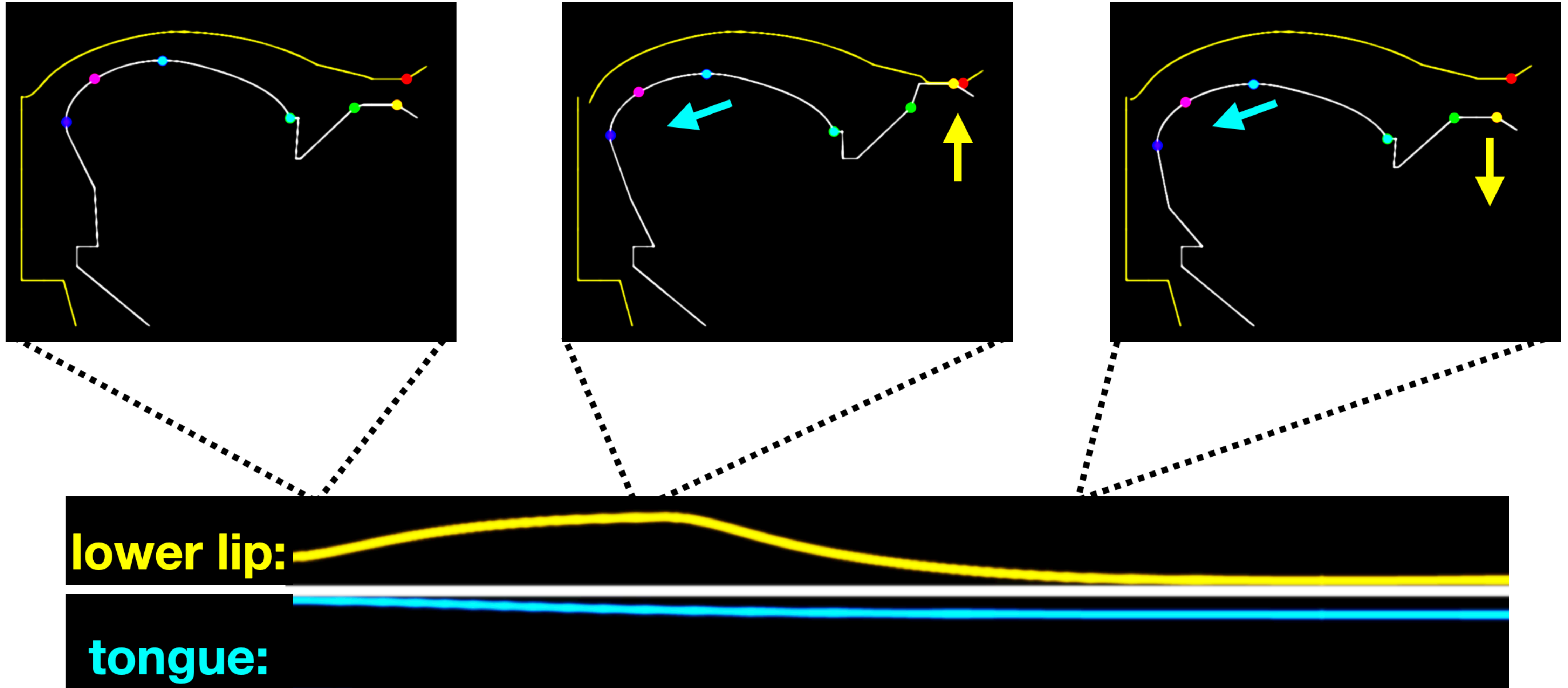


Articulatory simulation

TADA: Task Dynamics Application *(Nam et al. 2004)*

Images from a different study
sanity-checking the Tibetan
experiment results

(Geissler 2022)



Analysis-by-synthesis: <five>

- Diphthong targets can't be separated with kinematic data
- Make a simulation, then tweak it, → 34,000 simulations
Compare to 525 tokens from X-ray Microbeam Database

