# Towards large(r)-scale cross-linguistic study of speech:
## prosodic case studies

*Morgan Sonderegger*

McGill University

Northwestern University
6/1/2018

# Introduction

- Speech is highly variable

- Structure and sources of variability
  - Central Qs in linguistics/speech sciences
  - Decades of work $\rightarrow$ much known
  - Scale: mostly handful of cues (VOT, formants), languages (English), hand measurement

- Most of what we know is from fine-grained studies
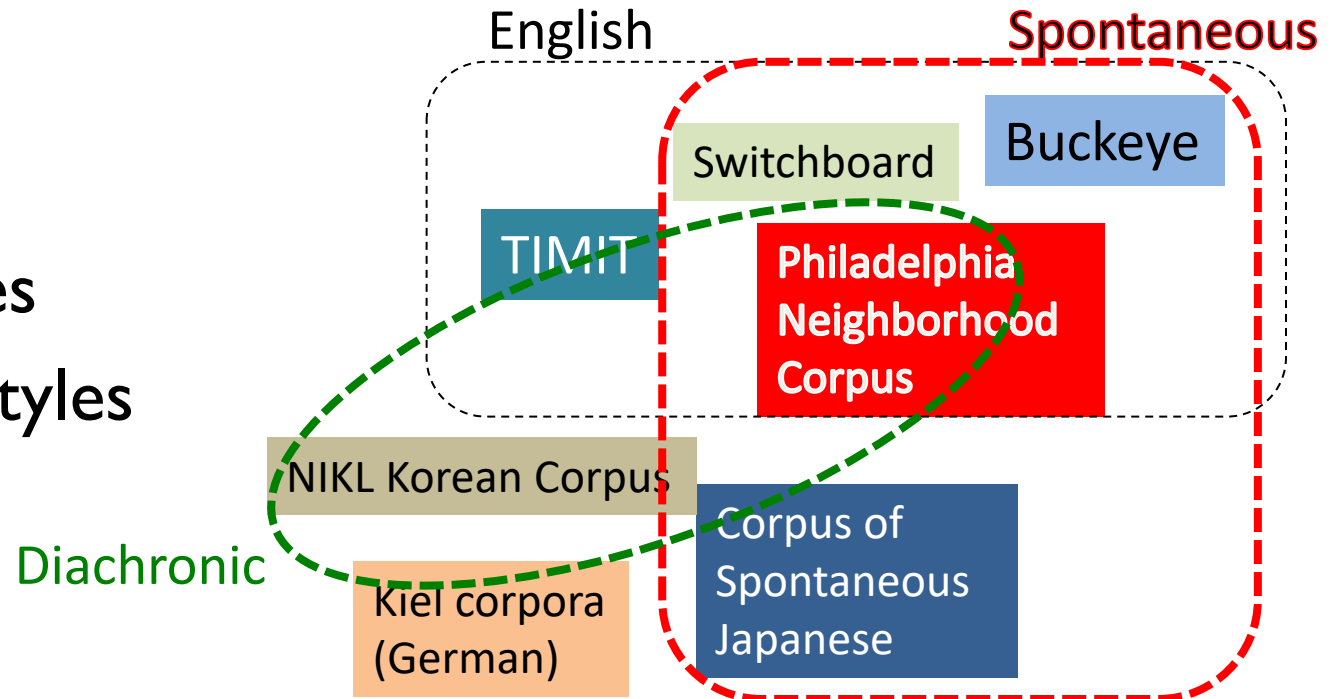
# Introduction

- Huge amount of annotated speech data exists
  - Corpora
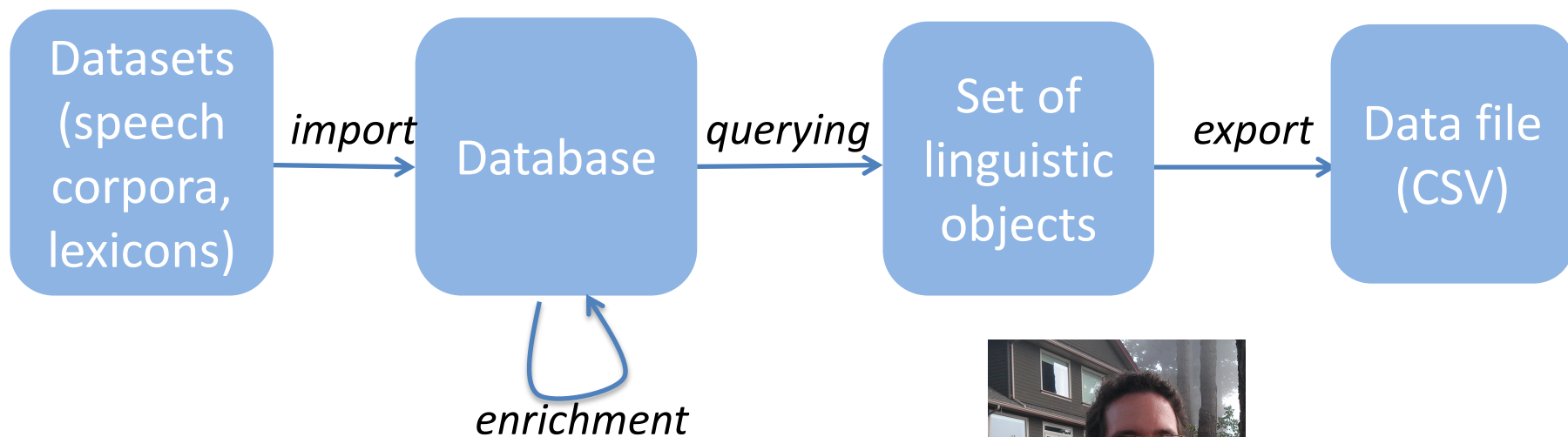  - Academic labs
  - Web

At least orthography + audio

- Across
  - Languages
  - Speech styles
  - Time

English          Spontaneous

Switchboard          Buckeye

TIMIT

Philadelphia Neighborhood Corpus

NIKL Korean Corpus

Diachronic

Kiel corpora (German)
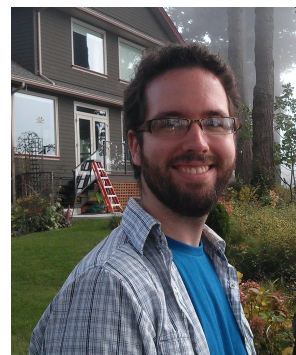
Corpus of Spontaneous Japanese

# Introduction

- Large(r)-scale studies
  - corpora + (semi) automatic analysis + statistical modeling
  - Scale up
  - Less careful

- Today: two case studies

- Enabled by software facilitating large-scale studies
- Claims:
  - New insights
  - Complementary to fine-grained studies

# Polyglot-Speech Corpus Tools

Datasets (speech corpora, lexicons) → *import* → Database → *querying* → Set of linguistic objects → *export* → Data file (CSV)

*enrichment*

- Implementation
  - Python module
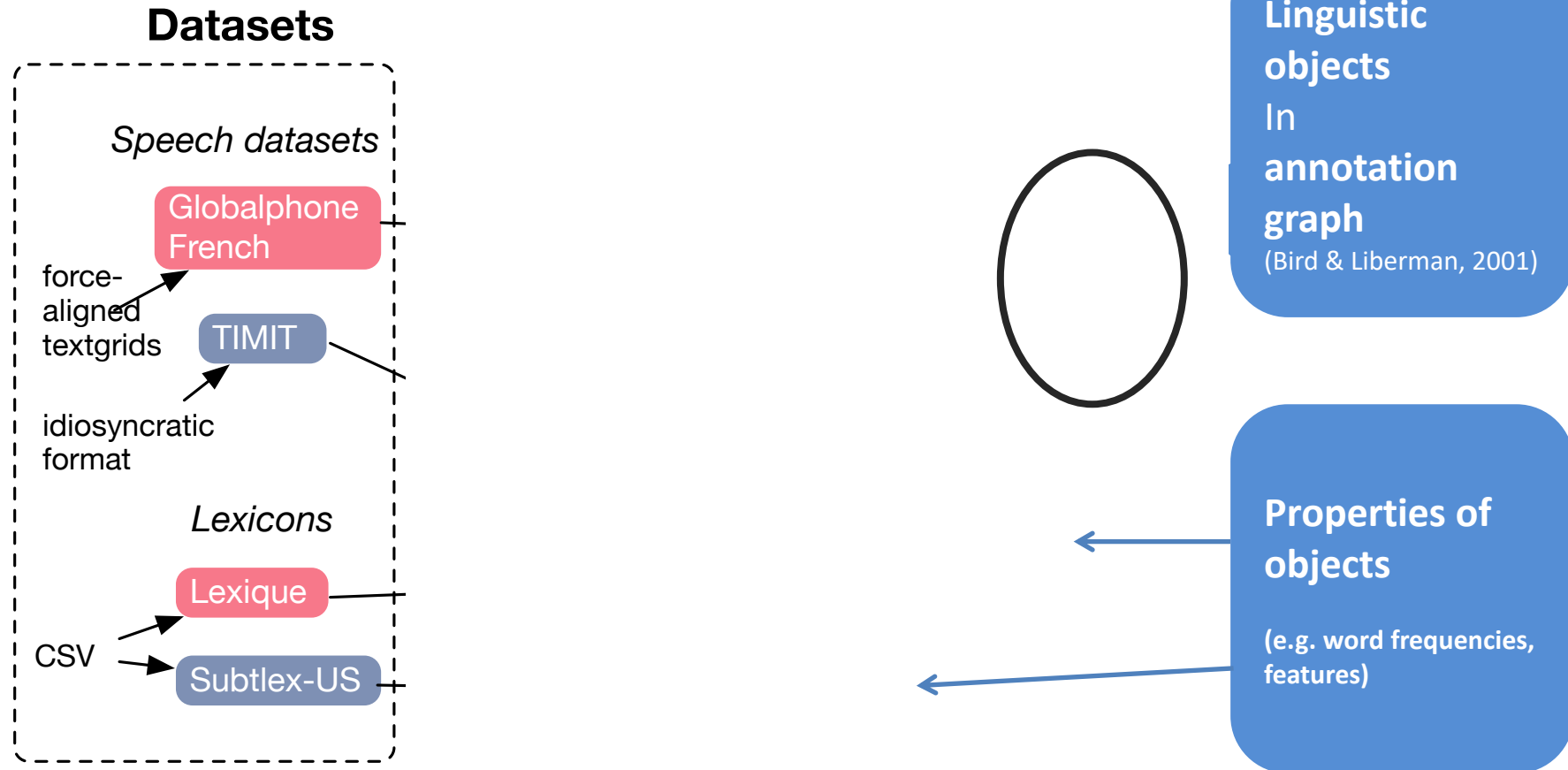  - Graphical interface (under redevelopment)

