

# The SPADE project:

large-scale analysis of a spoken language across  
space and time



Morgan Sonderegger,  
The SPADE Consortium



LSCP Language Group, Paris  
26 Mar, 2019

# SPADE

SPeECH ACROSS DIALECTS OF ENGLISH

August 2017 – July 2020

<http://spade.glasgow.ac.uk/>



THE UNIVERSITY  
of EDINBURGH



University  
of Glasgow



McGill

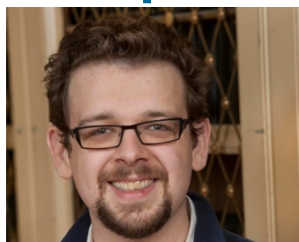


UNIVERSITY OF  
OREGON

# SPADE

SPeech Across Dialects of English

## investigators



<http://spade.glasgow.ac.uk/>

# SPADE

SPeech Across Dialects of English

## Postdocs



THE UNIVERSITY  
of EDINBURGH



Project  
manager +  
data



Software  
development

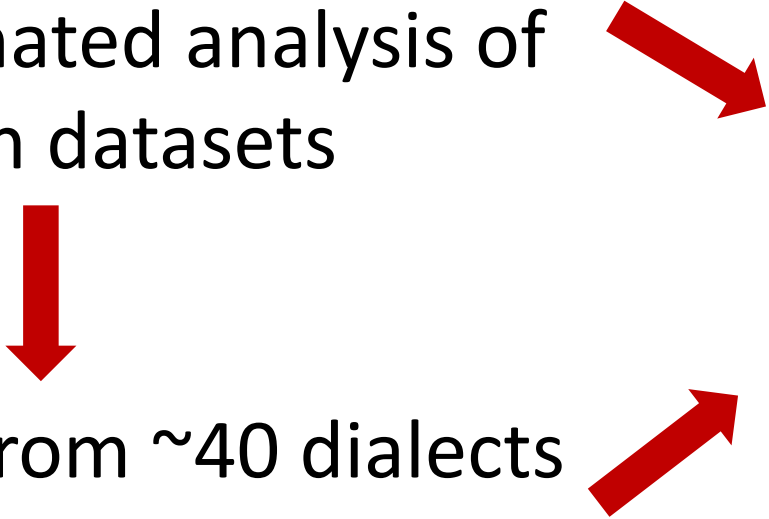
Graduates +  
undergrads: many

<http://spade.glasgow.ac.uk/>

# SPADE

SPeech Across Dialects of English

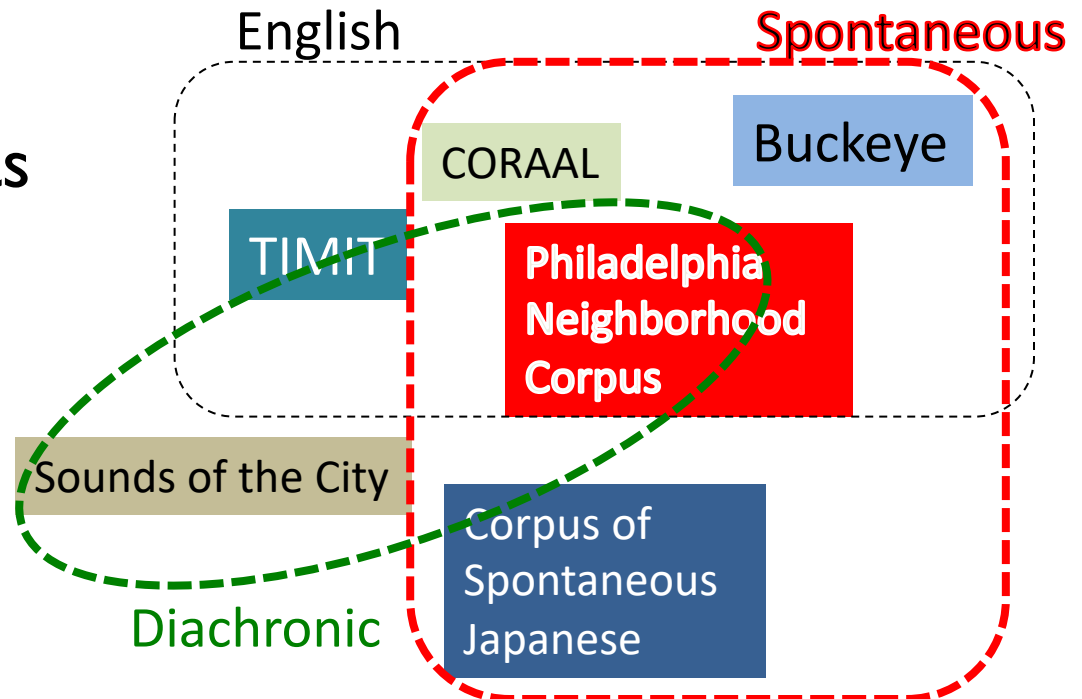
## Project goals

- **Software** for large-scale automated analysis of speech datasets
  - **Data** from ~40 dialects
    - public & private
    - Focus: sociolinguistic data
  - **Case studies:** investigate how “English” varies in time and space
- 

# Motivation

- Huge amount of annotated speech data exists
    - Corpora, academic labs...
- At least orthography + audio

- Across
  - Languages/dialects
  - Speech styles
  - Time



# Motivation

- **Huge amount of annotated speech data exists**
  - Corpora,  
academic labs,  
fieldwork...
- **Across**
  - Languages/dialects
  - Speech styles
  - Time
- **+ ever-better (semi)-automatic speech measurement tools**

# Motivation

- Great potential for speech analysis for different purposes
  - Bigger haystacks, same-sized needle...
  - ... need a bigger magnet
- Requires software for **unified corpus analysis**
  - Integrating speech datasets
  - Querying across them
- SPADE focus: sociolinguistic, phonetic datasets



# Barriers

- Speech datasets:
  - Large
  - Complex
  - Diverse formats

- Access to many speech datasets
  - Costly or ethically restricted

Most sociolinguistic,  
laboratory data

Switchboard: \$3000+

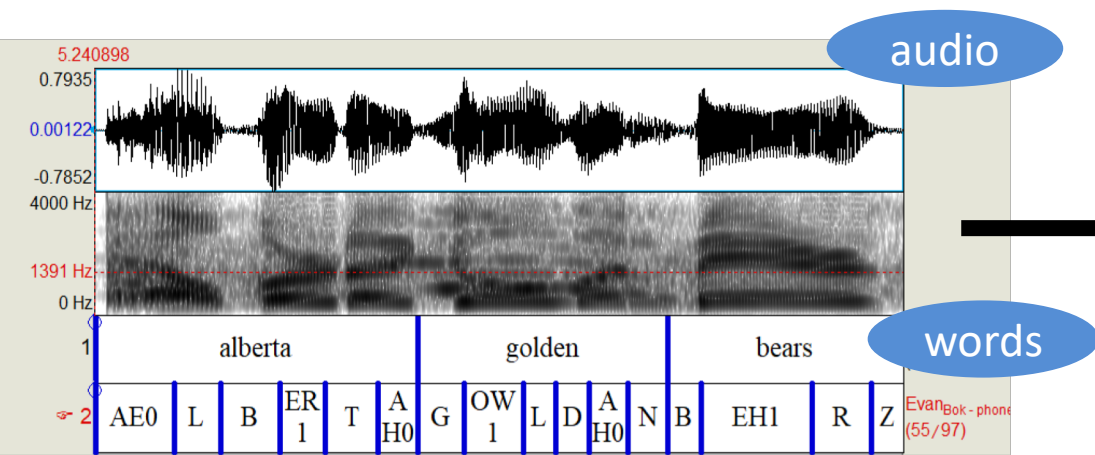
- Result: requires lots of specialized code, €€, effort, computational power

# Software goals

- Scalable & fast
- Require minimal technical skill from user
- Abstraction away from dataset format
- Querying dataset without access to raw data
- ⇒ Easier large-scale studies using speech corpora

- To motivate structure of software:
  - think about steps researcher goes through to do a (speech) corpus study
  - Running example: vowel formants
- Setting: sociolinguistic study, or laboratory phonology, phonetics, etc.

# Raw data



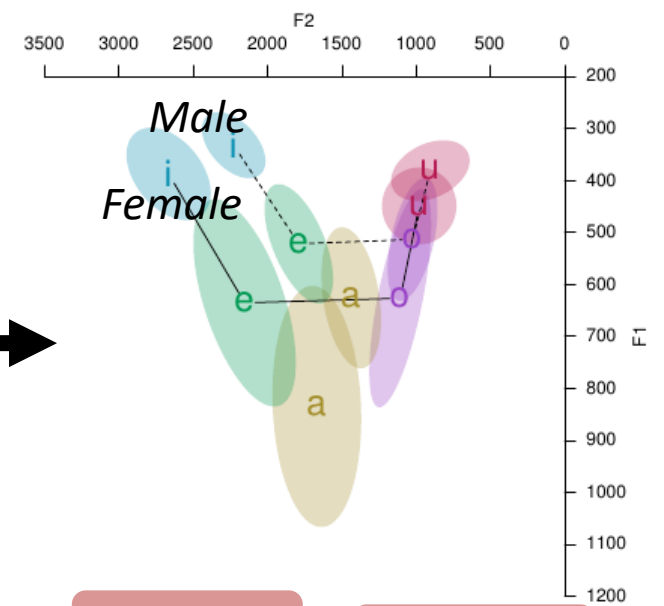
phones

utterances

speaker questionnaire

Speaker M01  
Gender: M  
Age: 35  
...