Bad fit

Good fit

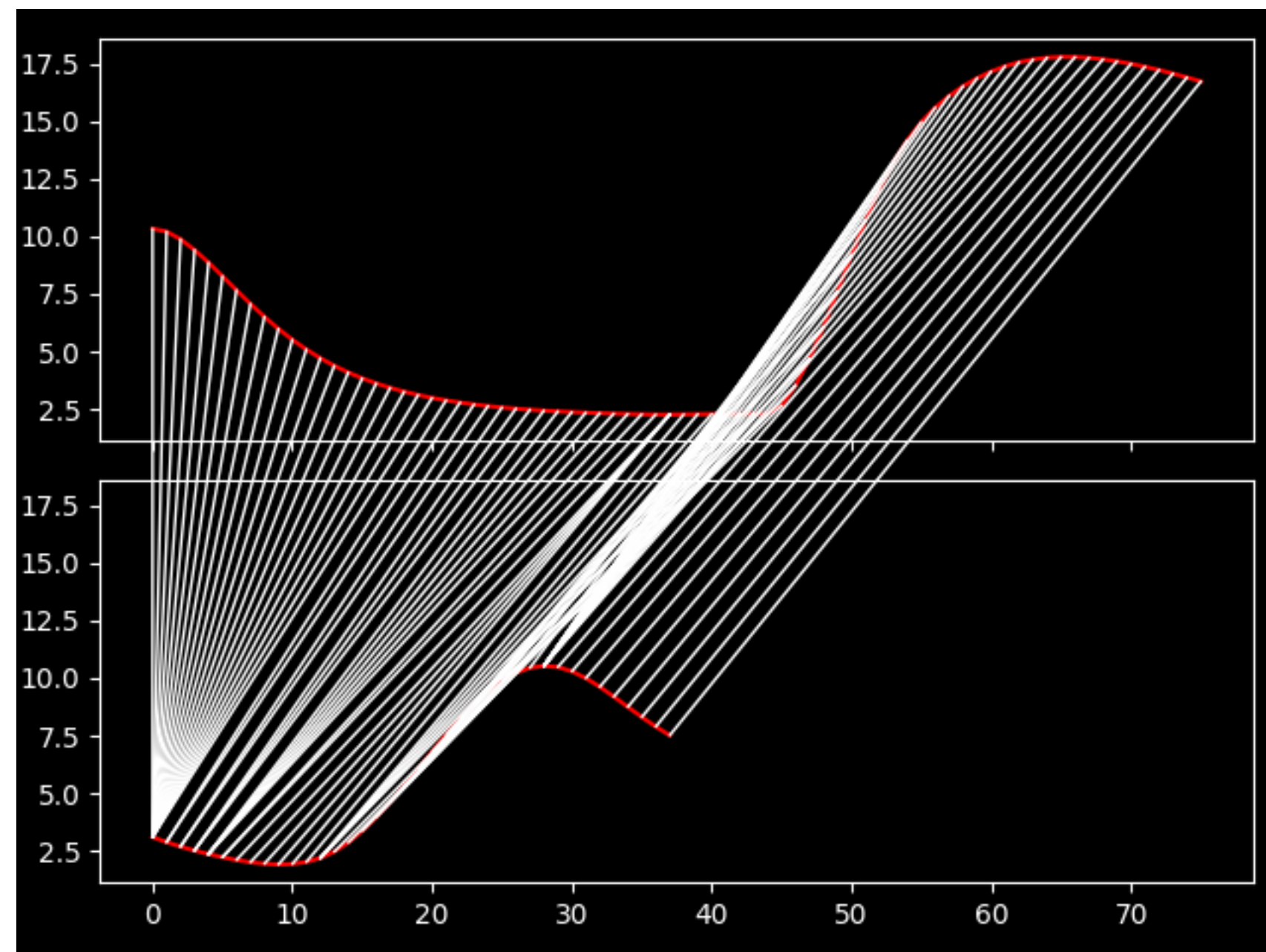
Simulated

Real

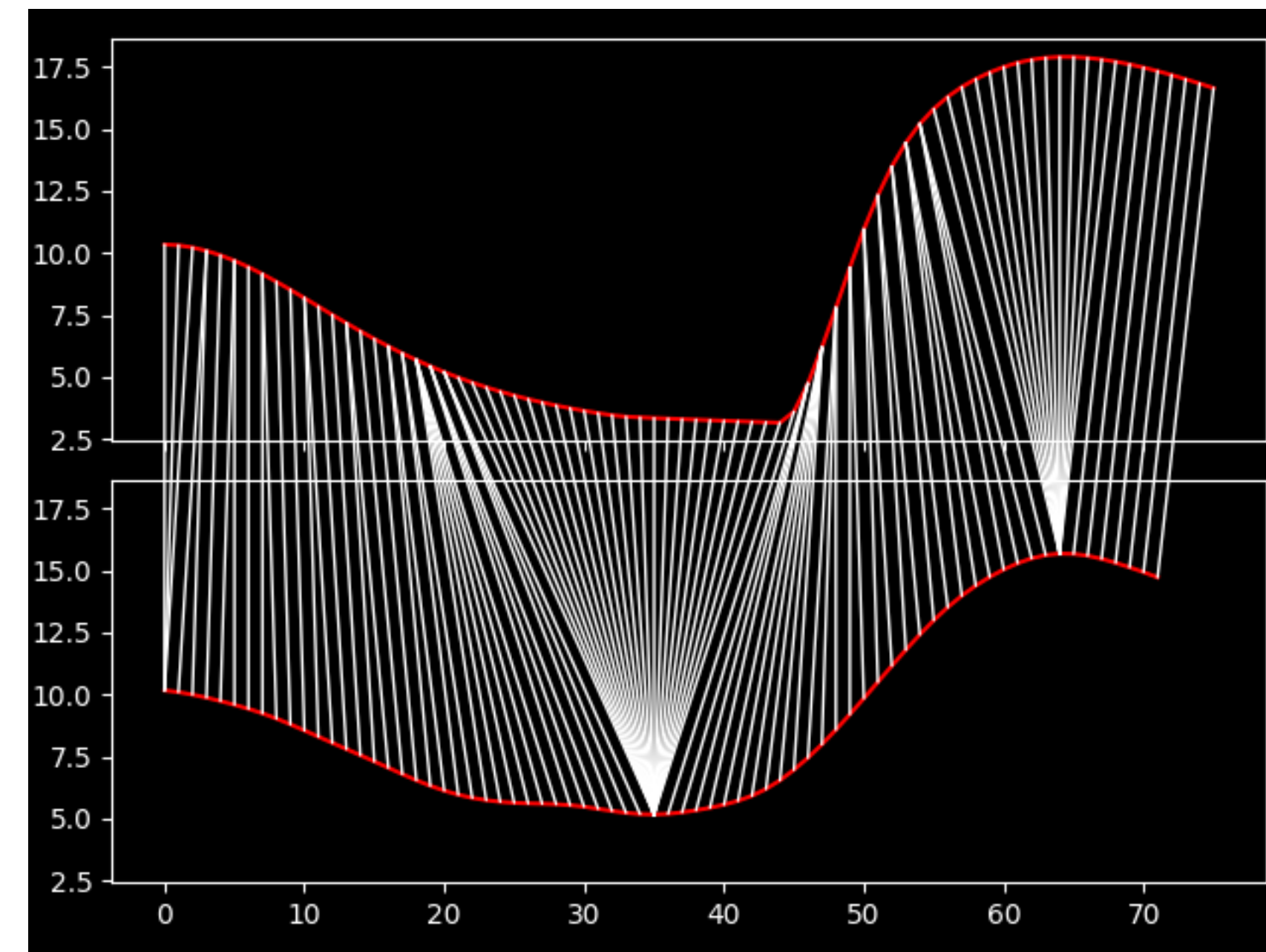
Analysis-by-synthesis: <five>

- Diphthong targets can't be separated with kinematic data
- Make a simulation, then tweak it, → 34,000 simulations
Compare to 525 tokens from X-ray Microbeam Database

Bad fit



Good fit



Simulated

Real

Interim findings

Analysis-by-Synthesis of <five>

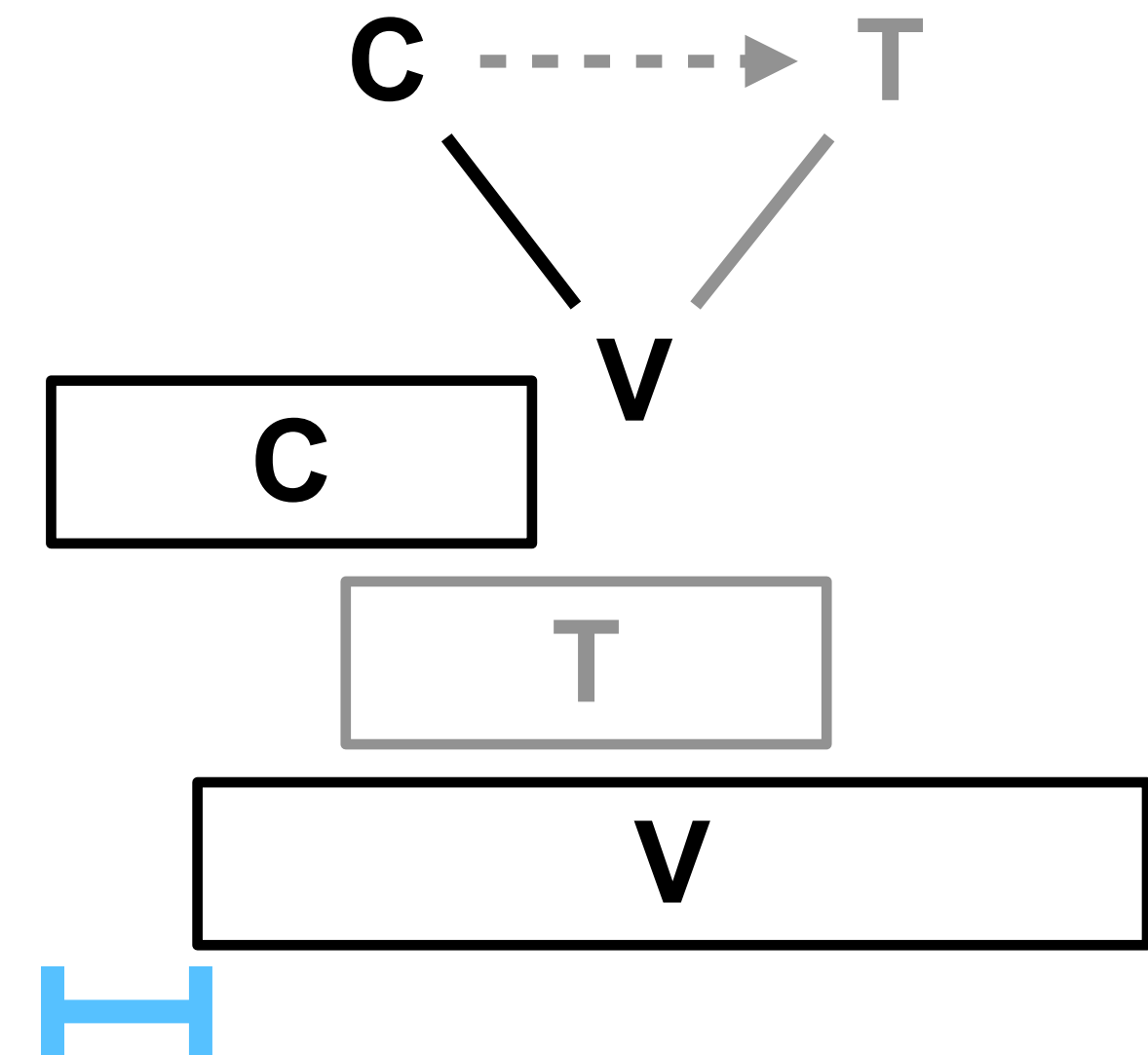
- We got some results!
 - [a] portion of diphthong timed to rest of word
 - [ɪ] portion more free to vary across tokens
- Still a lot to do
 - Extremely computationally-intensive
 - Which dimensions of variation? How much to vary?
 - What's the best way to compare curves?

Roadmap

- Phonology and articulatory gestures
- Coordinating gestures: the Coupled Oscillator Model
- Problems
 - Tibetan tone study
- Toward solutions: Analysis-by-synthesis
- **Conclusion**

What have we learned

- Tibetan: some speakers have tone, others don't
BUT they all have the same C-V lag
→ problem for Competitive Coupling of Tone
- Studying coordination requires consistent, reliable, practical ways to identify gestures
→ Analysis-by-Synthesis might help



Theory ↔ Data

- Observation: Gestures!
 - Theory: Oscillators!
- Observation: overlap in clusters
 - Theory: Coupled Oscillators!
What else can this do? Tone!
- Observations: or can it?
 - Theory: ...

Theory ↔ Data

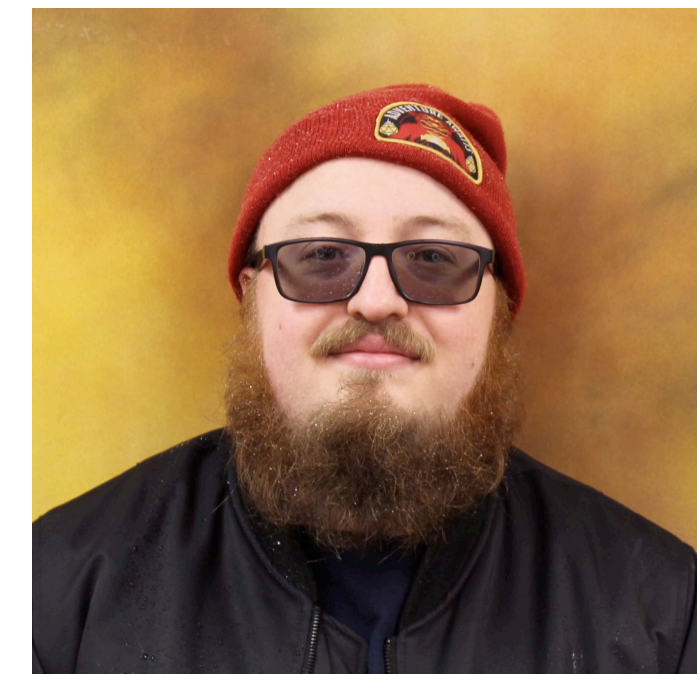
- Observation: Gestures!
- Theory: Oscillators!
- Observation: overlap in clusters
- Theory: Coupled Oscillators!
What else can this do? Tone!
- Observations: or can it?
- Theory: ...

For now:

- Gather new observations
& reevaluate old ones
→ descriptive generalizations
- Gather pieces that might help us with the next iteration
→ tools (e.g. simulators)
→ insights from other fields

If you liked this talk...

- Take classes!
 - CGSC 253: Philosophy of Cog Sci (Jonathan McKinney)
 - CGSC/PSYC 232: Cognitive Processes (Cathie Galotti)
 - LING 275: First Language Acquisition (Cati Fortin)
 - LING 318: Laboratory Phonology (Chris Geissler)
 - MUSC 110&204: Music Theory I&II (Justin London)
 - PSYC 216: Behavioral Neuroscience (Lawrence Wichlinski)
 - PSYC 220: Sensation & Perception (Julia Strand)
 - ... and more, incl. with Mija Van Der Wedge, Cherlon Ussery...
- Research! (Possibly this December?)
- Talk with linguists! (Snacks after...)



सुभासहेके

Thank you!