McAuliffe et al. (2017) *Interspeech*

montrealcorpustools.github.io/speechcorpustools/
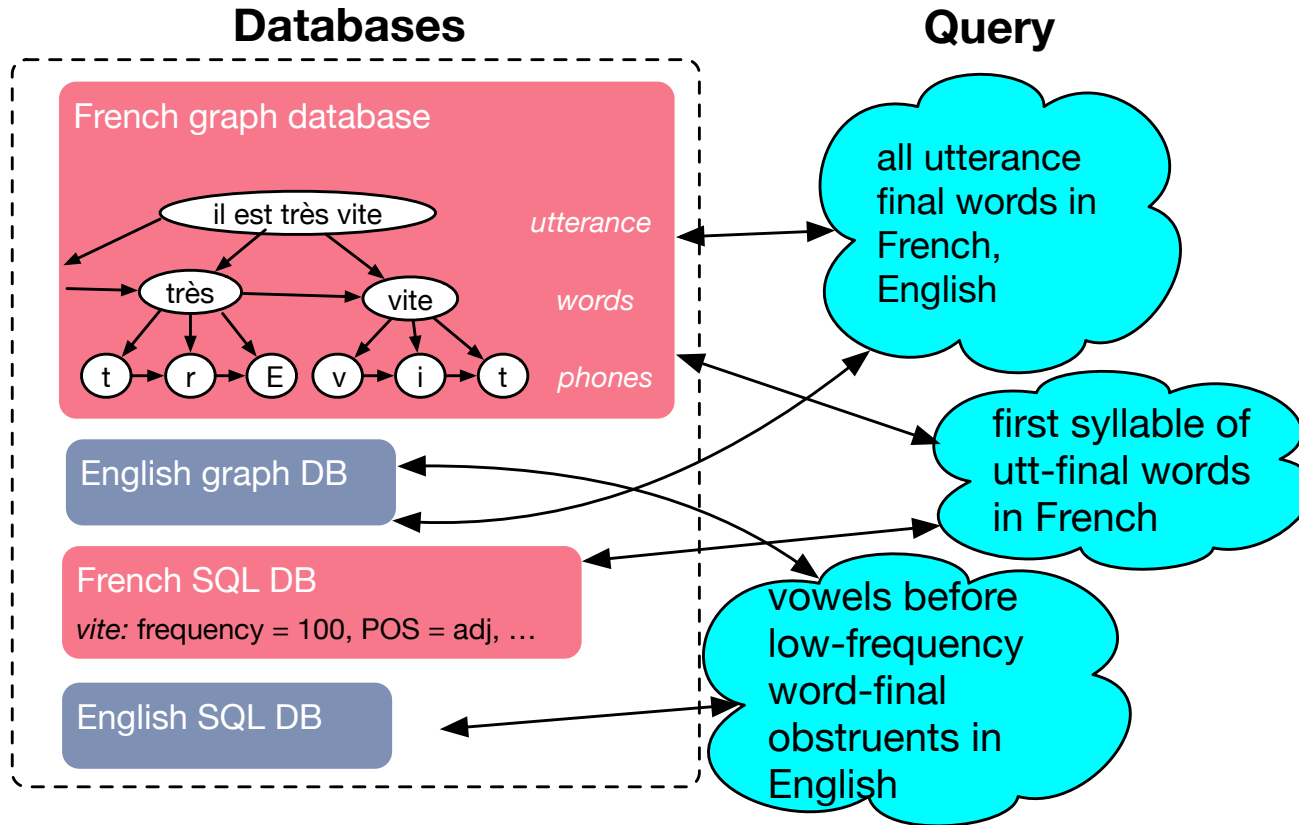
# Polyglot-SCT: Goals

1. Scalable

2. Require minimal technical skill from user

3. Abstraction away from dataset format

4. Querying dataset without access to raw data

• Aim to address barriers to large-scale corpus studies

# Polyglot-SCT: Import

**Datasets**

*Speech datasets*

Globalphone French

force-aligned textgrids

TIMIT

idiosyncratic format

*Lexicons*

Lexique

CSV

Subtlex-US

**Linguistic objects**
In **annotation graph**
(Bird & Liberman, 2001)

**Properties of objects**

(e.g. word frequencies, features)

- Speech, text datasets → queryable databases

# Polyglot-SCT: query

**Databases**

**Query**

French graph database

il est très vite    *utterance*

très     vite    *words*

t r E v i t    *phones*

English graph DB

French SQL DB

*vite:* frequency = 100, POS = adj, …

English SQL DB

all utterance final words in French, English

first syllable of utt-final words in French

vowels before low-frequency word-final obstruents in English

- Find subset of linguistic objects

# Polyglot-SCT: export



- Properties of objects → spreadsheet
  - (→ R, Excel)

# Study 1: intrinsic F0 effects

# Introduction

- Where does sound change come from?
- Most common:

  phonetic effect   →   phonological pattern

  "phonetic precursors"

- Ex: tones
  – Often (e.g. Chinese):
    F0 perturbations   →   lexical tone

pá [ — ]        pá [ — ]        pá [ — ]
bá [ ⌐ ]        bǐ [ ⌐ ]        pǎ [ ⌐ ]

# Introduction

But: most phonetic precursors never lead to sound change!

What kind of precursor can be a source of change?

- robust
  - Across speakers, languages
- ... but variable
  - Individual differences, language-specific phonetics

tension

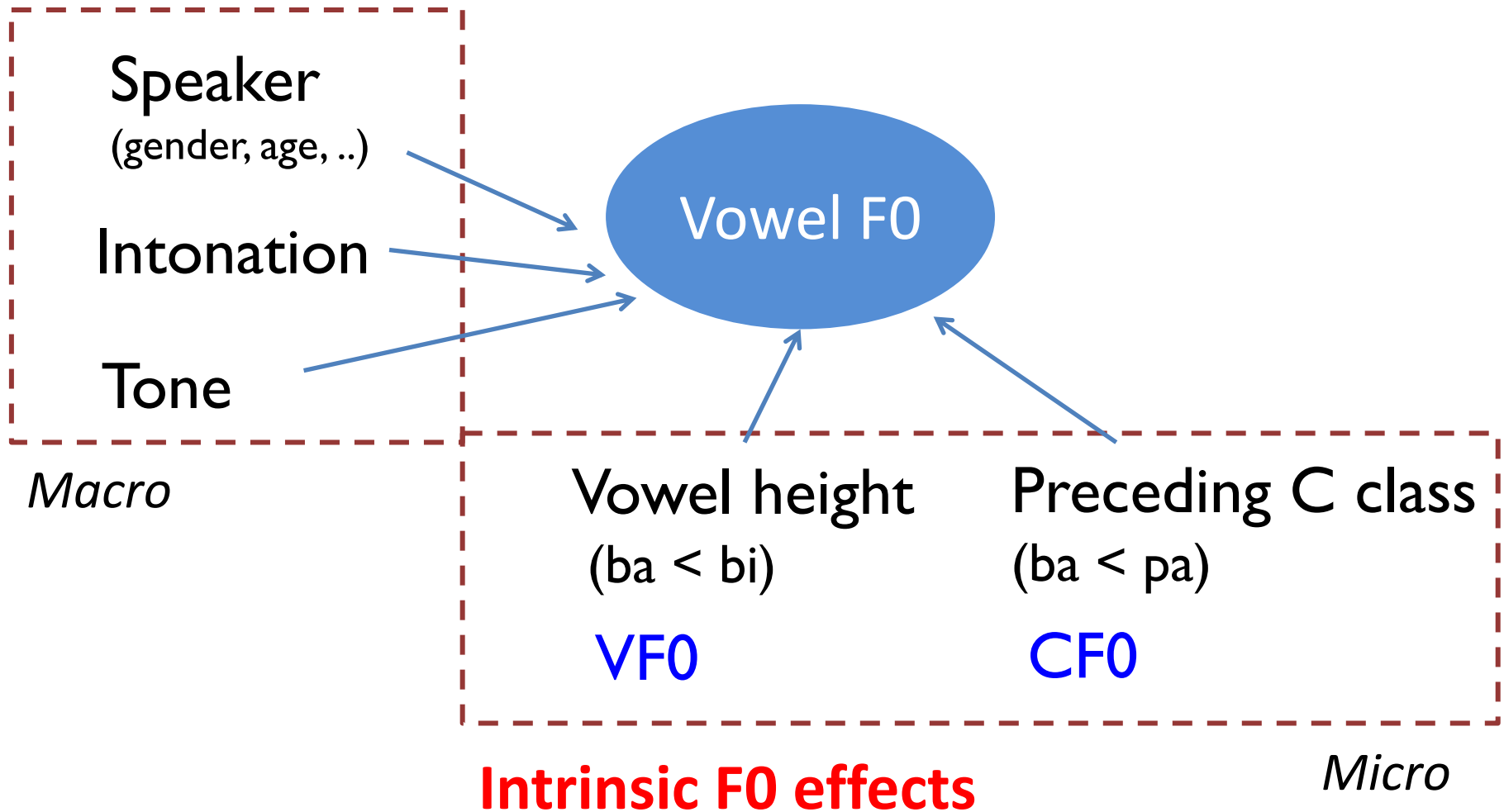(e.g. Hombert et al. 1979, Ohala 19XX; Baker et al., 2011; Labov, 1967; Kingston, 2007; Yu, 2013)

# Introduction

How robust/variable is each phonetic precursor, across languages and individuals?

# Introduction

- Methodologically hard
  - Need: big and comparable data: many languages, speakers
  - small effects, big confounds

- Approach:
cross-linguistic corpora + automatic analysis + statistical modeling

- Q1: can a "phonetic precursor" be detected in corpus data across languages & speakers?

# Influences on vowel F0



Speaker
(gender, age, ..)

Intonation

Tone

*Macro*

Vowel F0

Vowel height
(ba < bi)

VF0

Preceding C class
(ba < pa)

CF0

**Intrinsic F0 effects**

*Micro*

(e.g. Chen, 2011; Connell 2002; Fischer-Jørgenson, 1990; Hanson, 2009; Hoole & Honda, 2011; House & Fairbanks, 1953; Kingston & Diehl, 1994; Kirby & Ladd, 2016; Kingston, 2007; Ladd & Silverman, 1994; Meyer, 1896; Whalen & Levitt, 1995)

# Intrinsic F0

- Huge literature
  - primarily: small *n*, lab speech
  - focus: mechanism (automatic vs. controlled)

Across languages:

- CF0
  - "voiced"<"voiceless": most languages
- VF0
  - [-high] < [+high] : (near-)universal

- Effect size: variable
  - Tonal ⇒ smaller effect?

Q2: How much variability in IF0 across 14 languages?

# Intrinsic F0

- Strongly affected by:
  - "Intonation"
  - Gender (VF0)
  
  ...

- Interspeaker variability:
  - Often noted

- Relationship to sound change:
  - CF0 $\Rightarrow$ sound change ("tonogenesis")
  - VF0 $\neq$ sound change
  - Why?

Q3: How much variability in IF0 across speakers?