# Analysis



Gender

formants

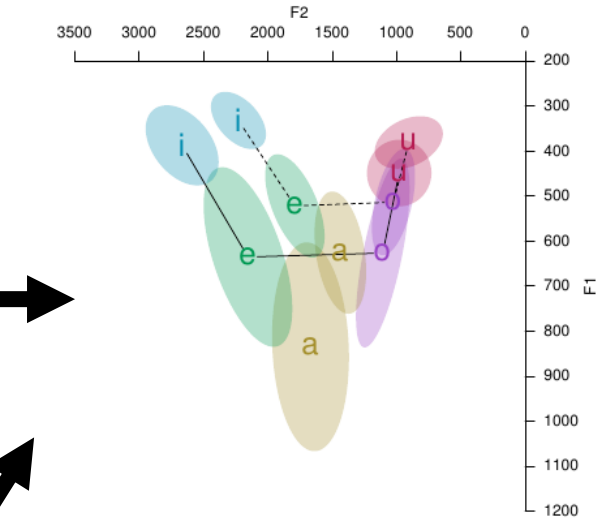
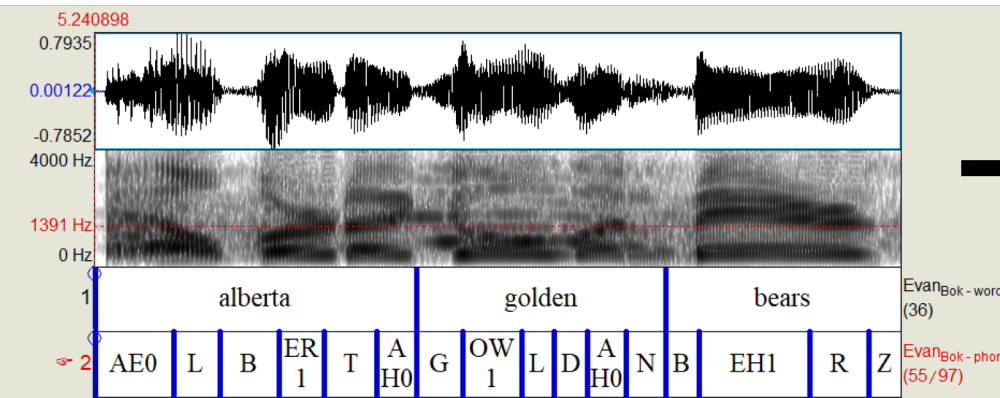
Age

adjacent segments

Vowel duration

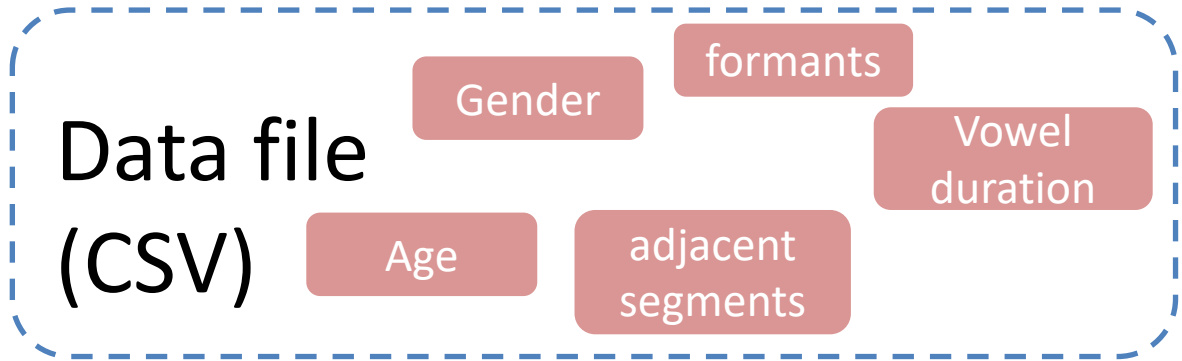
# Raw data

# Analysis



(R, Goldvarb...)

**How?**



1. Process raw data
2. Make measures
3. Find relevant tokens
4. End up with usable spreadsheet

# Raw data

speaker  
questionnaire

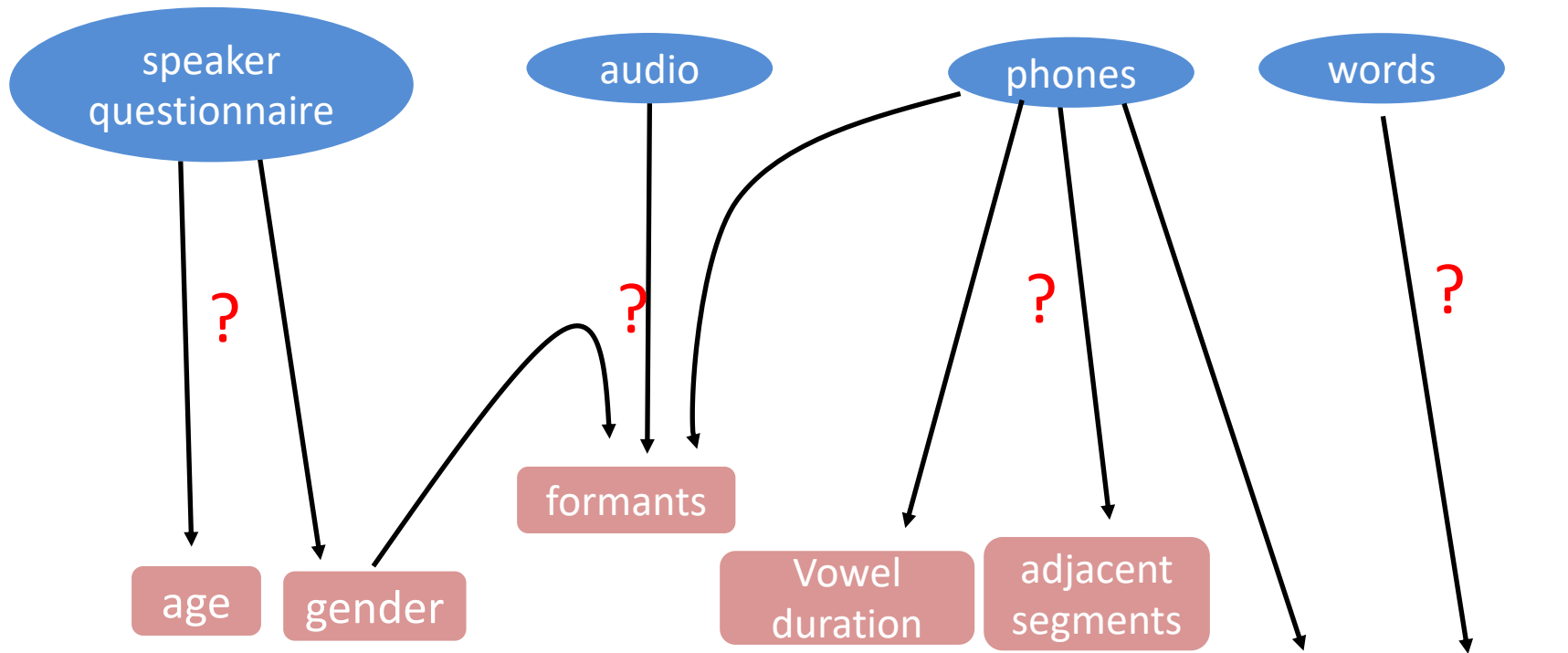
audio

phones

words

| <i>SpkrID</i> | <i>Age</i> | <i>Gender</i> | <i>F1</i> | <i>F2</i> | <i>V_dur</i> | <i>C_left</i> | <i>C_right</i> | <i>Phone</i> | <i>Word</i> |
|---------------|------------|---------------|-----------|-----------|--------------|---------------|----------------|--------------|-------------|
| M01           | 35         | M             | 340       | 1410      | 0.11 s       | [g]           | [s]            | UW           | goose       |
| F01           | 22         | F             | 480       | 1050      | 0.15 s       | [r]           | [r]            | UW           | truer       |
| ...           | ...        |               |           |           |              |               |                | ...          | ...         |

Data file(CSV)