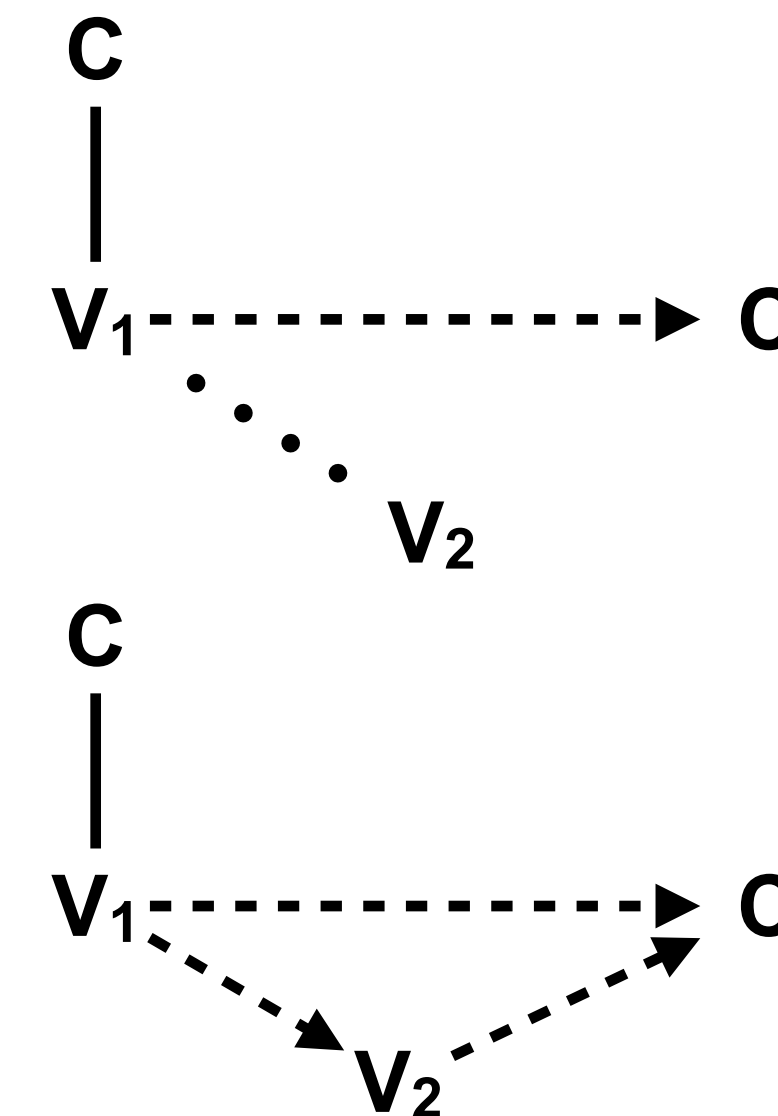
Pocket slides

What about diphthongs?

- Can approximately describe with in-phase/anti-phase
- How do diphthongs change when they get shorter?

< five > /faɪv/

LIPS	labiodent. critical	labiodent. critical
TONGUE TIP		
TONGUE BODY	pharyngeal wide	palatal narrow
VELUM		
GLOTTIS	wide	



Articulatory study

Geissler et al. (2021), Geissler (2021ch4)

- H1: variation in timing conditioned by presence/absence of lexical tone
 - speakers with tone contrast will have competitive coupling (pos. C-V lag)
 - speakers without tone contrast will have in-phase C-V timing (no C-V lag)
- H2: timing convergence:
 - all speakers will have similar coordination patterns despite interspeaker variation in presence/absence of tone
- What kind of tone contrast is there?
 - If H- \emptyset , then difference will be visible in high vs. low tone words
 - If H-L, then no difference in timing by tone.

EMA Study conclusions

- H1: variation in timing conditioned by presence/absence of lexical tone
 - speakers with tone contrast will have competitive coupling (pos. C-V lag)
 - speakers without tone contrast will have in-phase C-V timing (no C-V lag)
- **✓ H2: timing convergence:**
 - all speakers have similar coordination patterns despite interspeaker variation in presence/absence of tone
- What kind of tone contrast is there?
 - If H- \emptyset , then difference will be visible in high vs. low tone words
 - **✓ If H-L, then no difference in timing by tone.**