# Datasets

| | |
|---|---|
| English | Russian |
| French | Polish |
| German | Spanish |
| Korean | Turkish |

| |
|---|
| Hausa |
| Mandarin |
| Thai |
| Vietnamese |

- Read sentence corpora
  - ~20 hours each
  - Force-aligned

**Montreal Forced Aligner**: trainable for different languages

| |
|---|
| Swedish |

GlobalPhone (Schulz et al., 2013), Librispeech (Panyatov et al., 2015)

# Datasets

- "Utterance-initial"     C V

obstruent     /a/, /i/, /u/

- vowel F0  (Praat)
  – F0 histogram → speaker min, max → re-extract F0

- Controls : info about
  – Speaker
  – Utterance
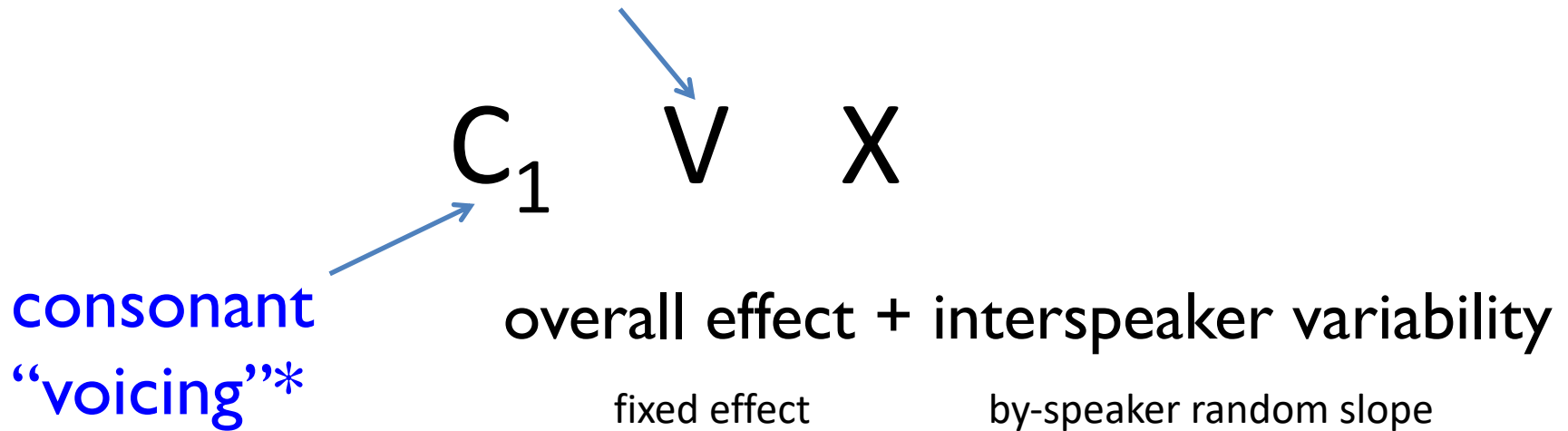  – Context
  – Word

Polygot-
Speech Corpus
Tools

# Datasets

- Data cleaning: minimize F0 errors, reduced vowels

- Exclusions:
  - "bad" speakers
  - "bad" tokens (e.g. too short)

- Data per language:
  - 1.9-9.5k tokens (~2000)
  - ~100 speakers

# CF0: Analysis

- One linear mixed effects model / language
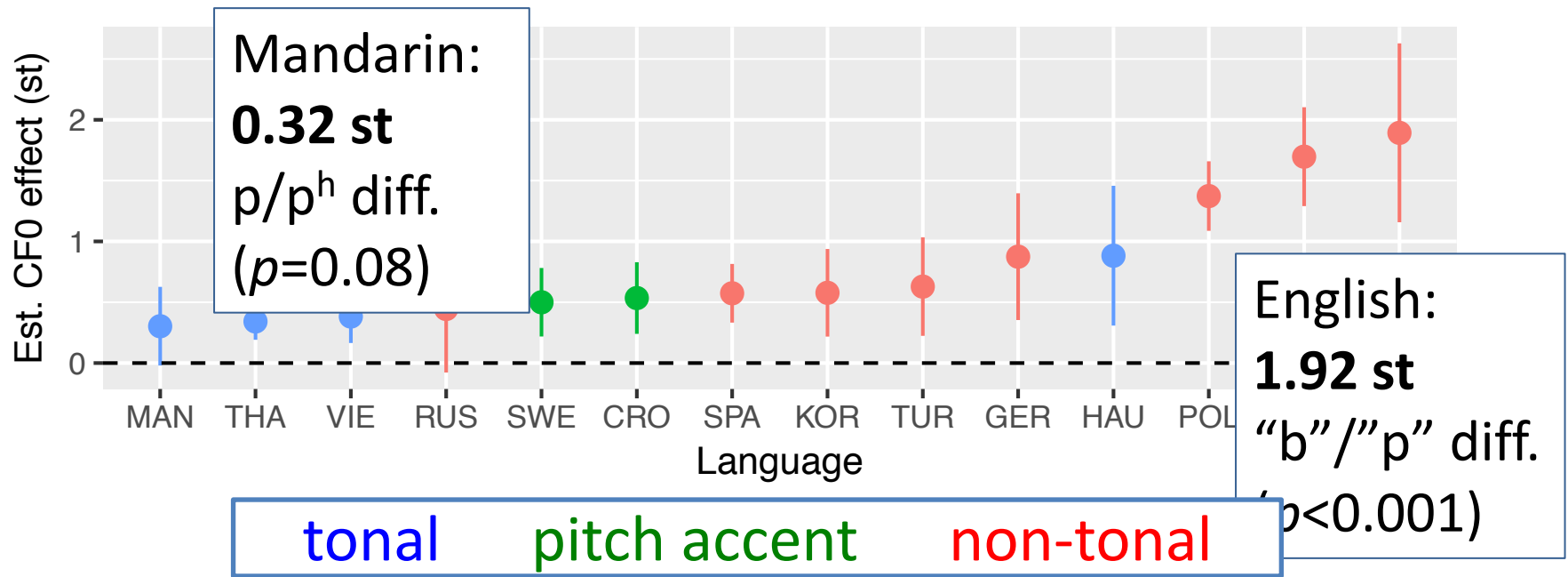- Main terms:

Response: <u>mean F0 in first 50 ms</u>

$$C_1 \quad V \quad X$$

consonant
"voicing"*

overall effect + interspeaker variability

fixed effect      by-speaker random slope

* Ex: French p/b, Mandarin p/p[h]

# CF0: analysis

- Other terms: extra slides
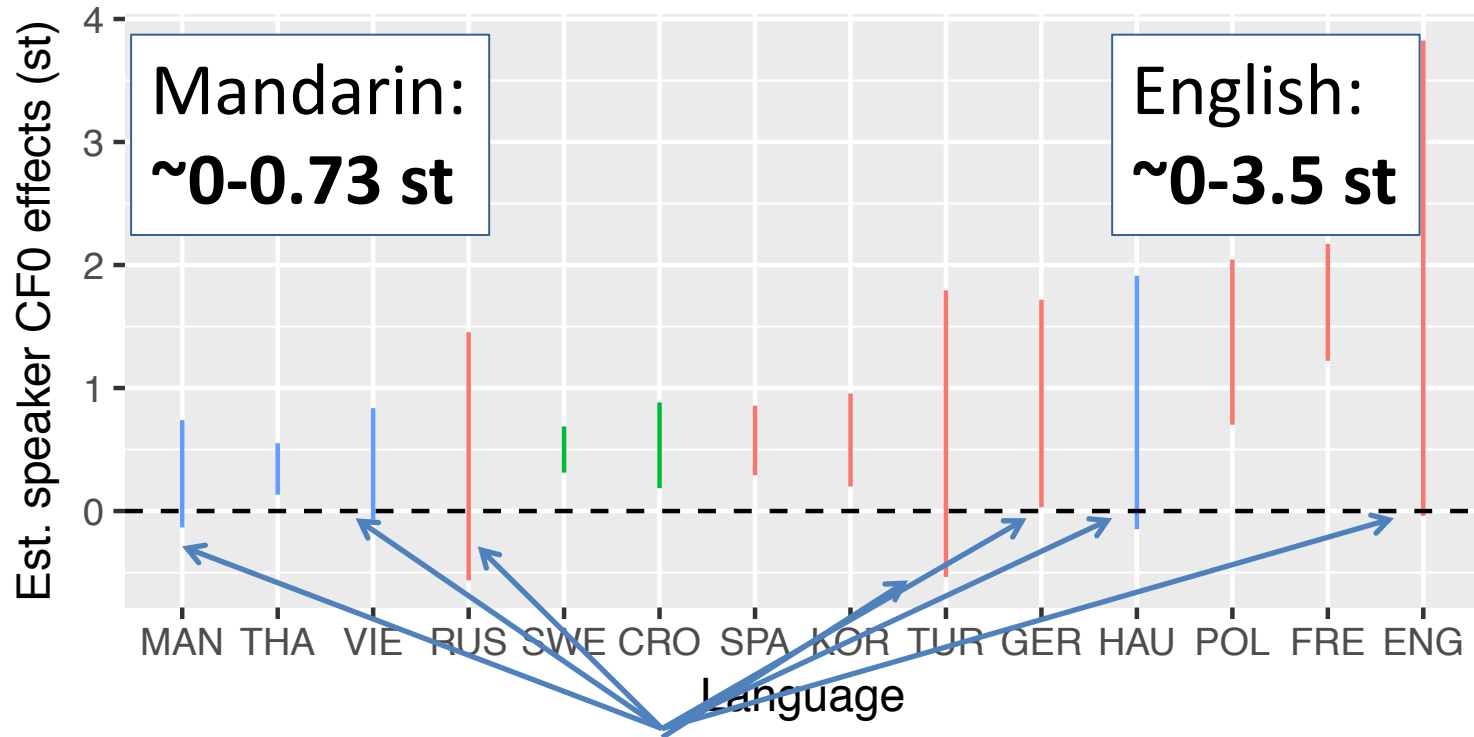

- Conservative model structure

# CF0: across languages

- "most voiceless" – "most voiced" effect:



Mandarin:
**0.32 st**
p/p$^h$ diff.
($p$=0.08)

English:
**1.92 st**
"b"/"p" diff.
($p$<0.001)

tonal    pitch accent    non-tonal

- Robust across languages

- Variable effect size
  – Non-tonal ⇒ larger effect

# CF0: across speakers

- Predicted effects for 95% of individuals:



- Common: large interspeaker variability

# VF0: Analysis

- One linear mixed effects model / language

- Main terms:

Response: mean F0

$$C_1 \quad V \quad X$$

Vowel identity
Height (a vs. i/u)
+ i vs. u
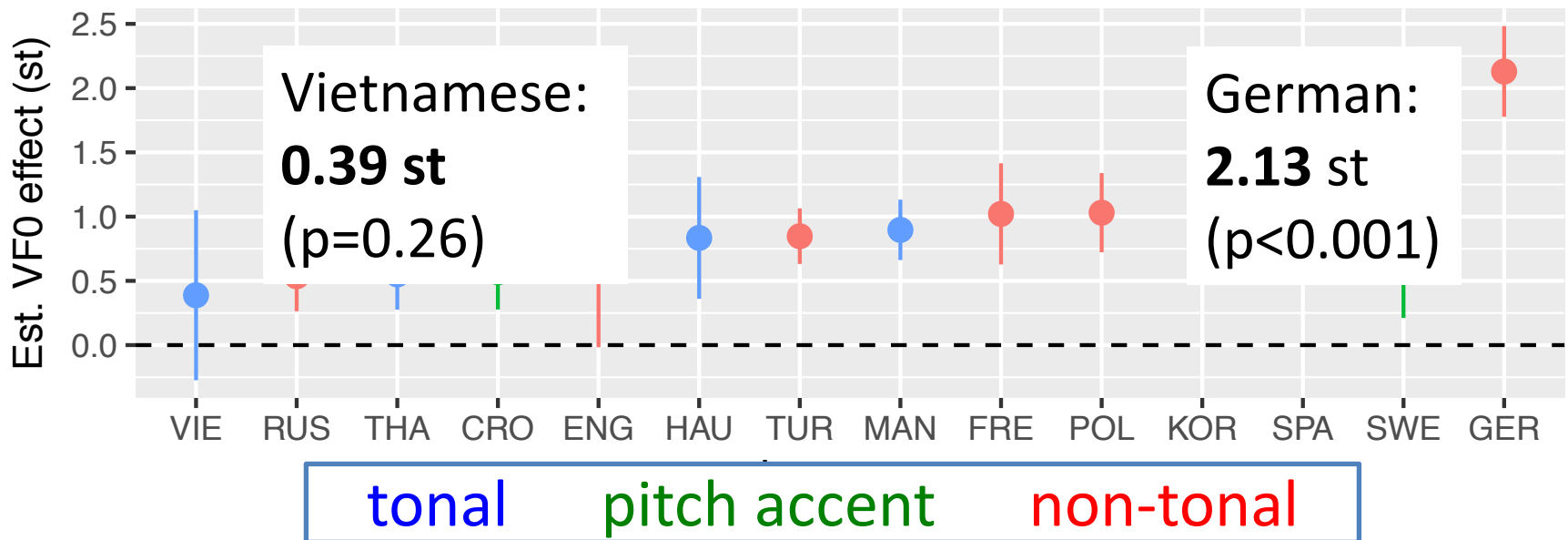
overall effect + interspeaker variability

fixed effect        by-speaker random slope

# VF0: analysis

- + various controls

- Conservative model structure
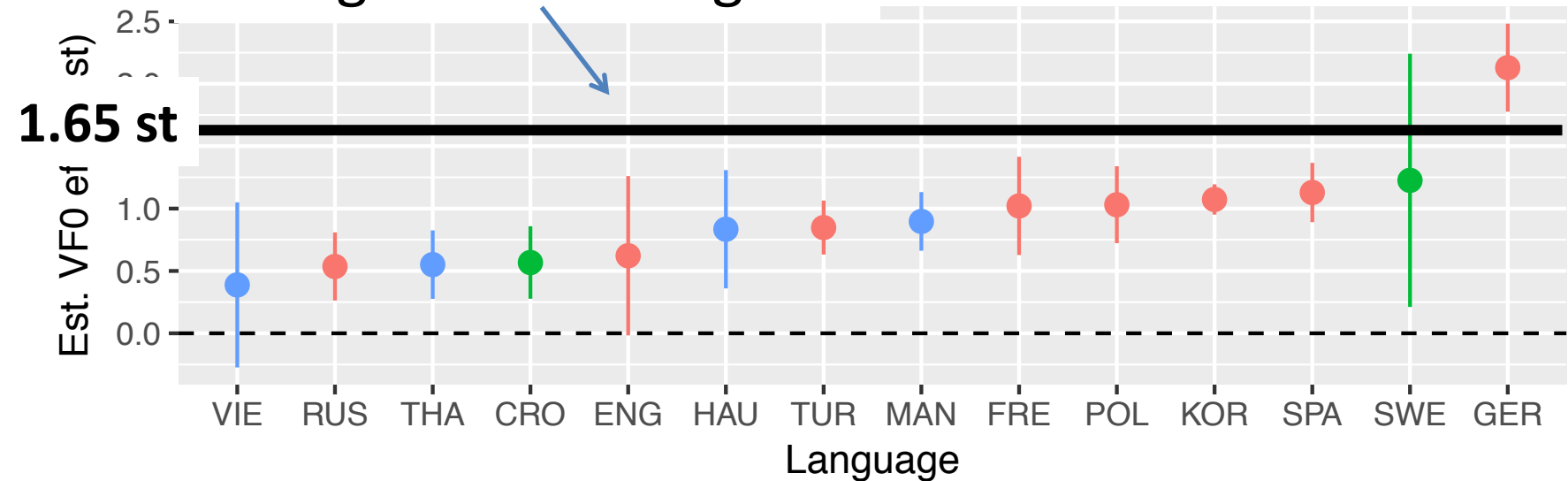
# VF0: across languages

- High – low vowel effect:



- mostly robust across languages

- variable effect size
  – Non-tonal ⇒ generally larger effect

Average effect across gender, tone, etc.

# VF0: across languages



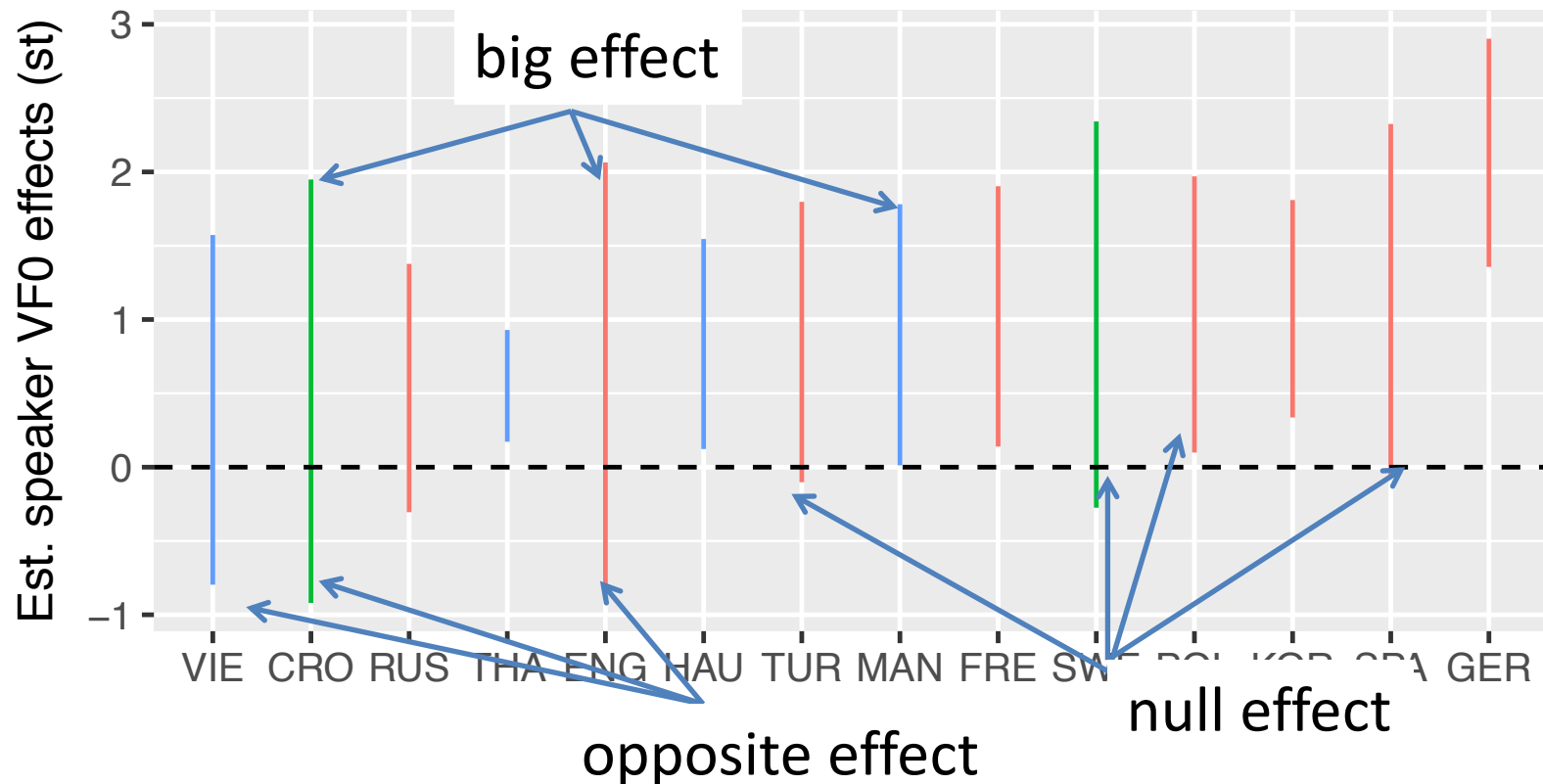- Sentences vs. lab speech?

- (or artifact of methodology?)

# VF0: across speakers

- Predicted effects for 95% of individuals



- Common: large interspeaker variability

# Discussion

- IF0 effects can be detected using
  - Corpus data
  - Fully automatic analysis
  - Basic statistical controls
  - $n = \sim 2\text{-}4k$
- Not obvious!

- Demonstrates feasibility of large-scale studies of phonetic precursors (involving F0)

# Discussion

- Robust group-level IF0 effects across languages
  - same direction
  - "universality" (Whalen & Levitt, 1995)

- Very different effect sizes
  - One reason: tonal/pitch accent language
    $\Rightarrow$ smaller IF0 more likely
    (hypothesized for VF0: Connell 2002)

- Fits with automatic + controlled mechanism
  (c.f. Hoole & Honda, 2011)

# Discussion

- Large interspeaker variability in IF0 magnitude common, within language
  - ⇒ there are some speakers with null/large effects
  - Still, most speakers show effect in same <u>direction</u>

- Overall: IF0 effects
  - robust across languages
  - variable across speakers
- Both important for sound change

- Related to actuation: why sound changes from IF0 possible, but rare? (Kingston, 2007)

# Study 2: duration compression effects





With:    Michael McAuliffe              Michael Wagner

# Introduction

- Major aspect of speech timing: longer linguistic unit $\Rightarrow$ compressed sub-parts

- Ex: _stick_, _sticky_, _sticki_ness (Lehiste, 1972)

- cover term: duration compression effects

# Introduction

- Menzerath's Law (Menzerath, 1928, 1954)
  - 'The longer the whole, the shorter the parts'
  - Domain-general (not just speech)
  - longer words ⇒ shorter average syllable duration
  - phonetic Menzerath effect

- Related: Polysyllabic shortening (Lehiste, 1972)
  - Syllable/V durations shorter in bigger words/prosodic domains

# Introduction

- Extensive work on DCEs
  - individual languages, controlled settings

- Unclear:
  - Are DCEs <u>universal</u>?
    (Siddins et al., 2014; Suomi, 2007; White & Turk, 2010)

Q1: can we observe duration compression effects across typologically-diverse languages?

- Today: test for phonetic Menzerath effect

# Introduction

- Unclear: are DCEs just reducible to other factors?

- Fewer segments per second:
  – Speech rate
  – Longer words ⇒ fewer segments/syllable ("Structural Menzerath effect")

- Prosodic effects on syllable duration:
  – Accent
  – Initial position
  – Final position

(e.g. Sluijter, 1995; Fougeron & Keating, 1997; Oller, 1973; Klatt, 1973, 1975; White & Turk, 2010; Windmann et al., 2015)

# Introduction

- Q2: can DCEs be reduced to fewer segments/second?

- Q3: can DCEs be reduced to a prosodic lengthening effect?