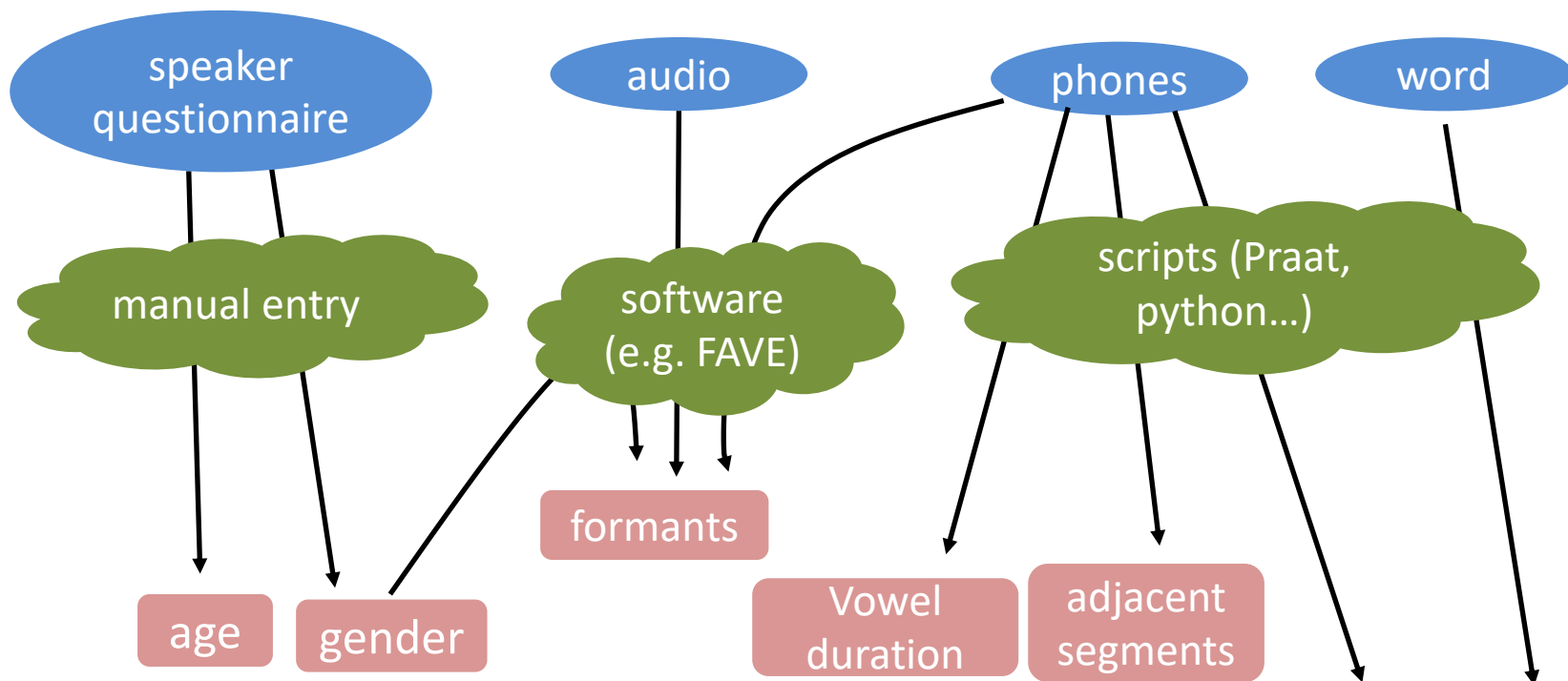
| <i>SpkrID</i> | <i>Age</i> | <i>Gender</i> | <i>F1</i> | <i>F2</i> | <i>V_dur</i> | <i>C_left</i> | <i>C_right</i> | <i>Phone</i> | <i>Word</i> |
|---------------|------------|---------------|-----------|-----------|--------------|---------------|----------------|--------------|-------------|
| M01           | 35         | M             | 340       | 1410      | 0.11 s       | [g]           | [s]            | UW           | goose       |
| F01           | 22         | F             | 480       | 1050      | 0.15 s       | [r]           | [r]            | UW           | truer       |
| ...           | ...        |               |           |           |              |               |                | ...          | ...         |

social factors

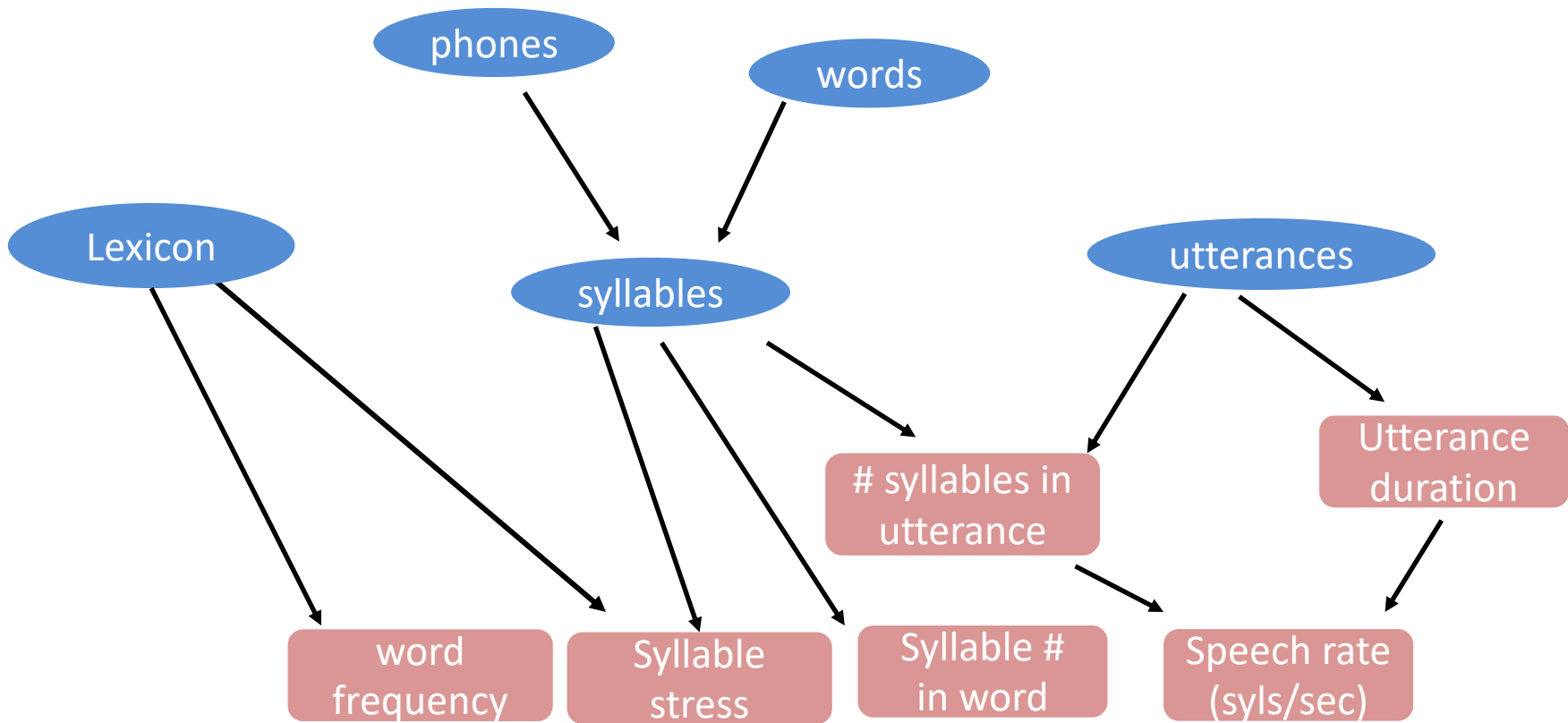
acoustic measures

phone info





| <i>ID</i> | <i>Age</i> | <i>Gender</i> | <i>F1</i> | <i>F2</i> | <i>V_dur</i> | <i>C_left</i> | <i>C_right</i> | <i>Phone</i> | <i>Word</i> |
|-----------|------------|---------------|-----------|-----------|--------------|---------------|----------------|--------------|-------------|
| M01       | 35         | M             | 340       | 1410      | 0.11 s       | [g]           | [s]            | UW           | goose       |
| F01       | 22         | F             | 480       | 1050      | 0.15 s       | [r]           | [r]            | UW           | truer       |
| ...       | ...        |               |           |           |              |               |                | ...          | ...         |



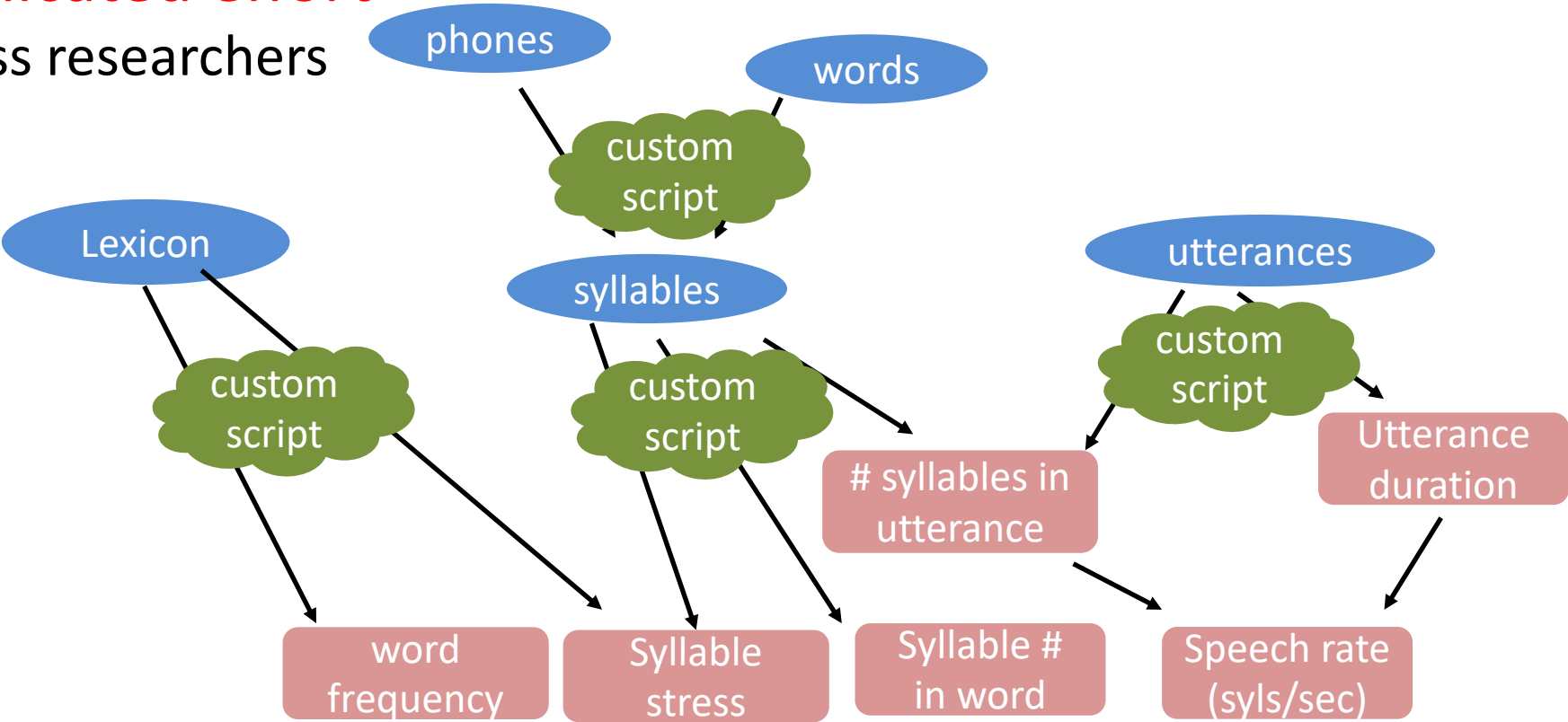
| <i>Phone</i> | <i>Word</i> | <i>Frequency</i> | <i>Stress</i> | <i>Syll_#</i> | <i>Speech rate (sylls/sec)</i> |
|--------------|-------------|------------------|---------------|---------------|--------------------------------|
| UW           | goose       | 15               | Y             | 1             | 0.11 s                         |
| UW           | truer       | 25               | Y             | 1             | 0.15 s                         |
| ...          | ...         |                  |               |               | ...                            |