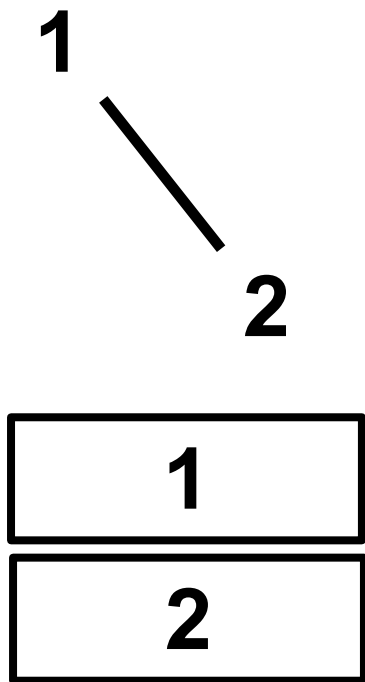
The temporal basis of complex segments

Shaw et al. 2019

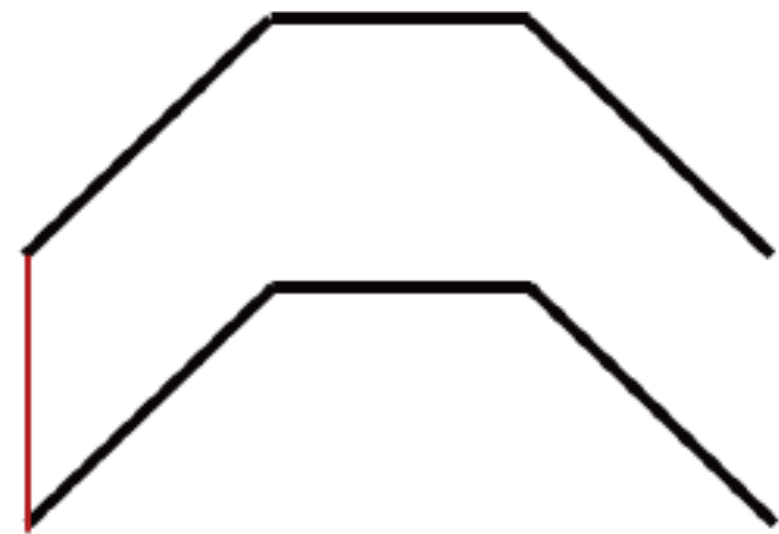
The temporal basis of complex segments

Shaw (2019): predictions

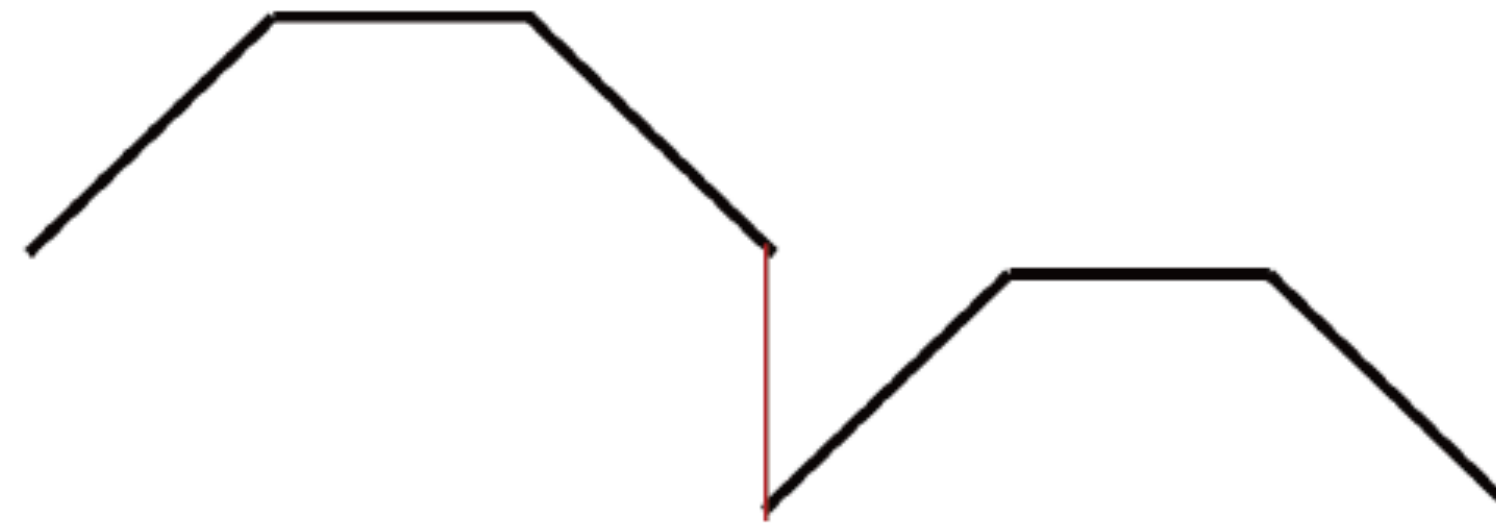
In-phase



(a) Complex segment—no lag



(b) Segment sequence—no lag

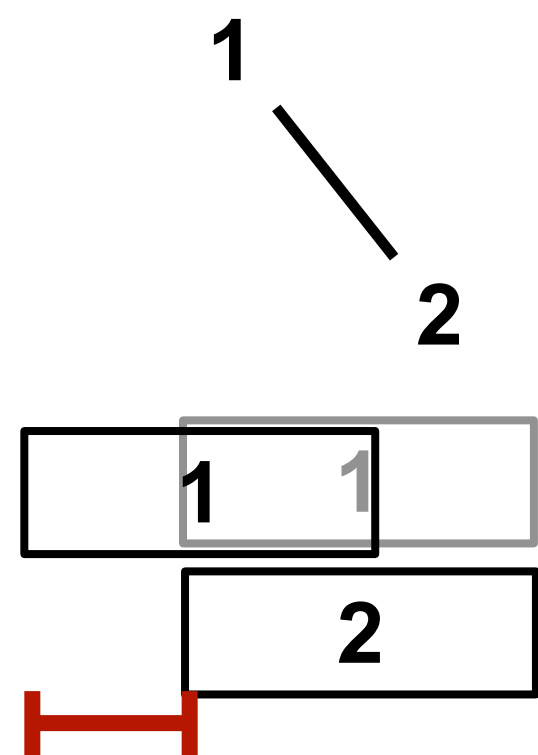


Anti-Phase

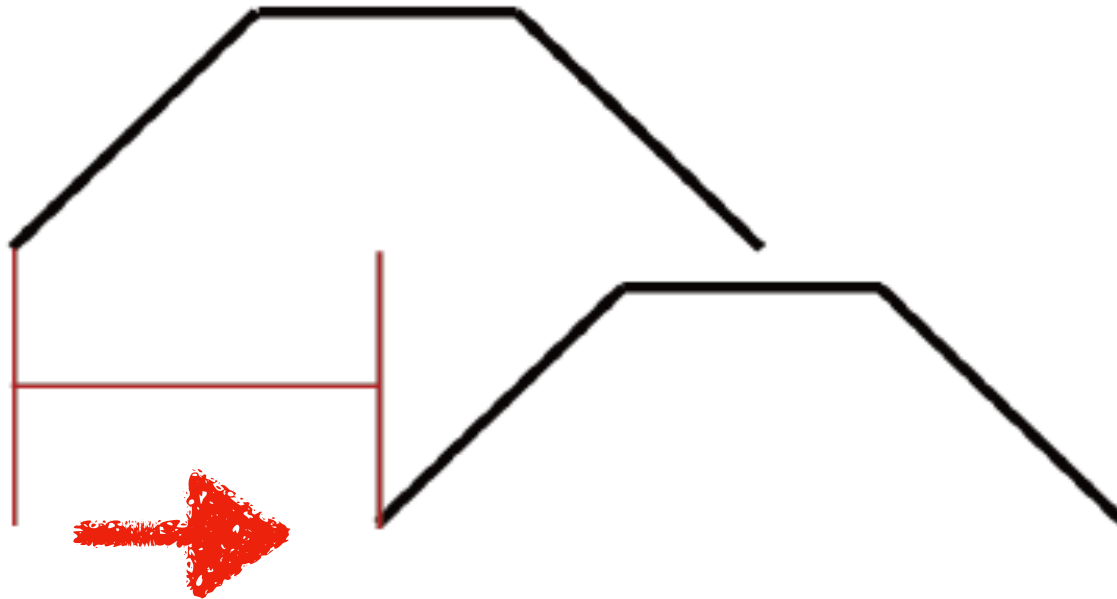
1→ 2



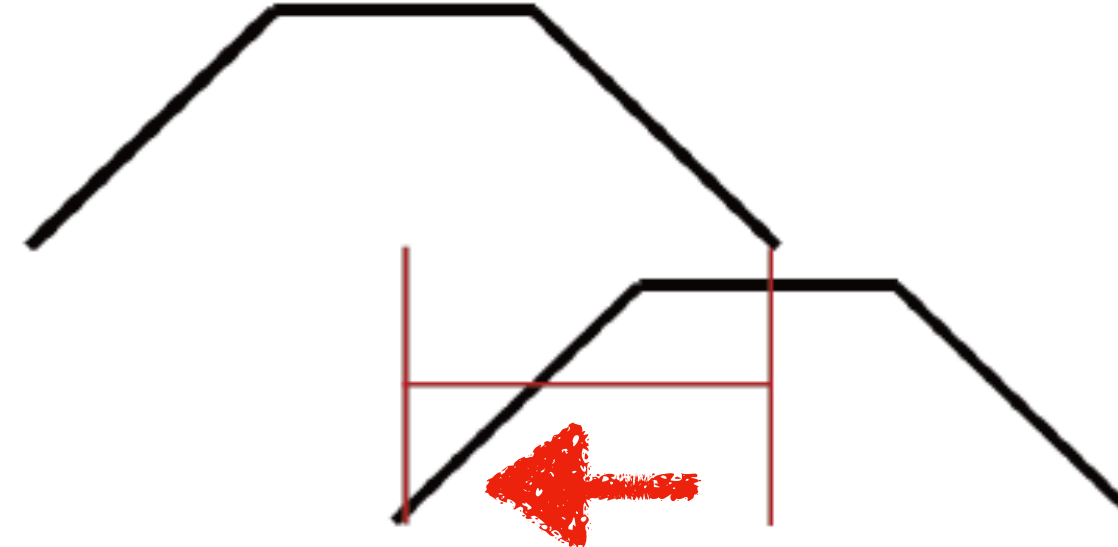
In-phase + lag
(offset)



(c) Complex segment—positive lag



(d) Segment sequence—negative lag



Anti-Phase - lag
(offset)

1→ 2

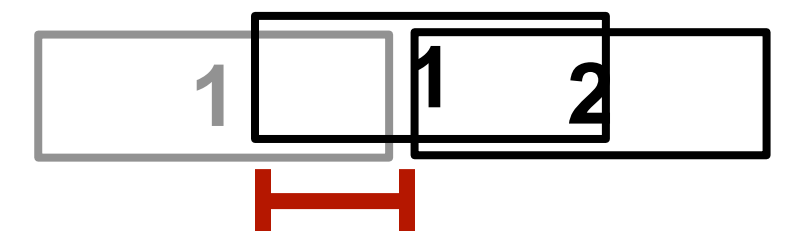


Figure 1: Hypothesized gestural coordination patterns for complex segments (a), (c) and segment sequences (b), (d)