# Datasets

English    Hausa
German    Polish
Russian    Portuguese
Swahili    Spanish
Ukrainian    Swahili
Bulgarian    Turkish

Mandarin
Vietnamese
Thai

Diverse (word) prosody

- Read sentence corpora
  - ~20 hours each
  - Force-aligned

**Montreal Forced Aligner**

French
Korean

Croatian
Swedish

, 2013), TIMIT (Garofolo et al., 1993)

# Datasets

- "Utterance" final "words"

- Measures
  – Word length (# syllables)
  – Mean syllable duration

- Controls
  – Speech rate
  – Expected syllable duration
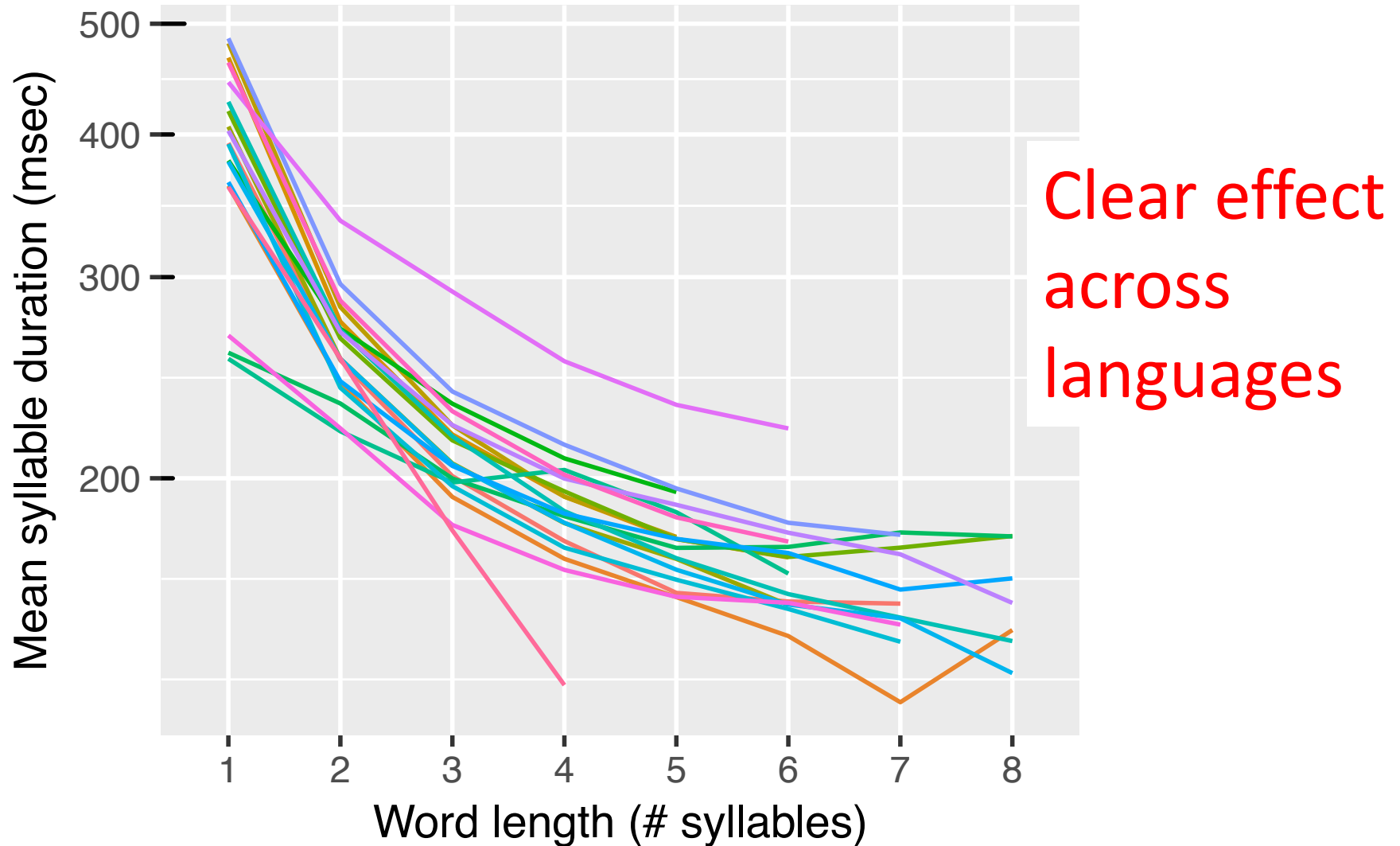  – Speaker, word ID
  – etc.

Polygot-Speech Corpus Tools

Given segmental content, for individual speaker
(Ernestus, Gahl)

# Datasets

- Pruning
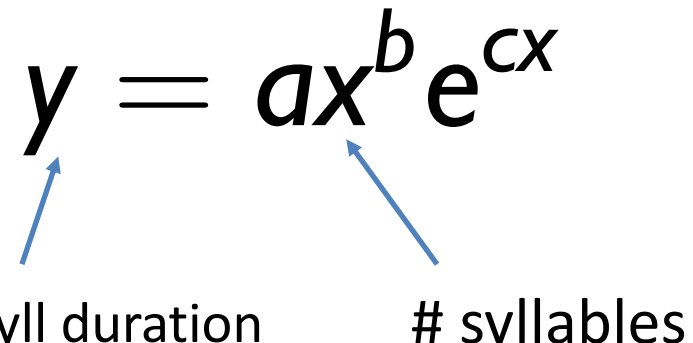  - Words above language-specific cutoff length

# Results: Mean syllable duration



Mean syllable duration (msec) vs Word length (# syllables)
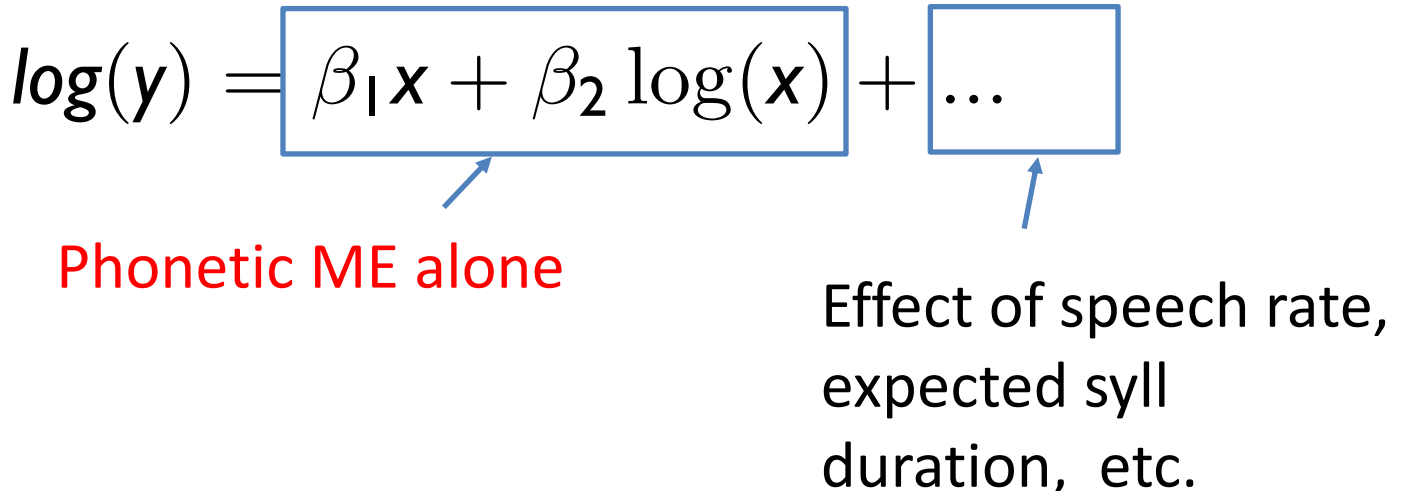
Clear effect across languages

# Analysis 1: controlling for segments/second

- Does mean syllable duration ~ word length, beyond effects of
  - Speech rate
  - Expected syllable duration ← Enhanced "# segments in syllable"
  - Who's talking
  - Particular words
  - Utterance length ?

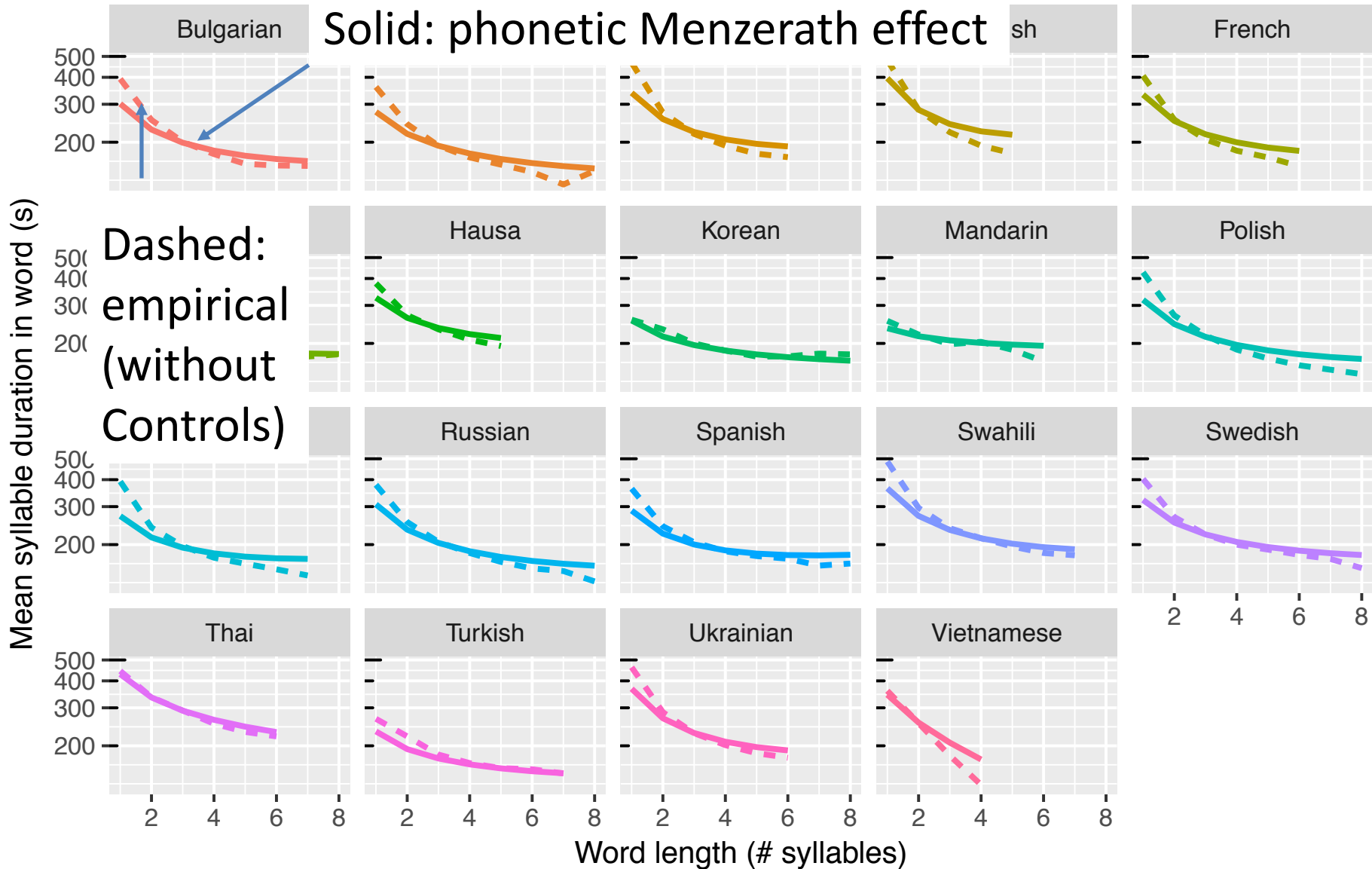# Analysis 1

- Menzerath-Altmann law: $y = ax^b e^{cx}$

  Mean syll duration      # syllables

- Our case:

$$log(y) = \boxed{\beta_1 x + \beta_2 \log(x)} + \boxed{\ldots}$$

  Phonetic ME alone

  Effect of speech rate, expected syll duration, etc.