Lexical info

Prosodic/suprasegmental

# Duplicated effort

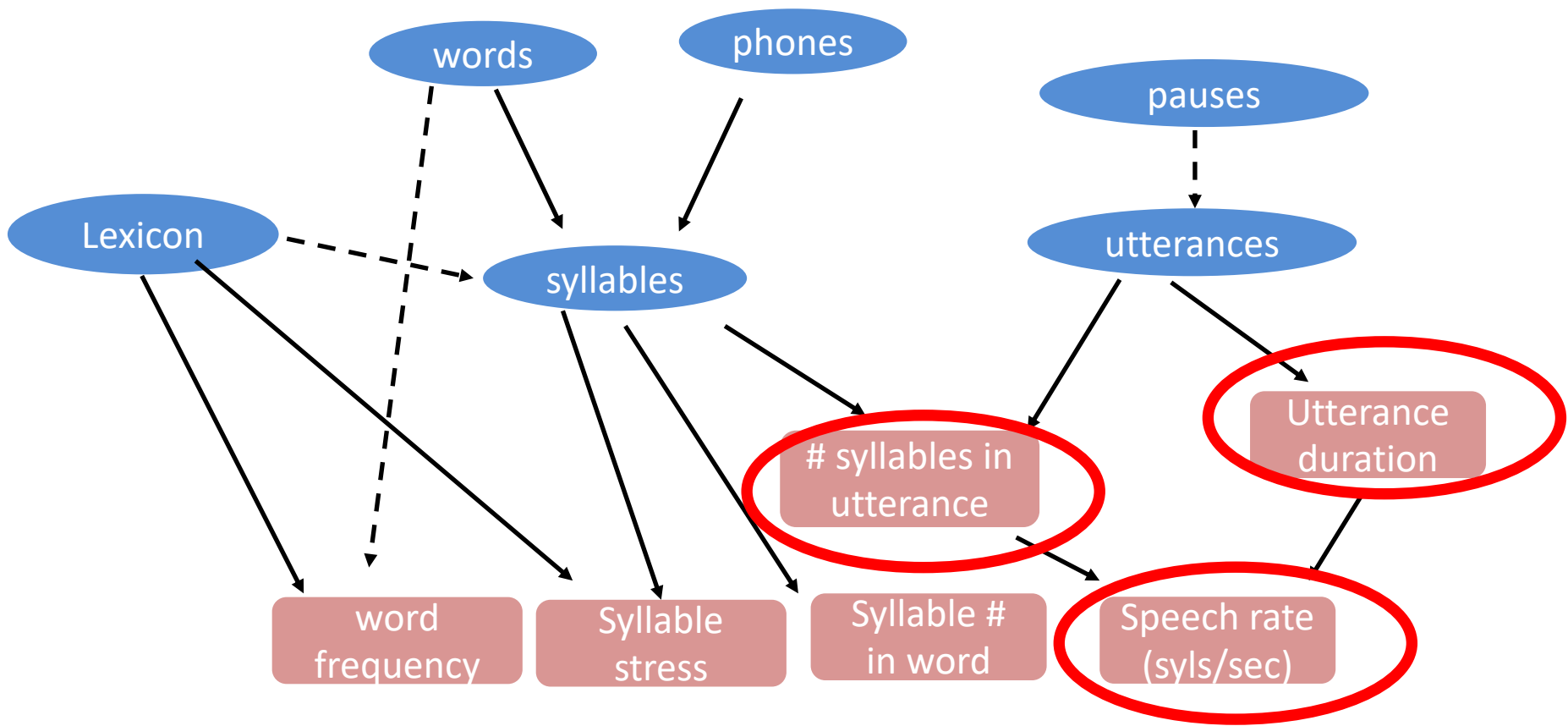
across researchers



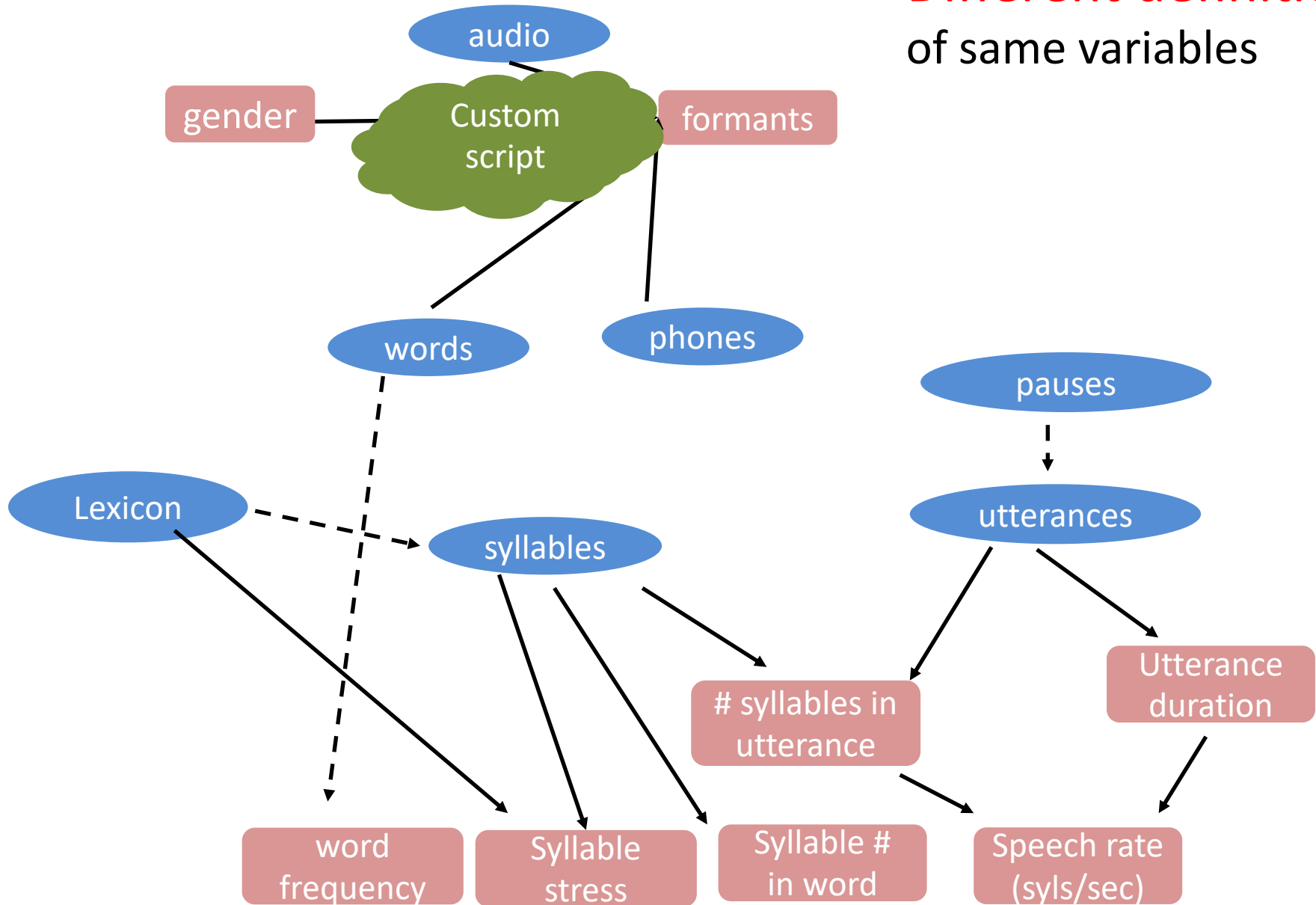
| <i>Phone</i> | <i>Word</i> | <i>Frequency</i> | <i>Stress</i> | <i>Syll_#</i> | <i>Speech rate (sylls/sec)</i> |
|--------------|-------------|------------------|---------------|---------------|--------------------------------|
| UW           | goose       | 15               | Y             | 1             | 0.11 s                         |
| UW           | truer       | 25               | Y             | 1             | 0.15 s                         |
| ...          | ...         |                  |               |               | ...                            |