The temporal basis of complex segments

Shaw (2019): results

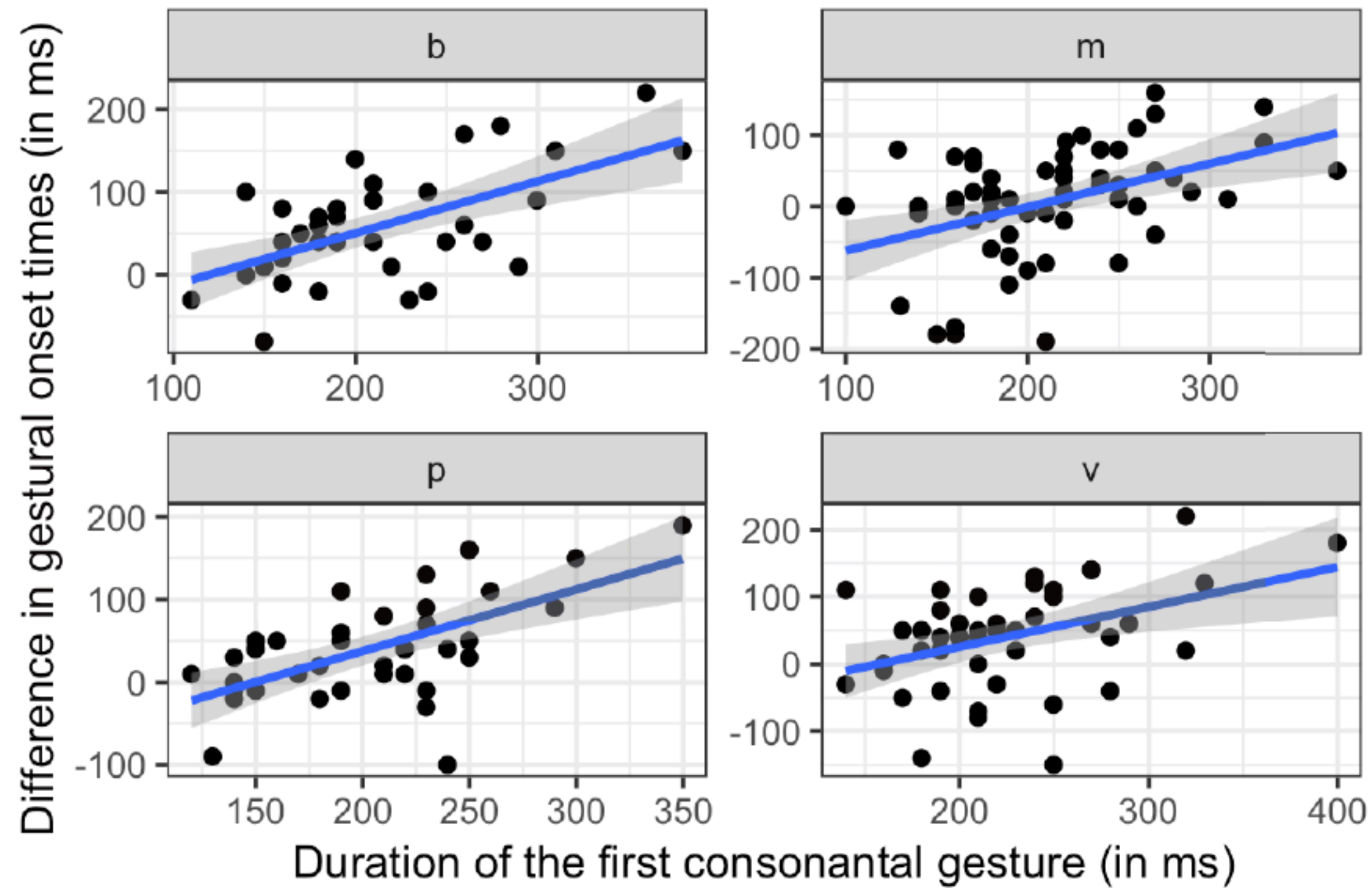


Figure 4: Correlations for the data from the English experiment

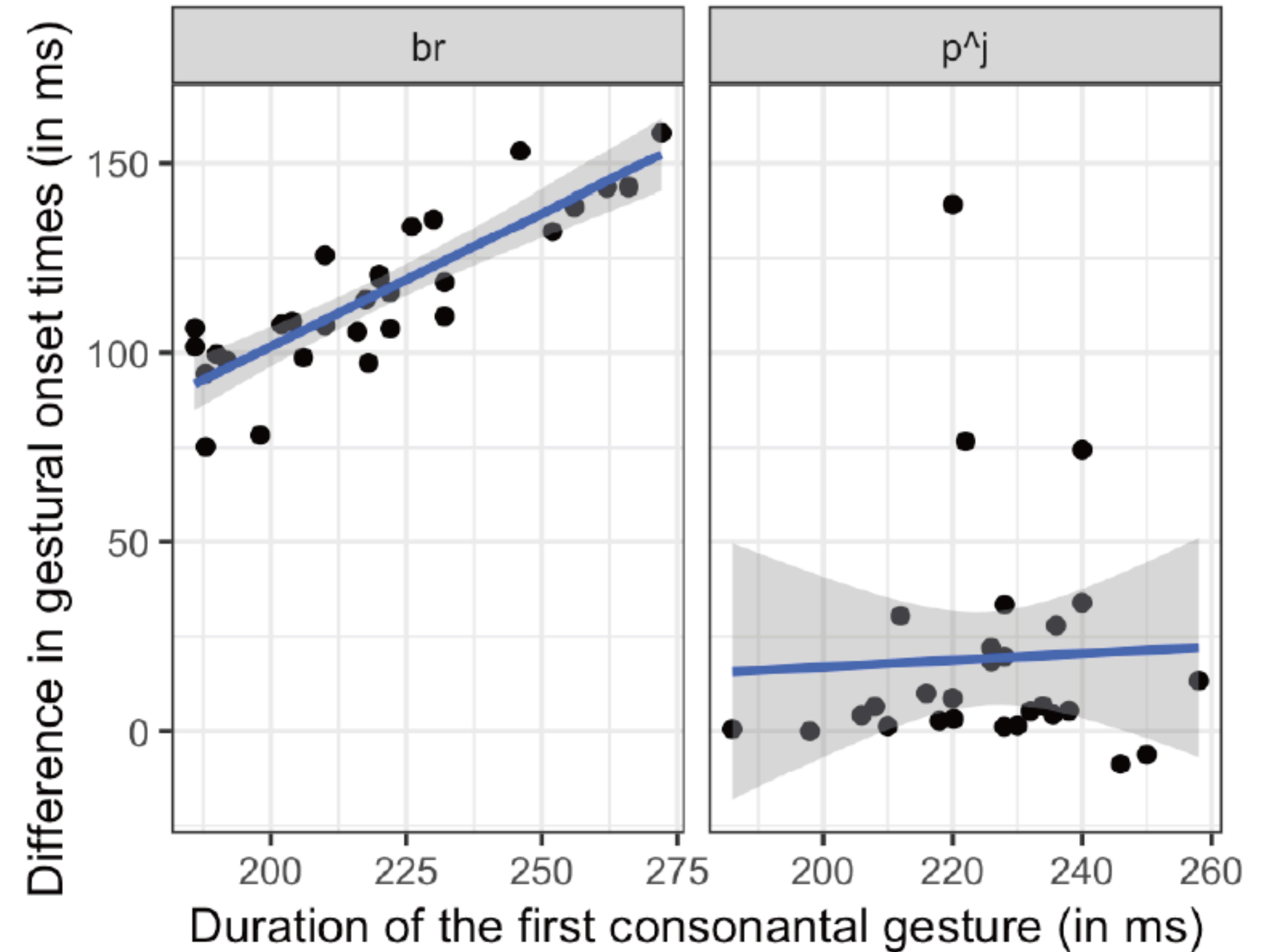


Figure 2: Correlations for the Russian data

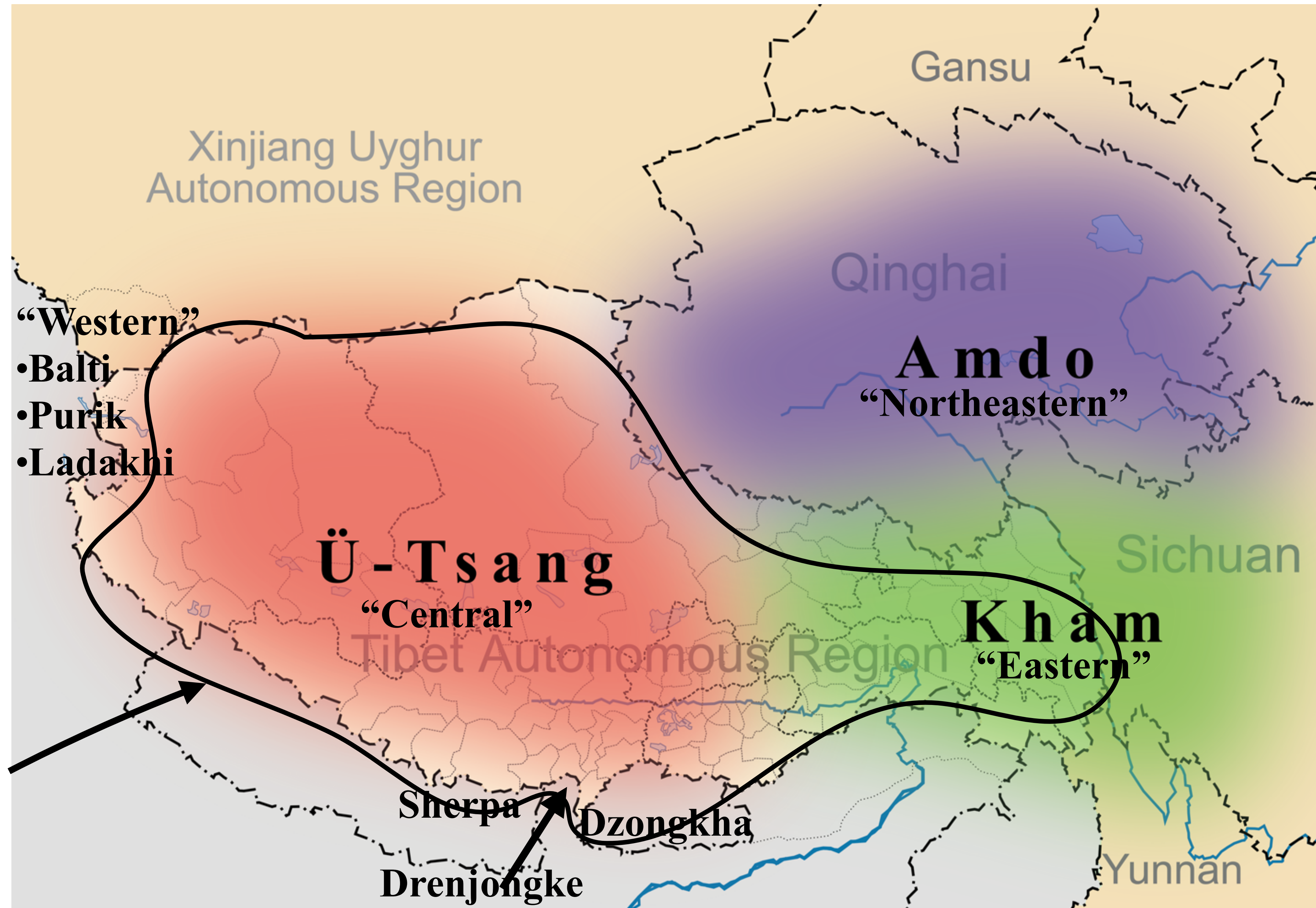
Tibetan dialects

Tibetan

བོད་སྐད་

- “archaic” / “cluster”
- “innovative” / “non-cluster”
- dialect continuum
- post-1959 diaspora

Approx.
extent of
tone



Dialects: Natural laboratory

- tonogenesis
- laryngeal variation
- cluster simplification
- vowel shifts, spirantization, retroflexion, palatalization
- evidential, honorifics, modality, etc.

Written (Classical) Tibetan	Balti (Western)	Rebkong (Northeastern)	Tokpe Gola (Central)	Gloss
<i>khrag</i>	[kʂʌk]	[t̪ɕʁɣ]	[tʰʌk] ([tʰák])	‘blood’
<i>rtswa</i>	[xstsoa]	[xtsa]	[tsá]	‘grass’
<i>spyang ki</i>	[spjaŋ. 'ku]	[xt̪ɕaŋ. 'kʰɣ]	[t̪ʂáŋ.gú]	‘wolf’
<i>bcu bdun</i>	[t̪ɕub. 'dun]	[t̪ɕɣb. 'dɣn]	[t̪ʂúp.tũ] ([t̪ʂúp.tý])	‘seventeen’

(Adapted from Caplow 2013)