- Linear mixed-effects model

# Results



Solid: phonetic Menzerath effect

Dashed: empirical (without Controls)

Mean syllable duration in word (s)

Word length (# syllables)

Panels: Bulgarian, [Engli]sh, French, Hausa, Korean, Mandarin, Polish, Russian, Spanish, Swahili, Swedish, Thai, Turkish, Ukrainian, Vietnamese

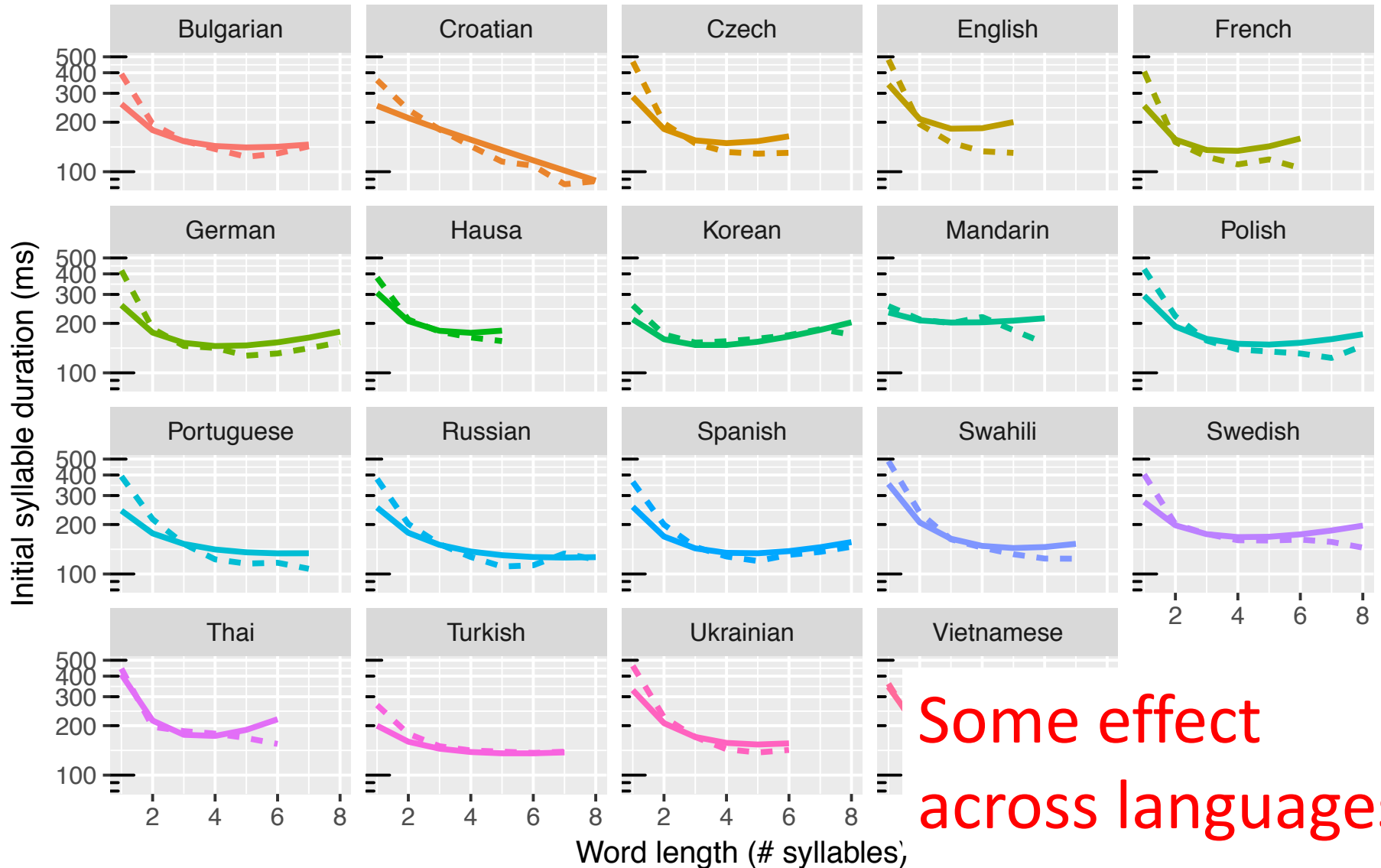# Results: analysis 1

- <span style="color:red">Clear phonetic Menzerath effect across languages</span>

- Q2: are DCEs reducible to segments/second?
  - <span style="color:red">No</span>

- Empirical relationship "steeper":
  - Phonetic M effect (compressed syllables), <u>plus</u>
  - Structural M effect (compressed # segments)

# Analysis 2: prosodic lengthening effects

- Observed effect due to

(White & Turk, 2010;
Windmann et al., 2015)

  - Initial strengthening

  - Final lengthening

  - Accentual lengthening

  - …

  across languages?

- Check:

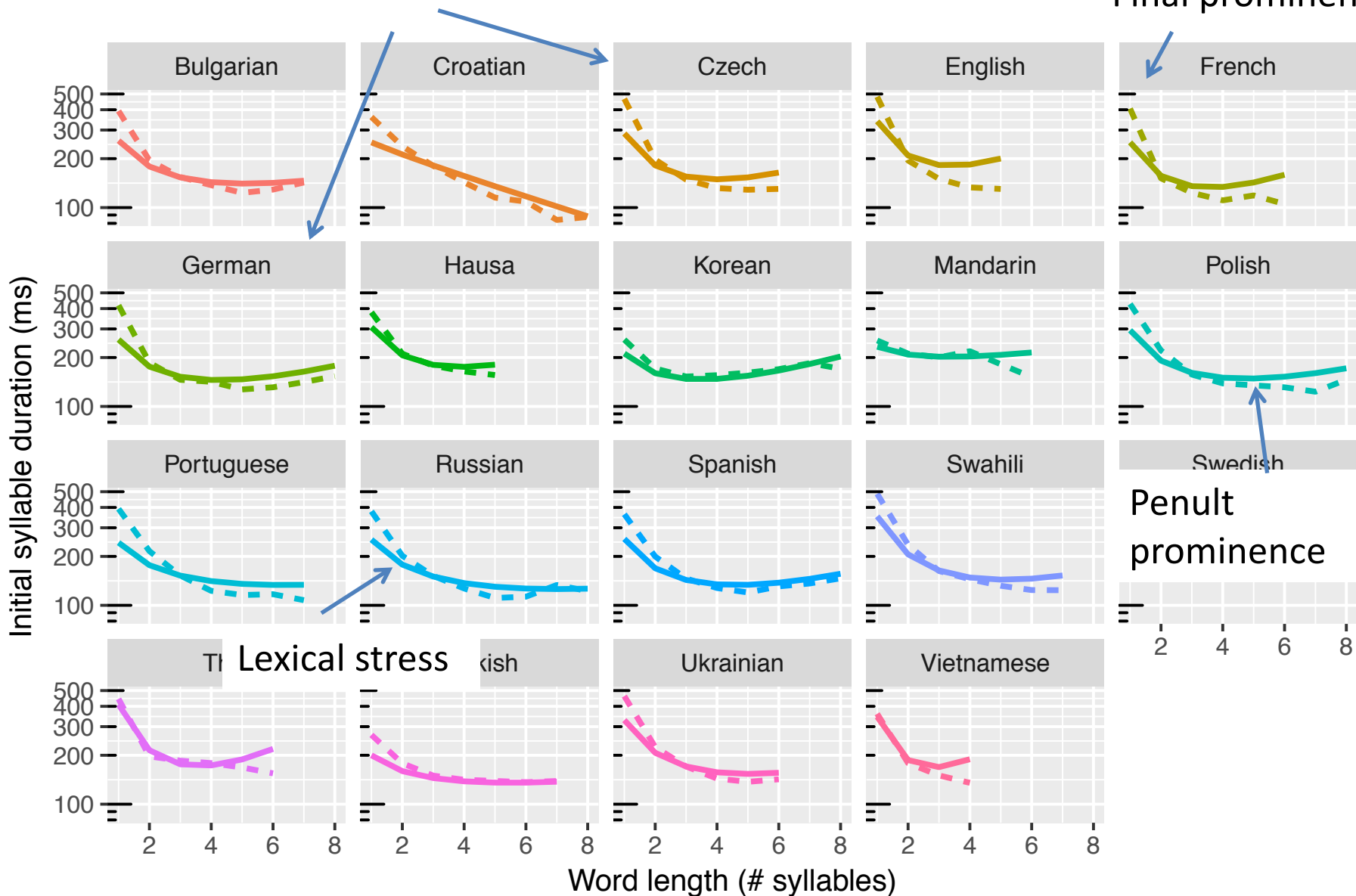  - Same analysis for <u>initial syllable duration</u> only, etc.
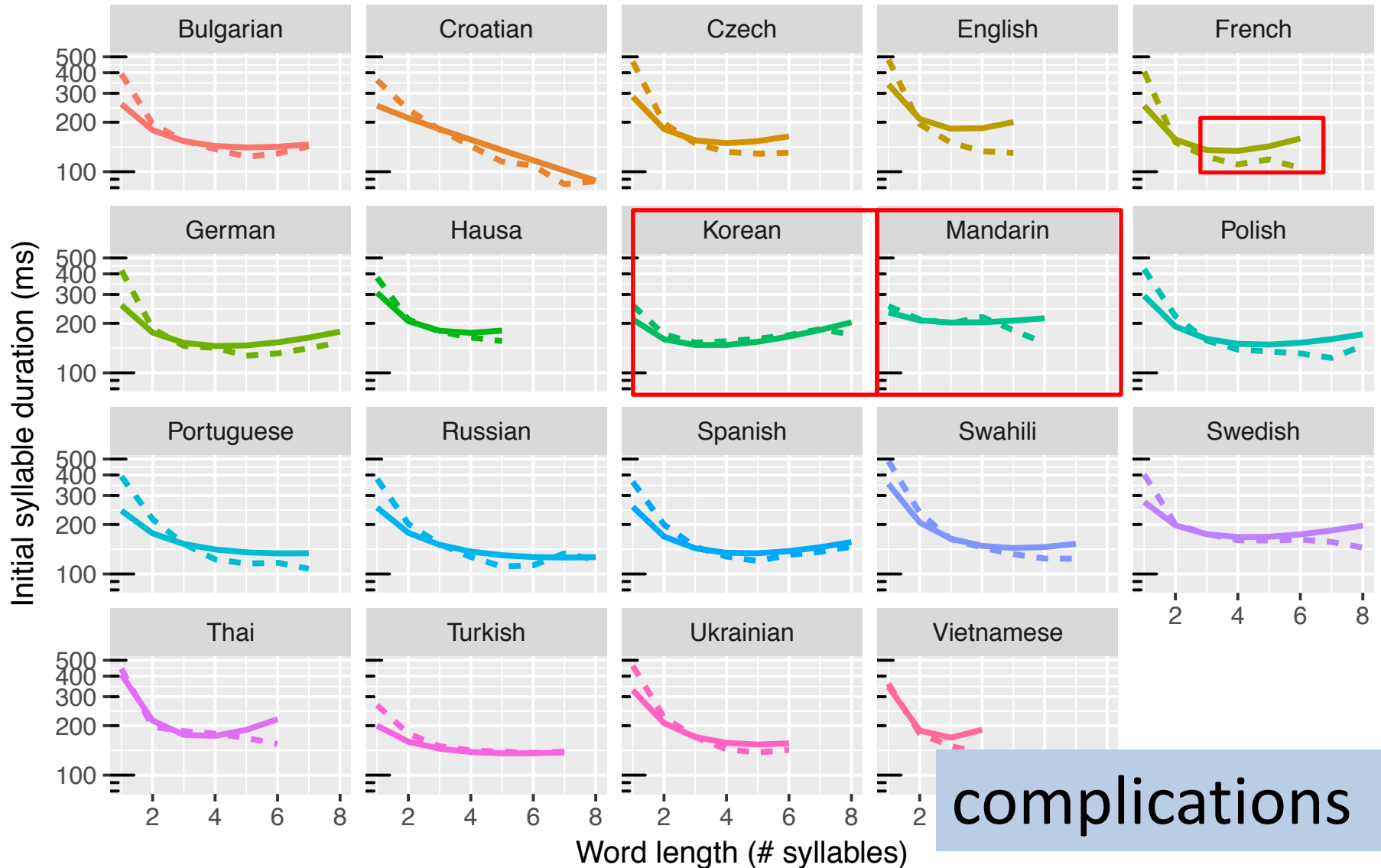
# Results: initial syllables