Lexical info

Prosodic/suprasegmental

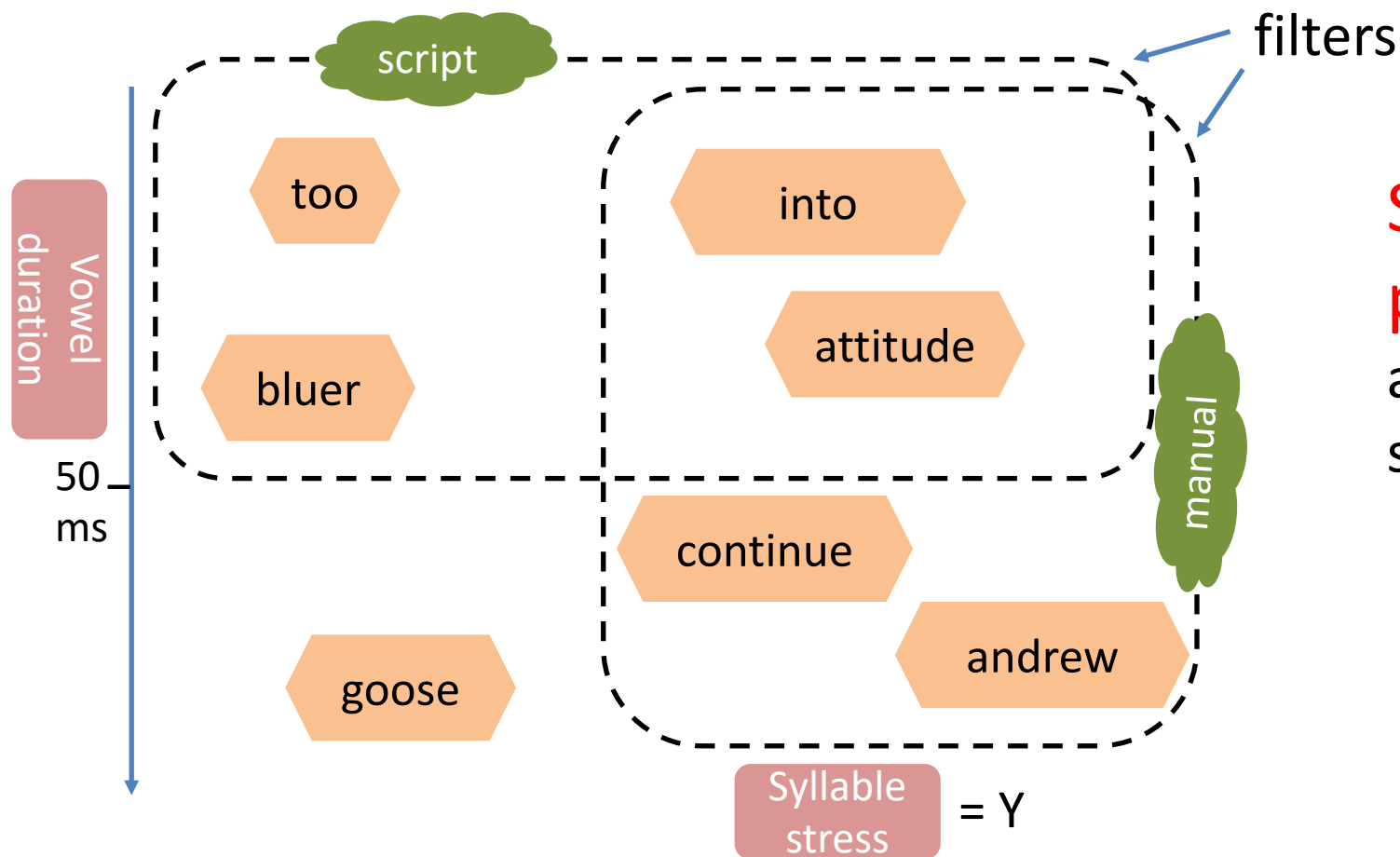


# Different definitions of same variables



1. Process raw data
2. Make measures
3. Find relevant tokens
4. End up with usable spreadsheet

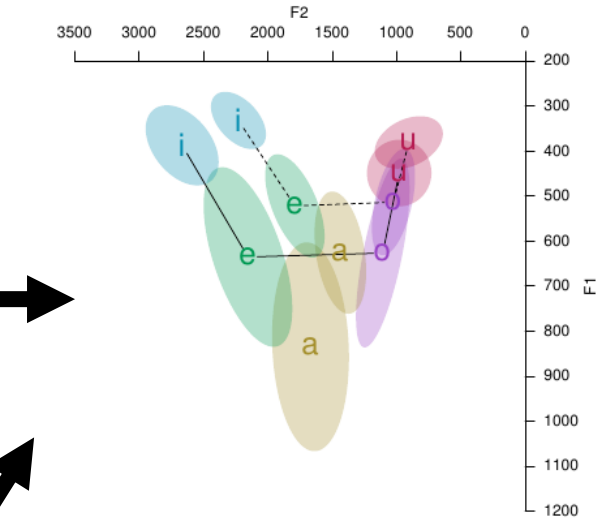
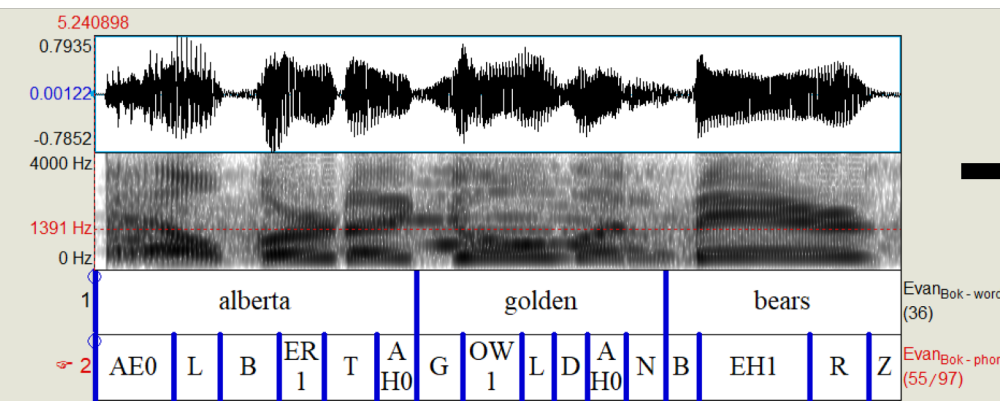
e.g. UW (GOOSE) tokens, stressed syllables > 50ms



Same  
primitives  
as other  
steps

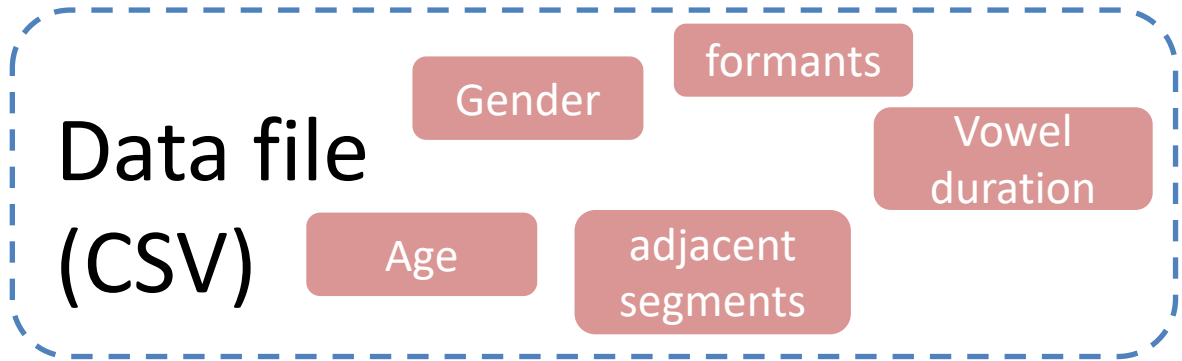
# Raw data

# Analysis



(R, Goldvarb...)