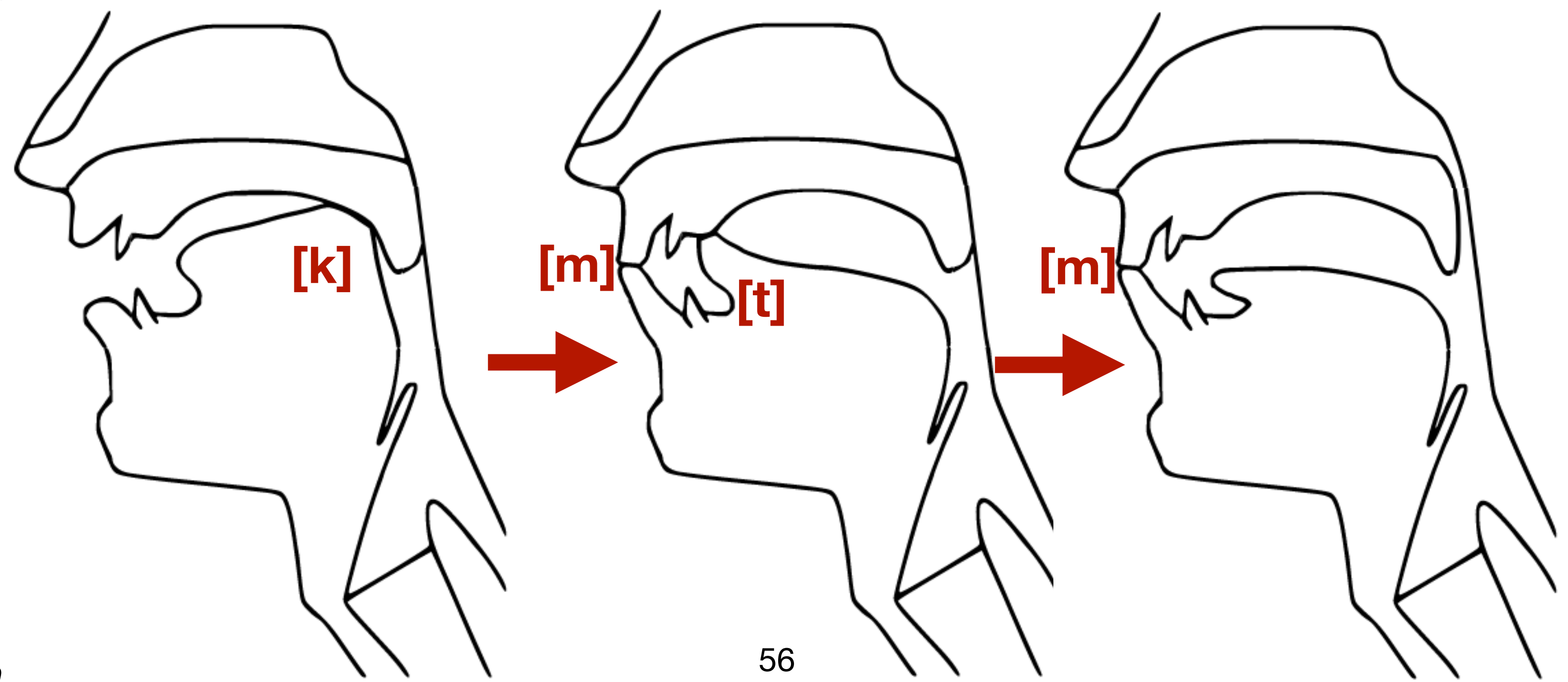
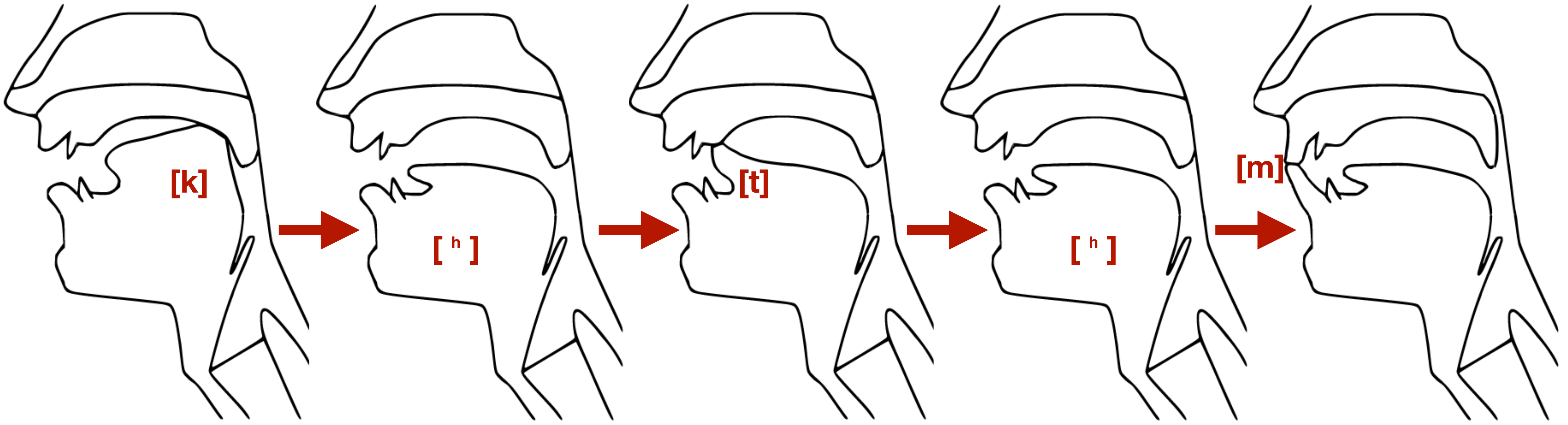
Tonogenesis

(tonal dialects only)

- Voiceless onsets > high tone
- Voiced onsets > low tone
- Sonorants with pre-initial > high tone
- *p^har ‘over there’ > H
*sa ‘earth’ > H
- *bar ‘between’ > L
*za ‘eat’ > L
*mar ‘butter’ > L
- *sman ‘medicine’ > H

Laryngeal contrasts

	Etymological onsets							Innovative features
Orthography	ཕ་	ཕ་	བ་	ཕ་	ཕ་	ཟ་	བཟ་	
Old Tibetan	s ^ə pa	p ^h a	ba	s ^ə ba	sa	za	b ^ə za	aspiration allphonic
Northeastern and Western dialects	spa	p ^h a	ba ~ wa	ɣba	sa	za	za	cluster simplification aspirated/unaspirated contrast
Eastern dialects	pá	p ^h á	pà	bà	sá	zà	zà	tonogenesis cluster simplification
Central dialects (Lhasa)	pá	p ^h á	p ^h à	pà	sá	sà	sà	voiced clusters > voiceless voiced simplex > aspirated

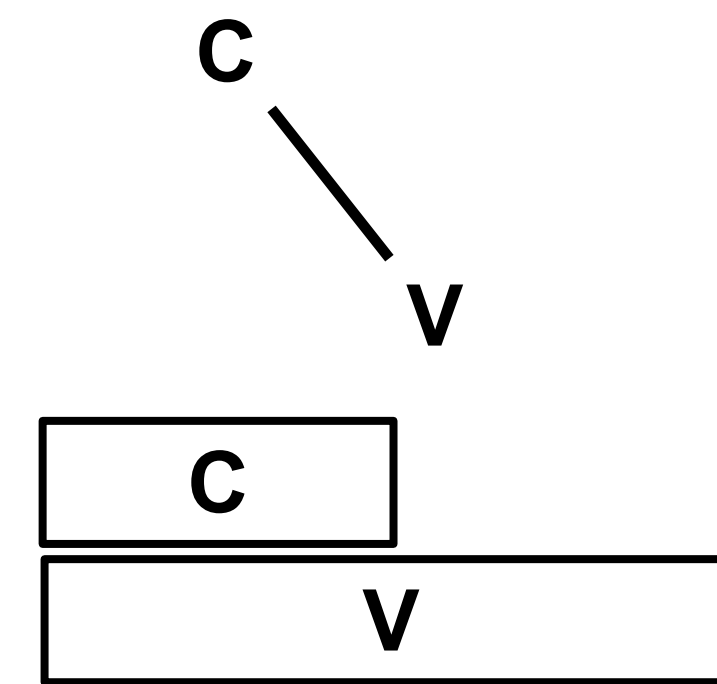


[back to slide 7](#)

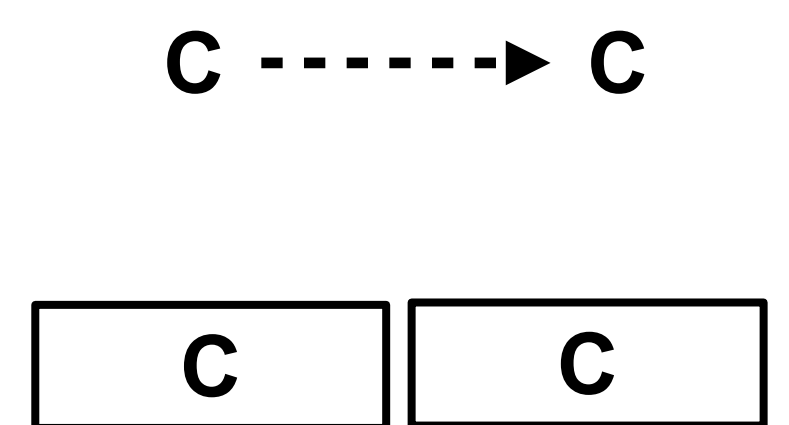
Coordinating gestures in time

- Gestural coupling modes:
 - *In-phase coupling*: (synchronous) and *Anti-phase coupling* (sequential) are most stable
 - *Competitive coupling*: combination of in-phase and anti-phase coupling relations
 - *Eccentric coupling*: one coupling relation, just not intrinsically stable