Initial syllable duration (ms)

Word length (# syllables)

Some effect across languages

# Results: initial syllables



Initial syllable duration (ms) plotted against Word length (# syllables) for each language: Bulgarian, Croatian, Czech, English, French, German, Hausa, Korean, Mandarin, Polish, Portuguese, Russian, Spanish, Swahili, Swedish, Thai, Turkish, Ukrainian, Vietnamese. Annotations: ~ initial prominence, Final prominence, Penult prominence, Lexical stress.

# Results: initial syllables



Initial syllable duration (ms) vs. Word length (# syllables), shown for Bulgarian, Croatian, Czech, English, French, German, Hausa, Korean, Mandarin, Polish, Portuguese, Russian, Spanish, Swahili, Swedish, Thai, Turkish, Ukrainian, Vietnamese.
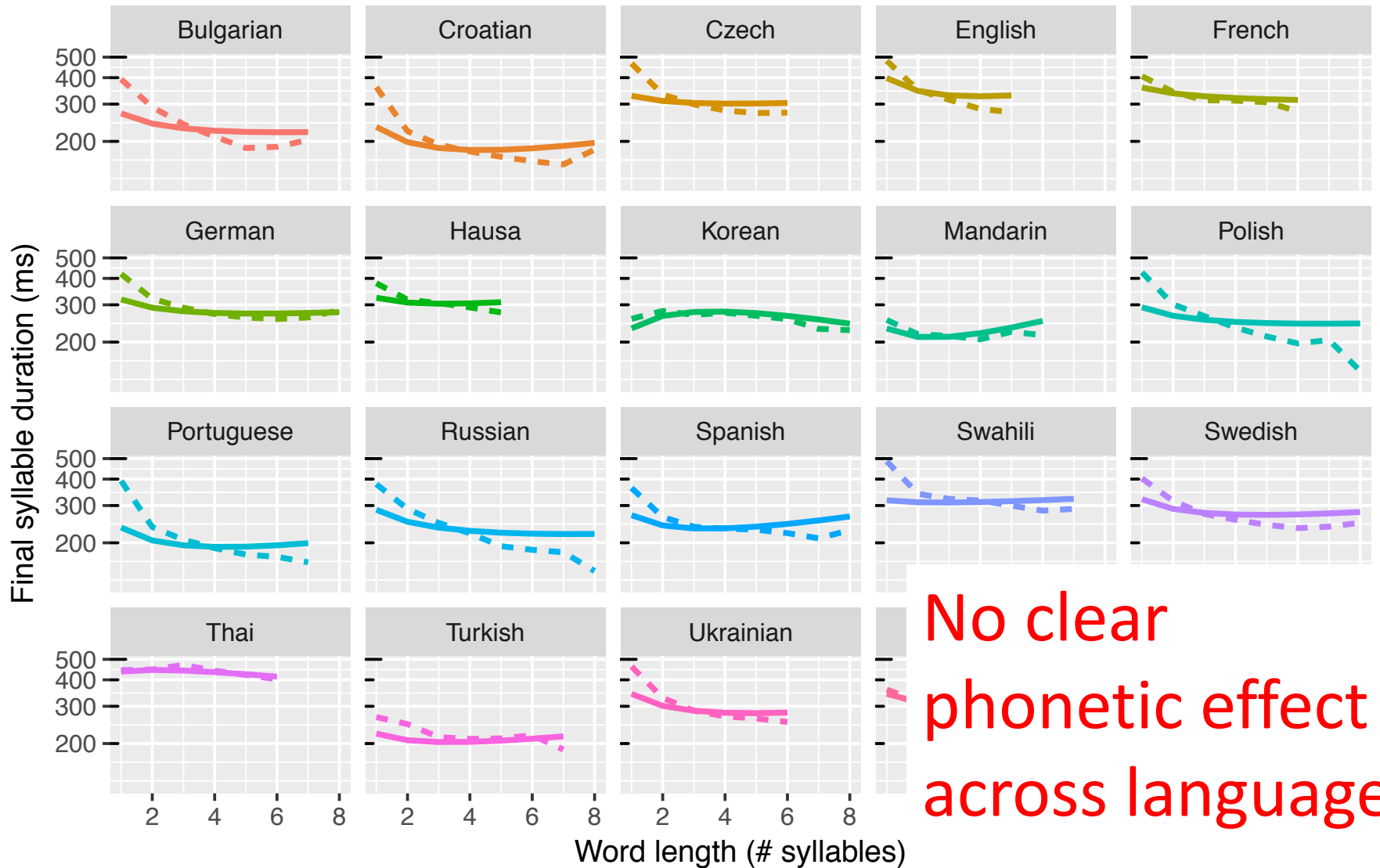
complications

# Results: initial syllables

- Consistent compression effect
  - (at least: 1-3 syllables)

- Very different prosodic systems

- Can't be just
  - Accentual lengthening
  - Initial strengthening
  - PSS on accented syllables only

# Results: final syllables



No clear phonetic effect across languages

# Results: final syllables

- No consistent phonetic compression effect
  - = phonetic ME

- Overridden by other factors?
  - final lengthening, language-specific prosody

- Aside: Much of <u>empirical</u> effect is actually due to fewer segments/syllable
  - = structural ME

# Discussion

1. Duration compression effects may be universal
   – At least phonetic Menzerath effect

2. DCEs not reducible to (some) other factors

- Not obvious!

- (1)+(2) $\Rightarrow$ DCEs reflect something deep about processing/planning
   – Mechanism?

# Thanks

- Michael McAuliffe, Elias Stengel-Eskin, Arlie Coles

- Comments: James Kirby, Simon King, Montreal Language Modeling Lab members

- Funding:

# Questions

# Extra slides

# Barriers to large-scale corpus studies

- Speech datasets:
  - Large
  - Complex
  - Diverse formats

- Access to many speech datasets
  - Costly or ethically restricted

- Result: requires lots of specialized code, $$, effort

# CF0: analysis

- Other terms
  - "Voicing" interactions: gender
  - Controls:
    - Speaker gender, mean F0
    - Utterance length
    - V identity (incl. height)
    - Speaker, word, preceding/following phone
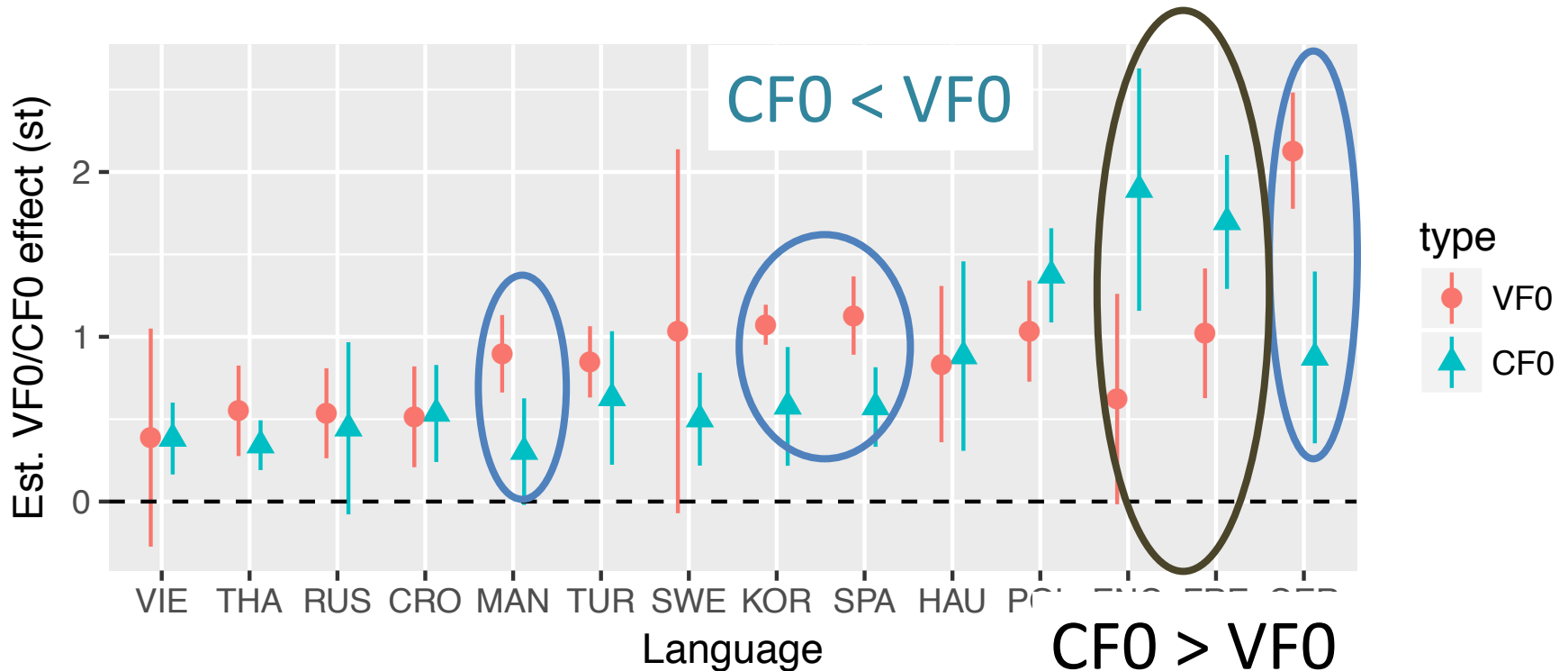
- Conservative model structure