- **scripts**
- **manual entry**
- **software**





# Why 'Integrated' Speech Corpus Analysis?

– Practical reasons

ISCAN

- Technical skill
- Time/duplication of effort
- Availability

– Methodological/theoretical reasons

- Standardized, customizable linguistic measures

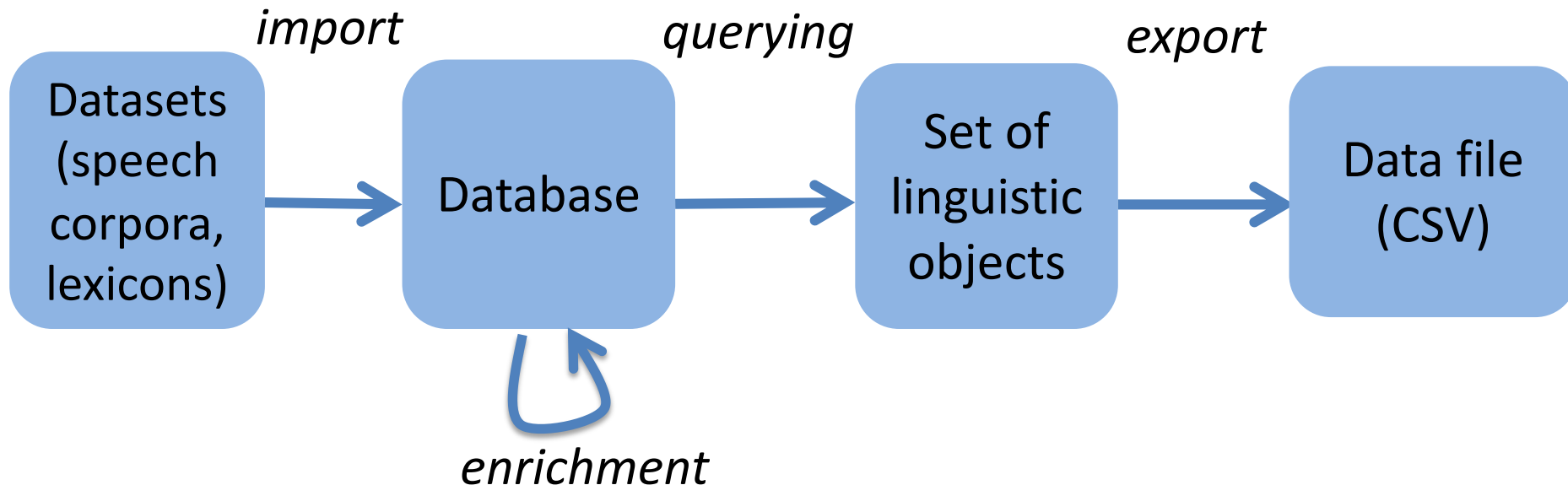
– much more difficult with I+ corpora...

# ISCAN: A SYSTEM FOR INTEGRATED PHONETIC ANALYSES ACROSS SPEECH CORPORA

Michael McAuliffe<sup>a</sup>, Arlie Coles<sup>a</sup>, Michael Goodale<sup>a</sup>, Sarah Mihuc<sup>a</sup>, Michael Wagner<sup>a</sup>, Jane

*Proc. ICPHS 2019*

Stuart-Smith<sup>b</sup>, Morgan Sonderegger<sup>a</sup>



- **Implementation**

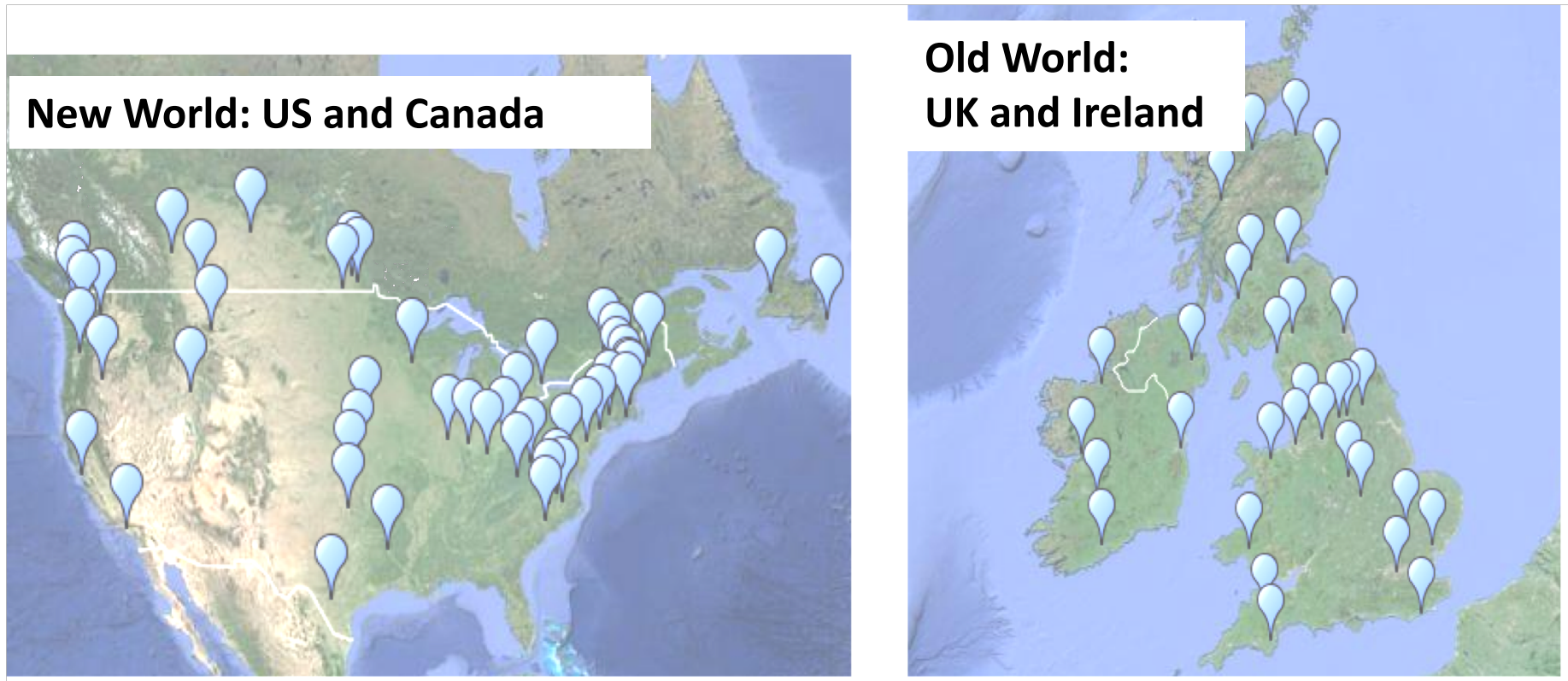
- Python API
- Graphical User Interface

<https://github.com/MontrealCorpusTools/iscan-spade-server>

<https://iscan.readthedocs.io/>

- (show GUI here)
- Note:
  - **Server-client architecture** enables analysis without access to raw data
  - **Permissions system** controls who can see/hear tokens
  - Can be installed on web server (default) or personal computer

# SPADE: datasets



- 43+ datasets, 4 countries, 115 years
- heterogeneous corpus formats
- public and private

# SPADE: datasets

To date:

- Acquired: 20
- Measurements generated: 10
- ~ 10 dialect regions
  - ~500 hours (?)

# SPADE: ethics and credit

- For private datasets (data guardians):  
**ethics** complex: GDPR + US laws
- Data transfer agreement
  - data use in keeping with original permissions, as far as is possible
- **We welcome new datasets!**
- “SPADE consortium” author on everything
- Datasets of measures -> data guardians at end of project

# Case studies

|             |   |   |
|-------------|---|---|
| Done        | Sibilants<br>Stuart-Smith et al.<br><i>Proc. ICPhS 2019</i> | Vowels<br>Mielke et al. <i>Proc.</i><br><i>ICPhS 2019</i> |
| In progress | Vowels 2  | Stops<br>Vowels 3   |
| Planned     | Vowels (dynamic)  | r, l  |

# /s/-retraction in English

- /s/ → [ʃ]-like sound in /str/
  - *string, street*
- Sound change, varies by dialect:
  - Ex: London, Philadelphia, NZ English
  - but not others, e.g. RP, Australian English  
(e.g. Baker et al, 2011; Stevens and Harrington, 2016)
- varies by individual speaker within dialect



# Research questions

- Q1: what is the evidence for /s/-retraction across English?
- Q2: Do English dialects show a dichotomous pattern of /s/-retraction
  - or a continuum?
- Received wisdom: there are “retracting”/non-retracting dialects and speakers

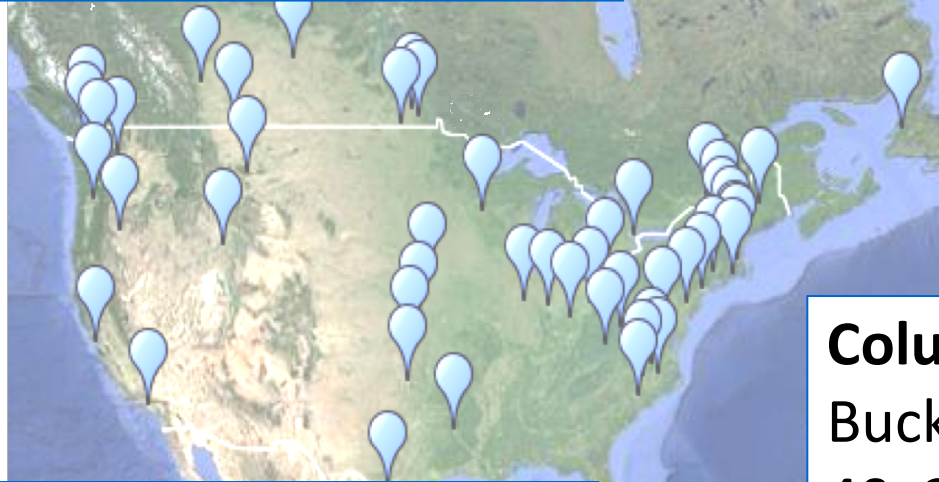
# SPADE

Sample for this study: New World

## Canada

ICECAN Corpus

28: 18m, 10f



## Northern Cities, e.g. New York, Philadelphia

Santa Barbara Corpus

20: 9m, 11f

## Columbus, Ohio

Buckeye Corpus

40: 20m, 20f

## West coast/California

Santa Barbara Corpus

46: 20m, 26f

## Raleigh, North Carolina

Raleigh Corpus

101: 50m, 51f

[www.google.com/maps/](http://www.google.com/maps/)