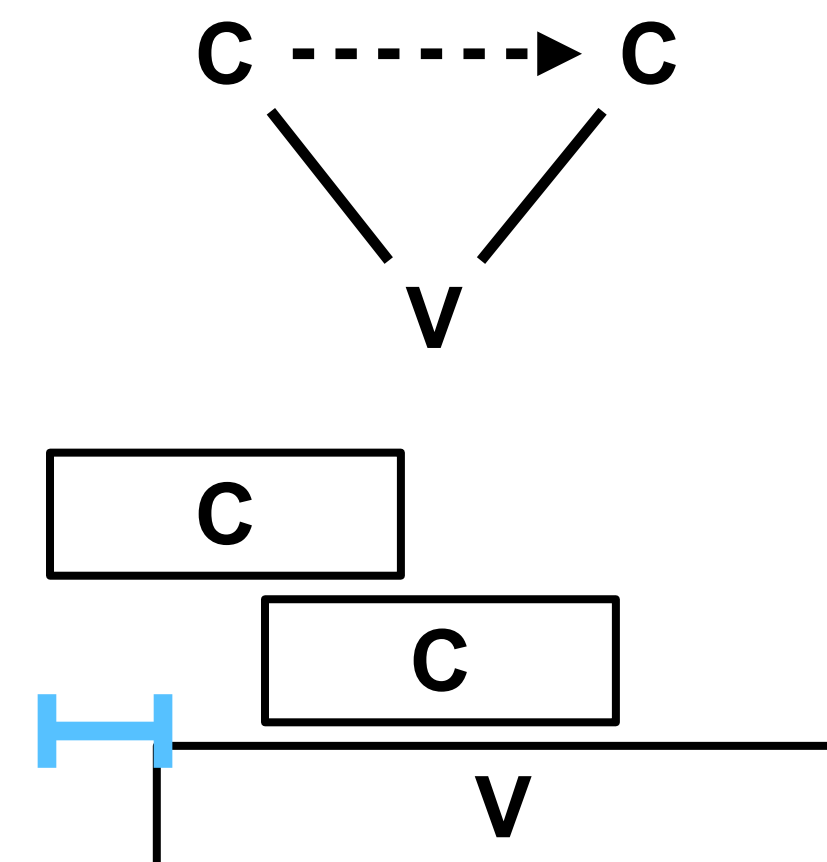
In-phase



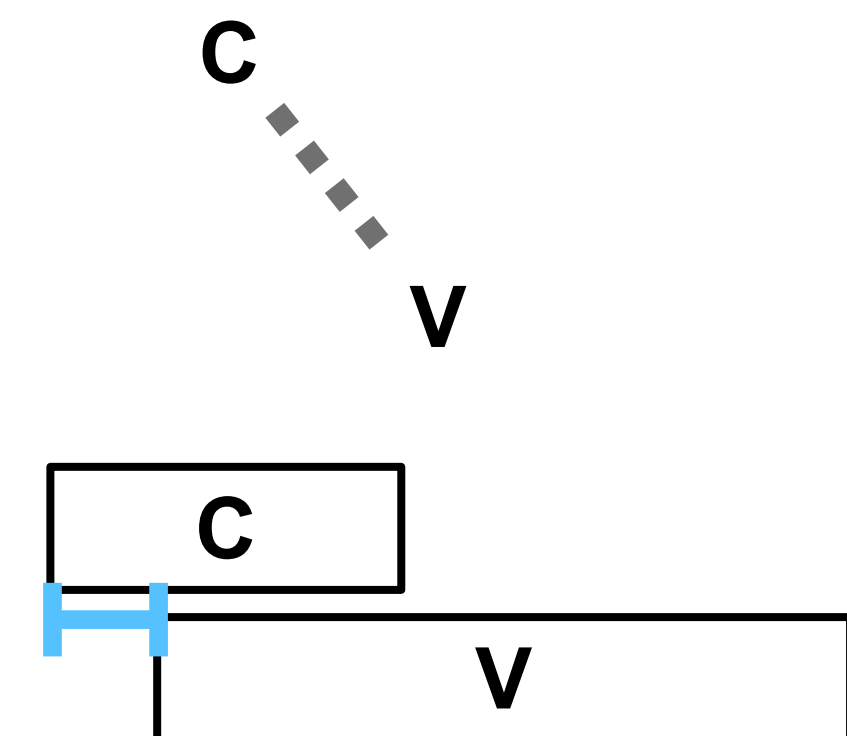
Anti-Phase



Competitive



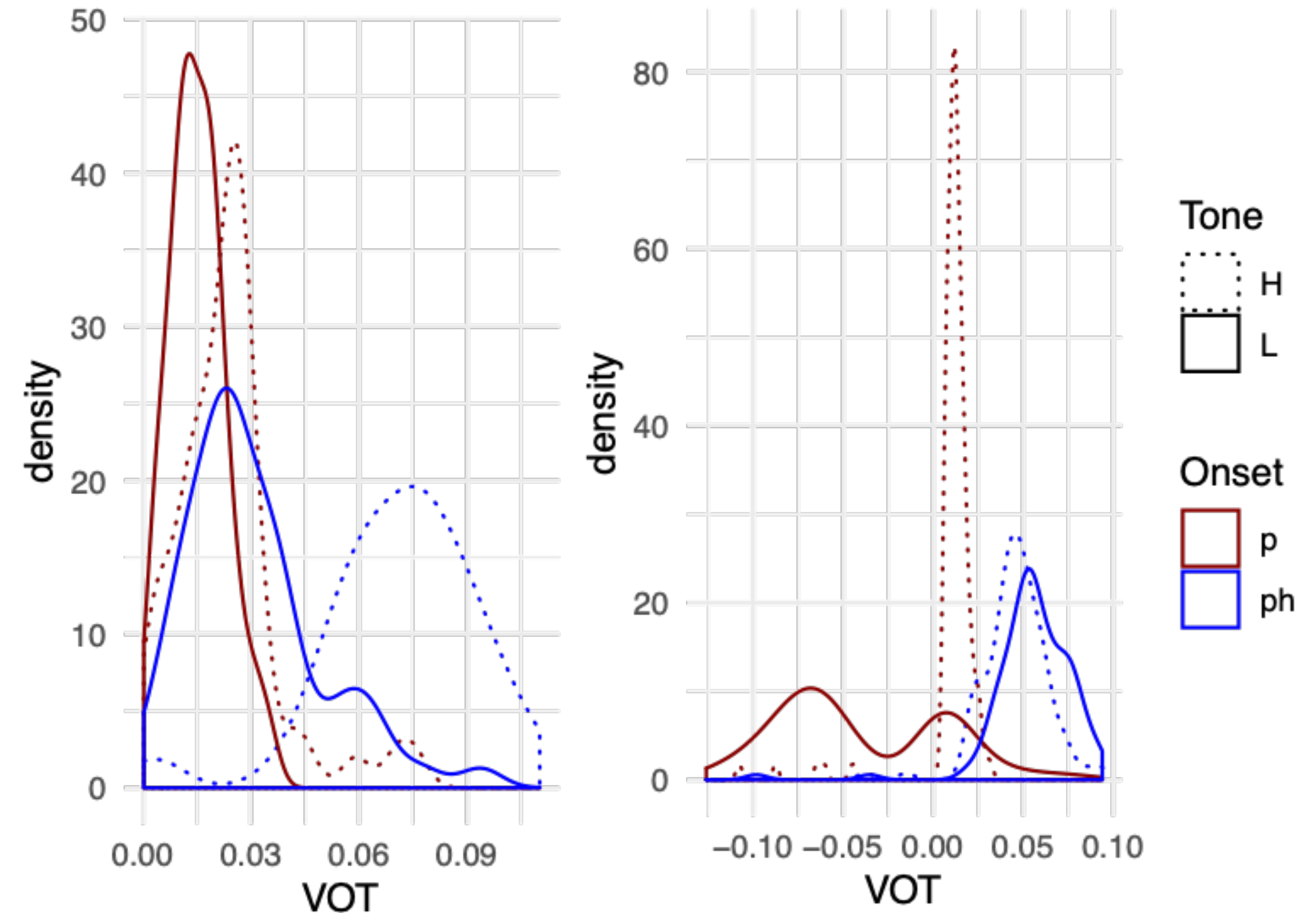
Eccentric



Two systems of laryngeal contrasts

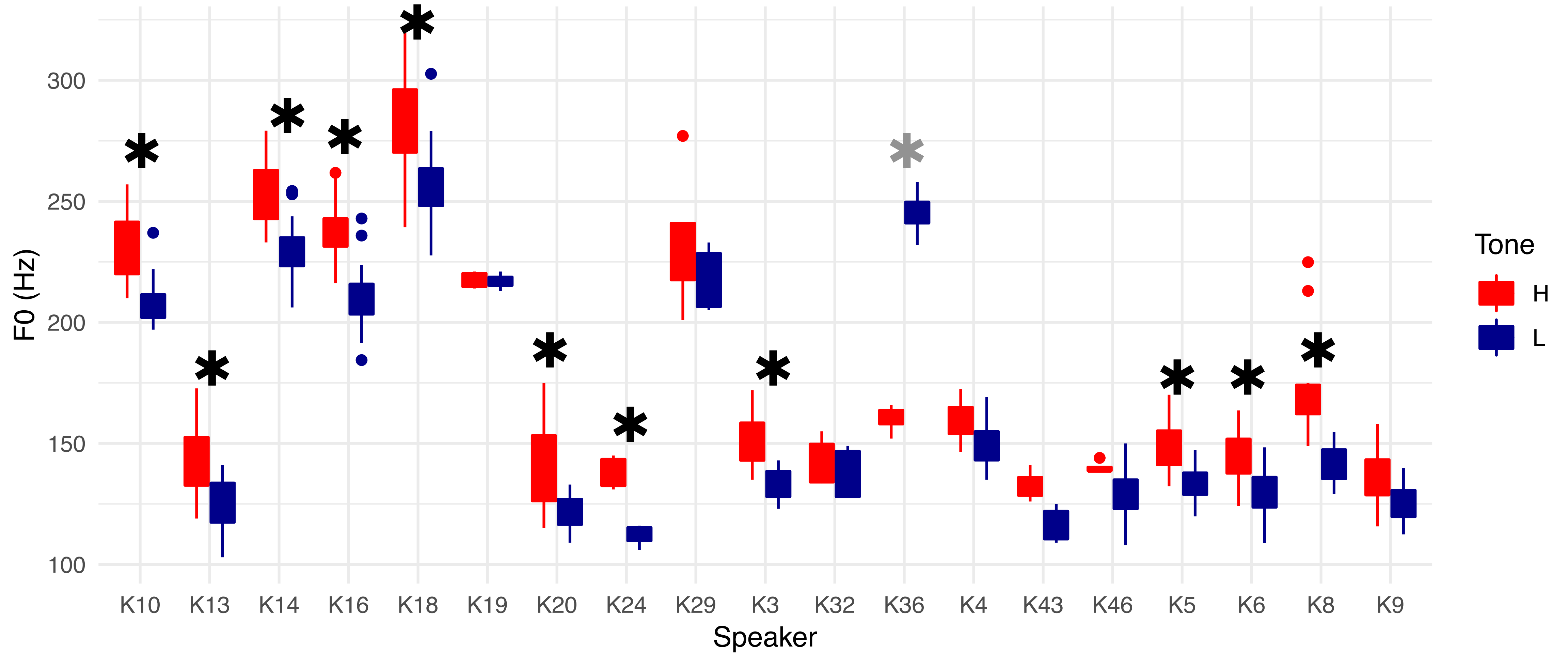
even in speakers with no F0 contrast (!!!)