# VF0: analysis

- Other terms
  - V height interactions: gender
  - Controls:
    - Speaker gender, mean F0
    - $C_1$ "voicing"
    - Utterance length
    - V identity
    - Speaker, word, preceding/following phone

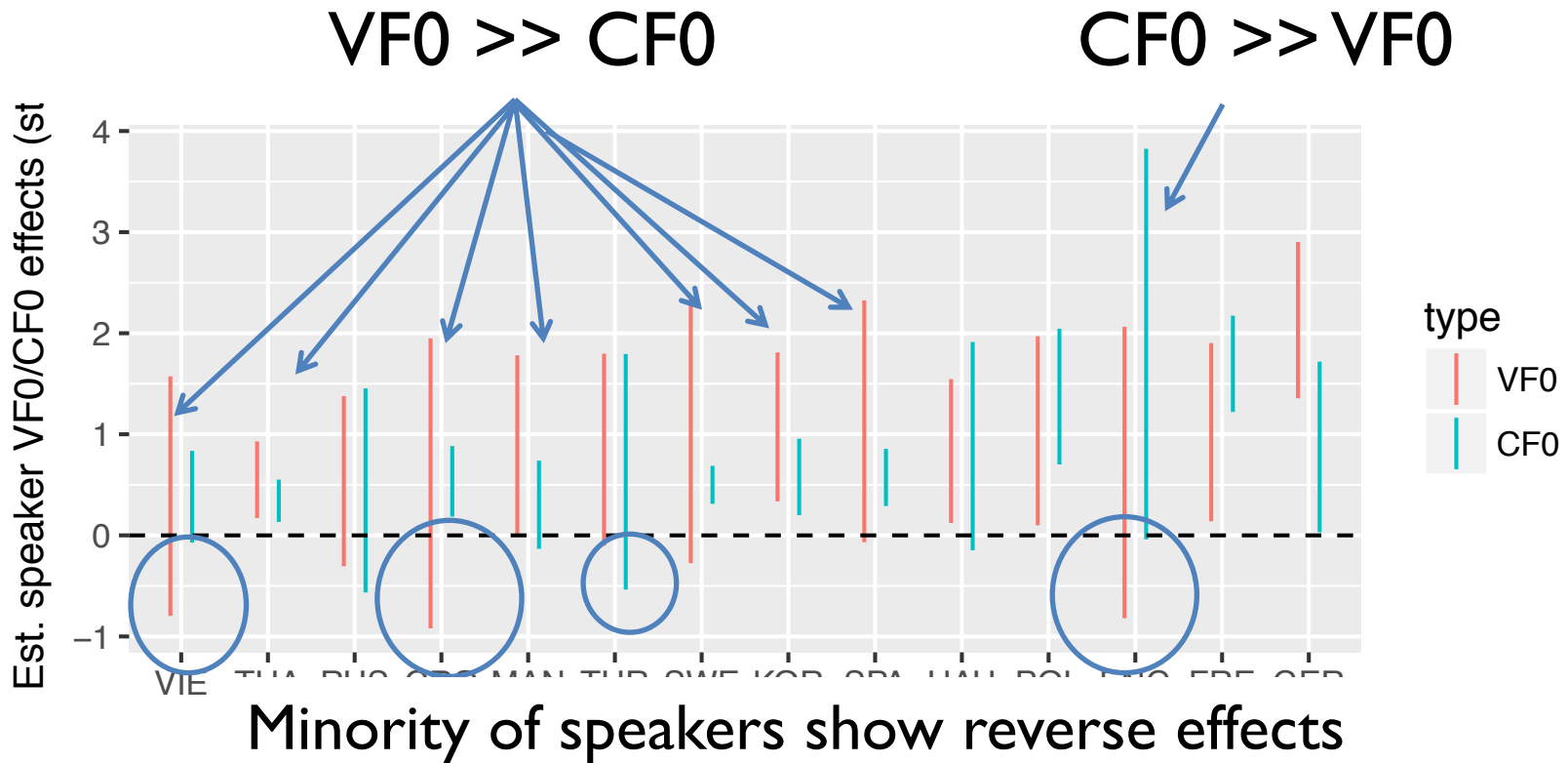- Conservative model structure

# Extra: VF0 vs. CF0

- Asymmetry between IF0 effects w.r.t. sound change:
    - CF0: many attested changes
    - VF0: ~none

- Why?
    - VF0/CF0 magnitude roughly similar? (Hombert et al., 1979)
    - Perhaps perception is different (Hombert, 1979)
    - VF0 effects show more variability? (Kingston, 2011)

- Q4: Relative magnitude, variability of CF0 & VF0 across languages?
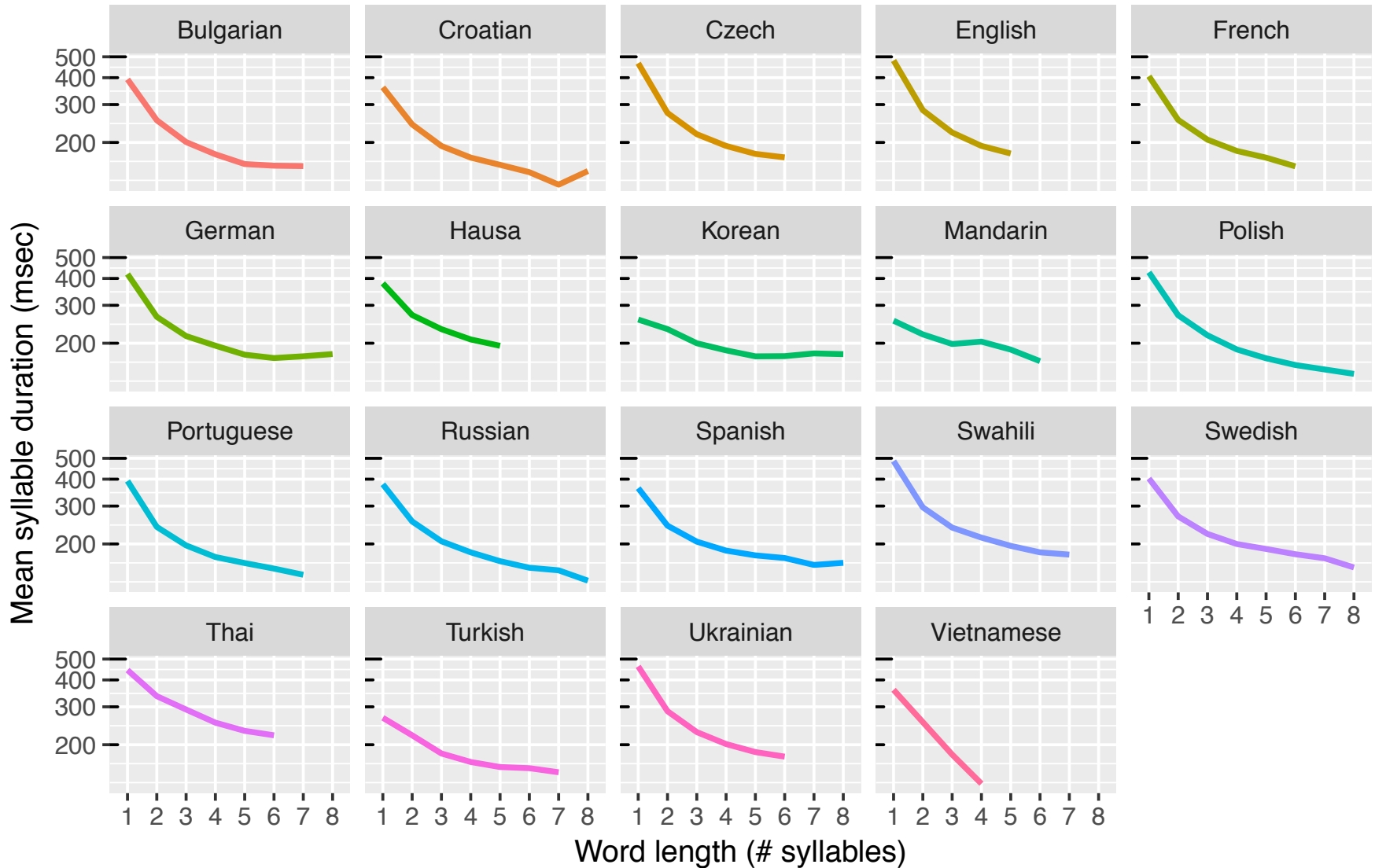
# VF0 vs. CF0: effect size



- <span style="color:red">No clear pattern</span>
- CF0, VF0 of ~comparable size

# VF0 vs. CF0: speaker variability



VF0 >> CF0            CF0 >> VF0

Minority of speakers show reverse effects

- Overall: no obvious pattern
- But: some evidence that VF0 "more variable" than CF0

# Mean syllable duration

# SCT: representation & enrichment

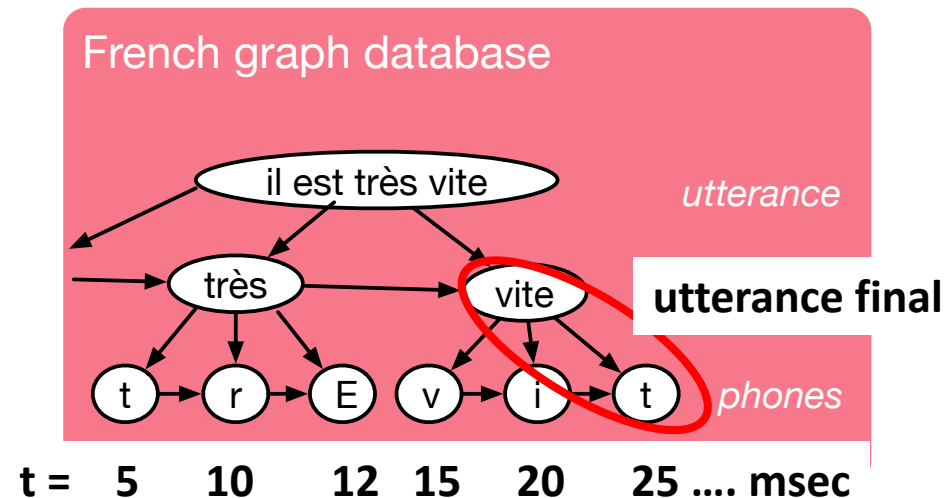- DBs: contains properties of objects, relationships between them:
  - Positional:
    - Ex: Utterance position
  - Hierarchical
    - Ex: containing word
  - Temporal
    - Begin, end, duration



French graph database

il est très vite

très          vite

t → r → E → v → i → t

utterance

**utterance final**

phones

t =   5     10     12   15    20     25 …. msec

- Enrich with additional information:
  - Suprasegmental: pauses, speech rate, ..
  - Acoustic: F0, formants..