235 speakers

# SPADE

Sample for this study: New World

## Canada

ICECAN Corpus

28: 18m, 10f

## West coast/California

Santa Barbara Corpus

43: 20m, 23f

**reported to show  
/s/-retraction**

## Northern Cities, e.g. New York, Philadelphia

Santa Barbara Corpus

19: 8m, 11f

## Columbus, Ohio

Buckeye Corpus

40: 20m, 20f

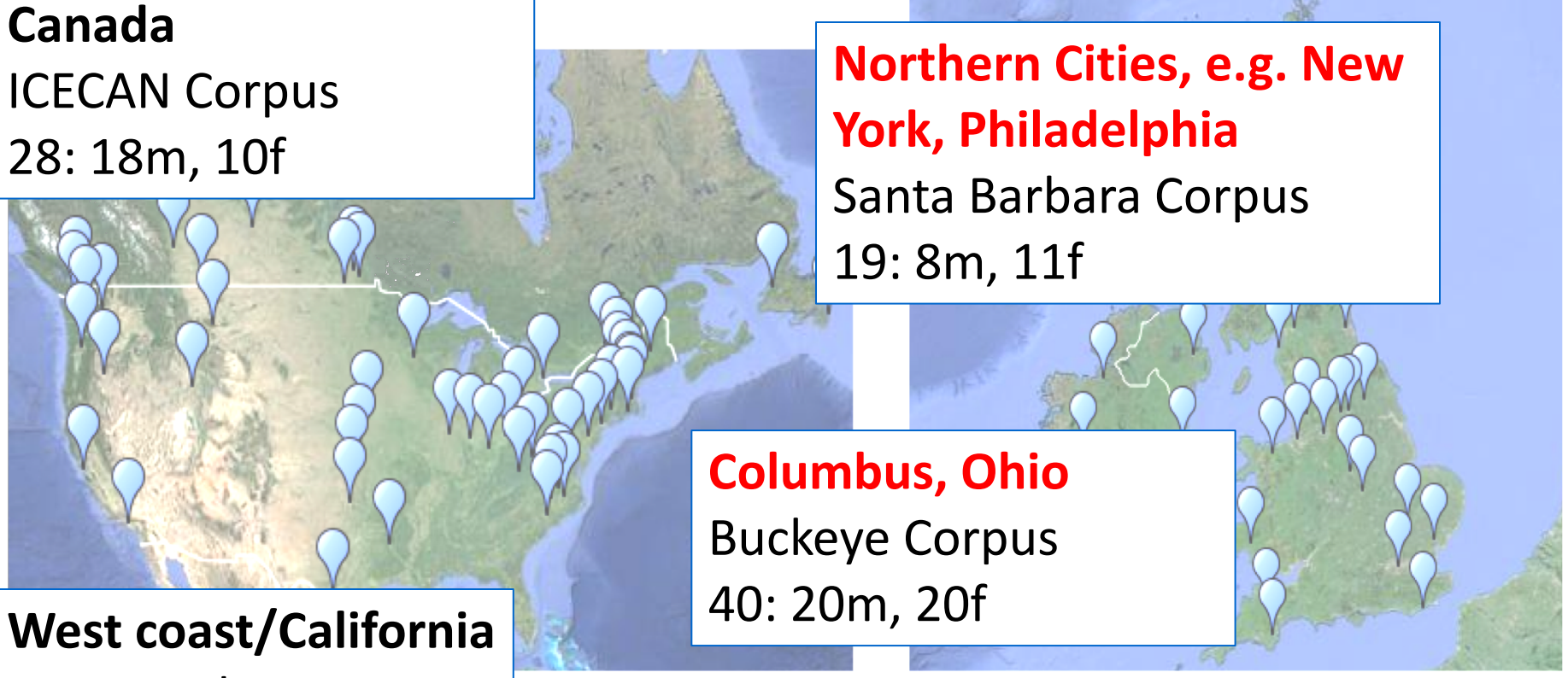
## Raleigh, North Carolina

Raleigh Corpus

101: 50m, 51f

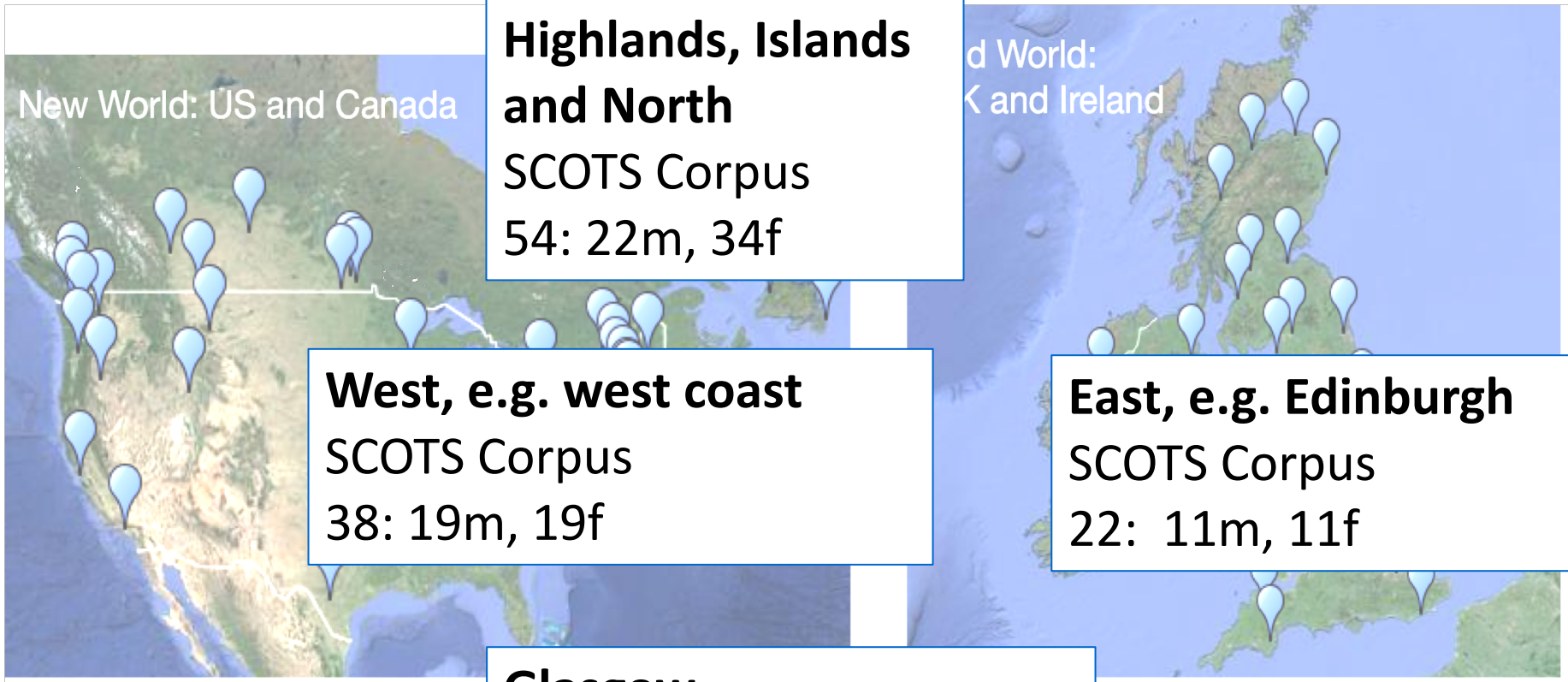
[www.google.com/maps/](http://www.google.com/maps/)