- Both conditioned by etymological tone category:
- Left speaker
 - no prevoicing
 - long VOT only with H tone
- Right speaker:
 - prevoicing with L tone
 - long VOT with both tones

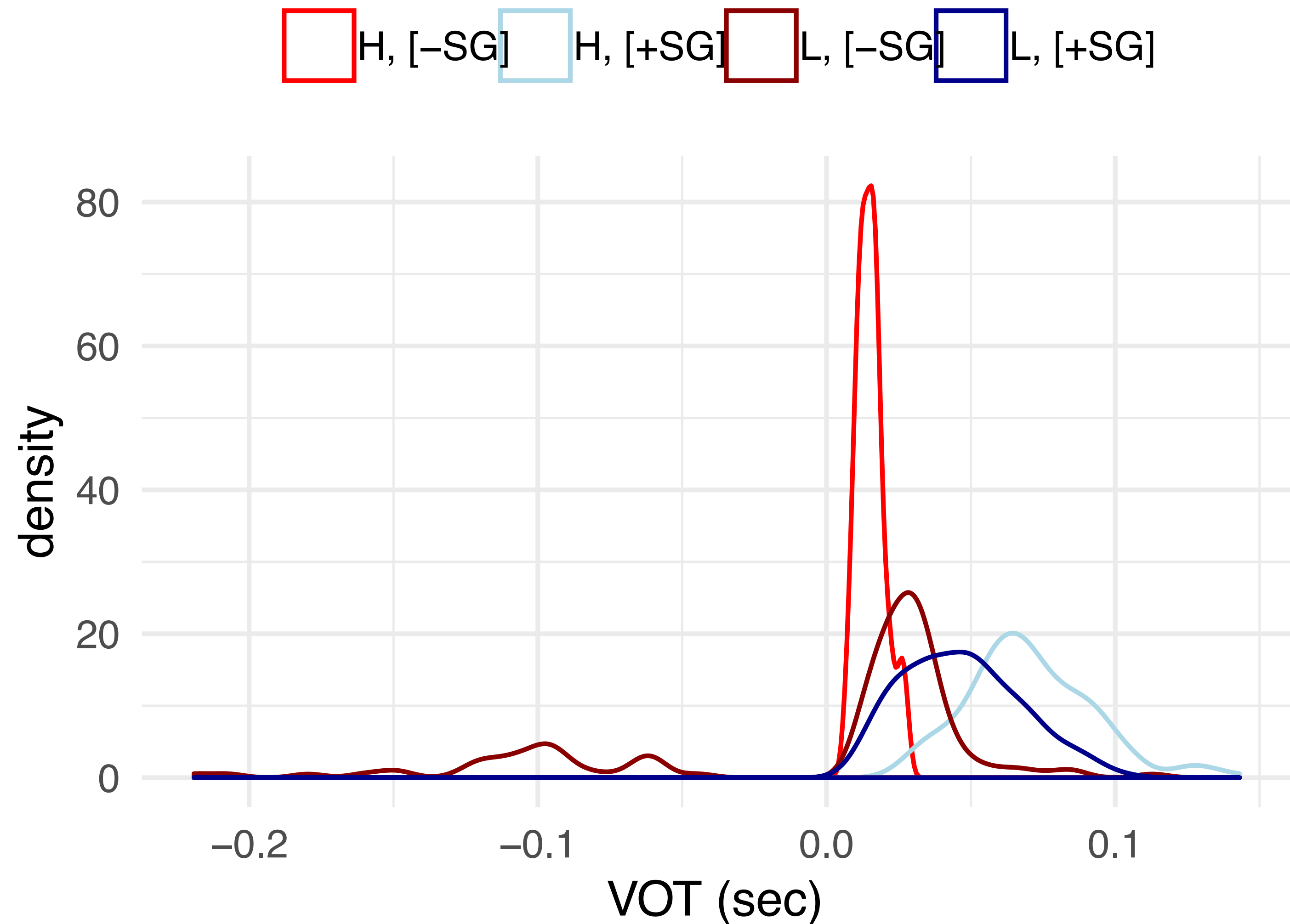


Does H have higher pitch than L?

Yes for 11/19, no for 7/19



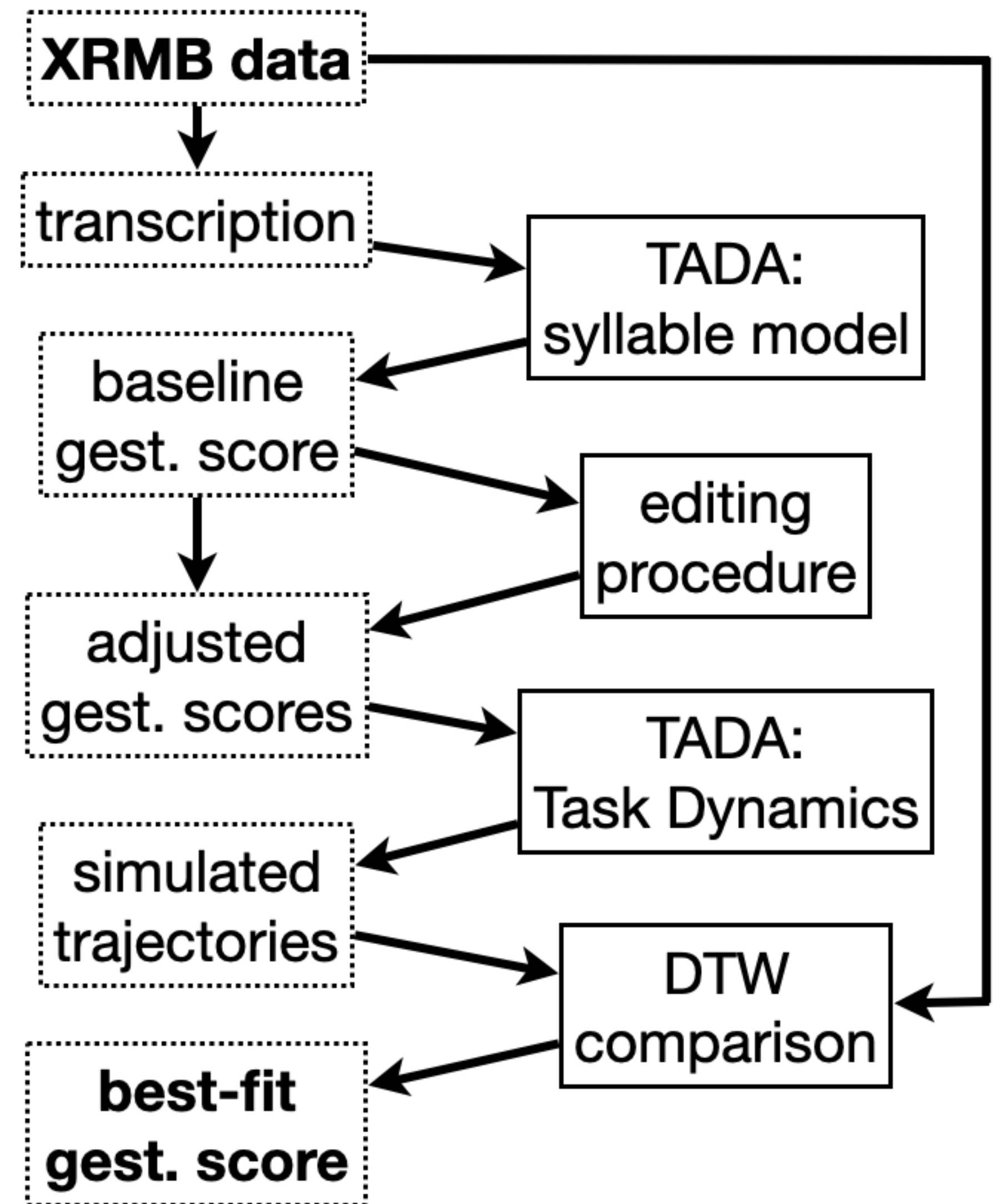
Consonant and tone categories



<five> study: methods

O'Reilly, Geissler, & Tang (2023)

- Ideal test case?
 - diphthongs: all four modes
 - C's with lips, V's with tongue
 - available data



Timing in phonology and/or phonetics?

- “Discrete Phonology” vs. “Gradient Phonetics”
- Speech timing as phonology
 - Is timing *intrinsic* or *extrinsic* to phonology?
 - Are gestures coordinated at *beginning* or *end*?
 - *Symbolic* vs. *phonetically-enriched* representations?