235 speakers



# SPADE

Sample for this study: Old World



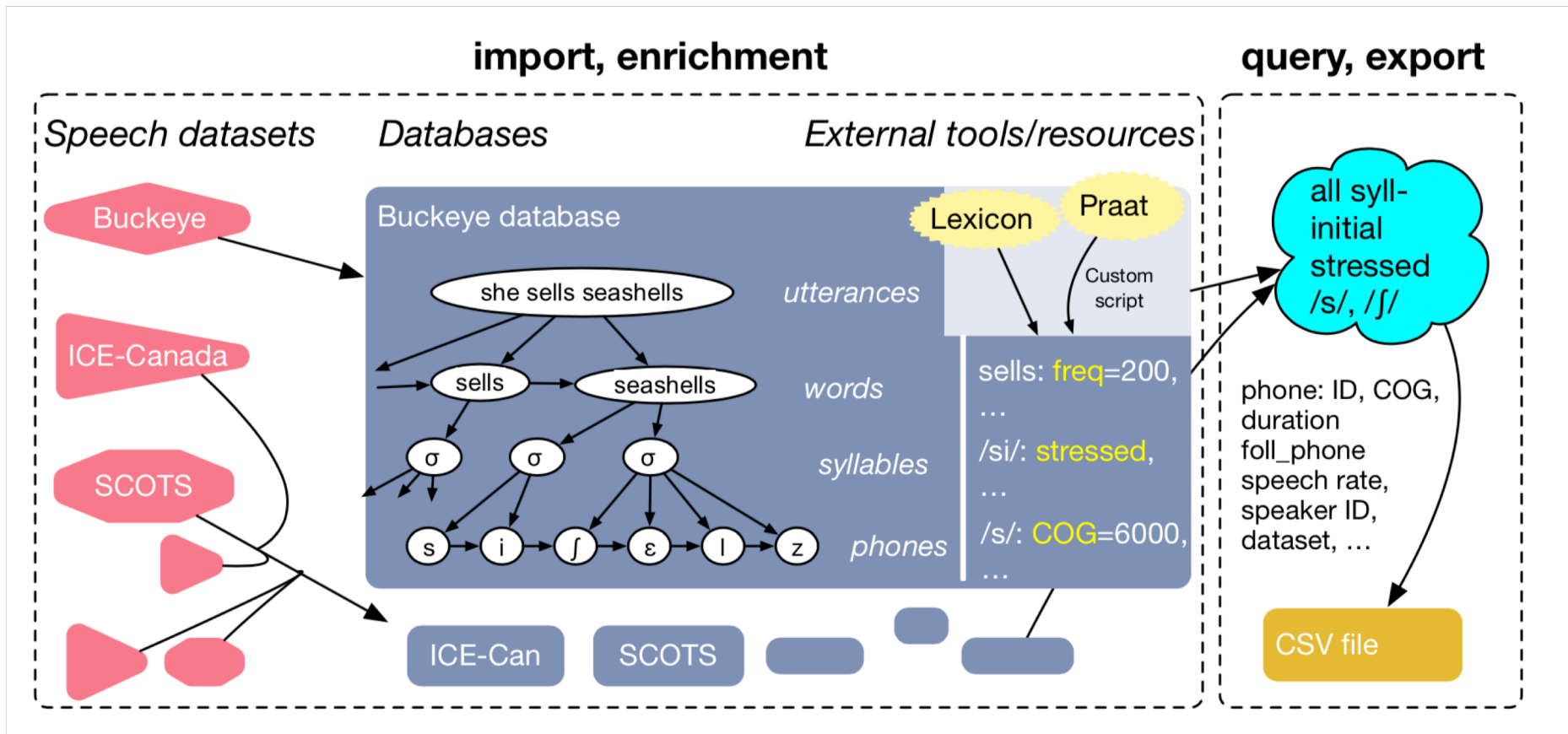
[www.google.com/maps/](http://www.google.com/maps/)

185 speakers

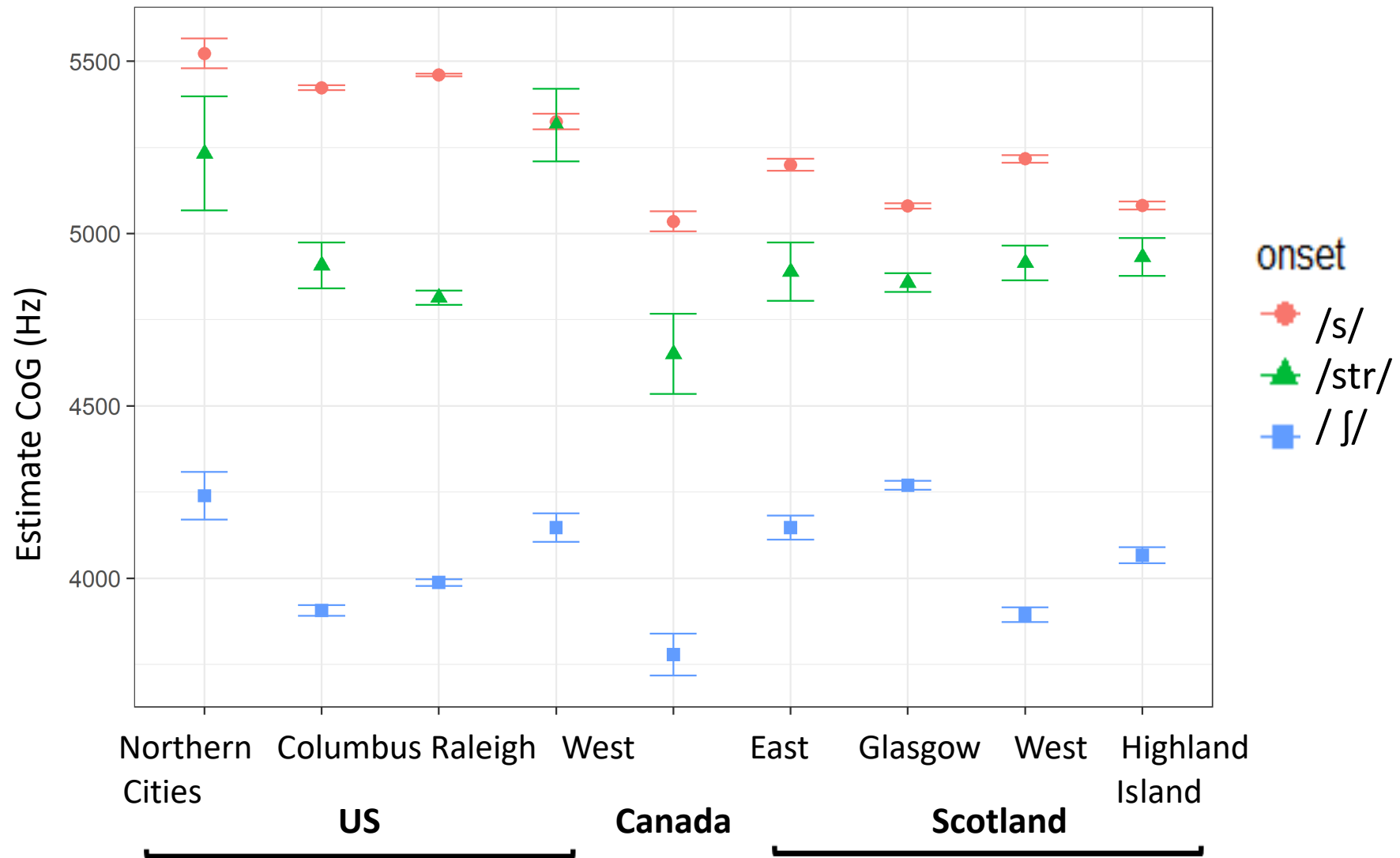
# Data

- All instances of stressed, word-initial /s/
- Acoustic measures: **peak**, spectral Centre of Gravity (**CoG**)
  - 1-16 kHz
  - Middle 50%
- Data cleaning
- N = 76,440
- Prediction: **/s/ > /str/ > /ʃ/**

# ISCAN usage

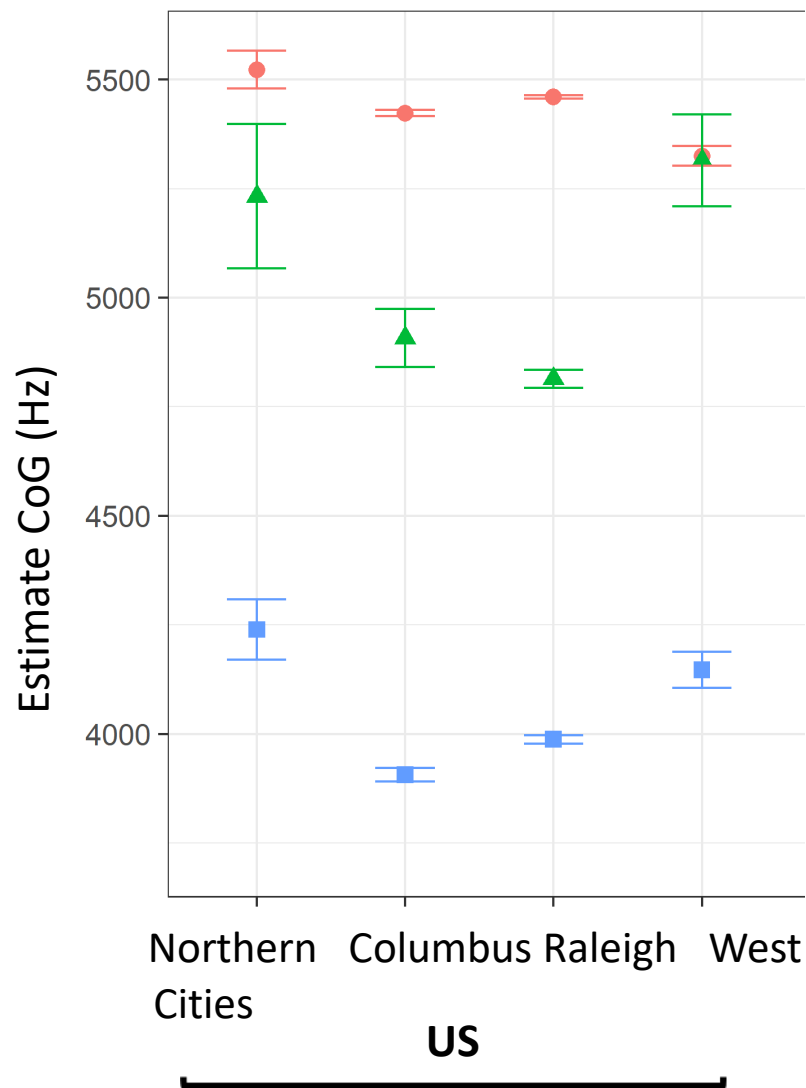


# Results



N = 76,440

• /str/ shows substantial variation across dialects



In **US** dialects, large differences in lowering of /str/ with respect to /s/

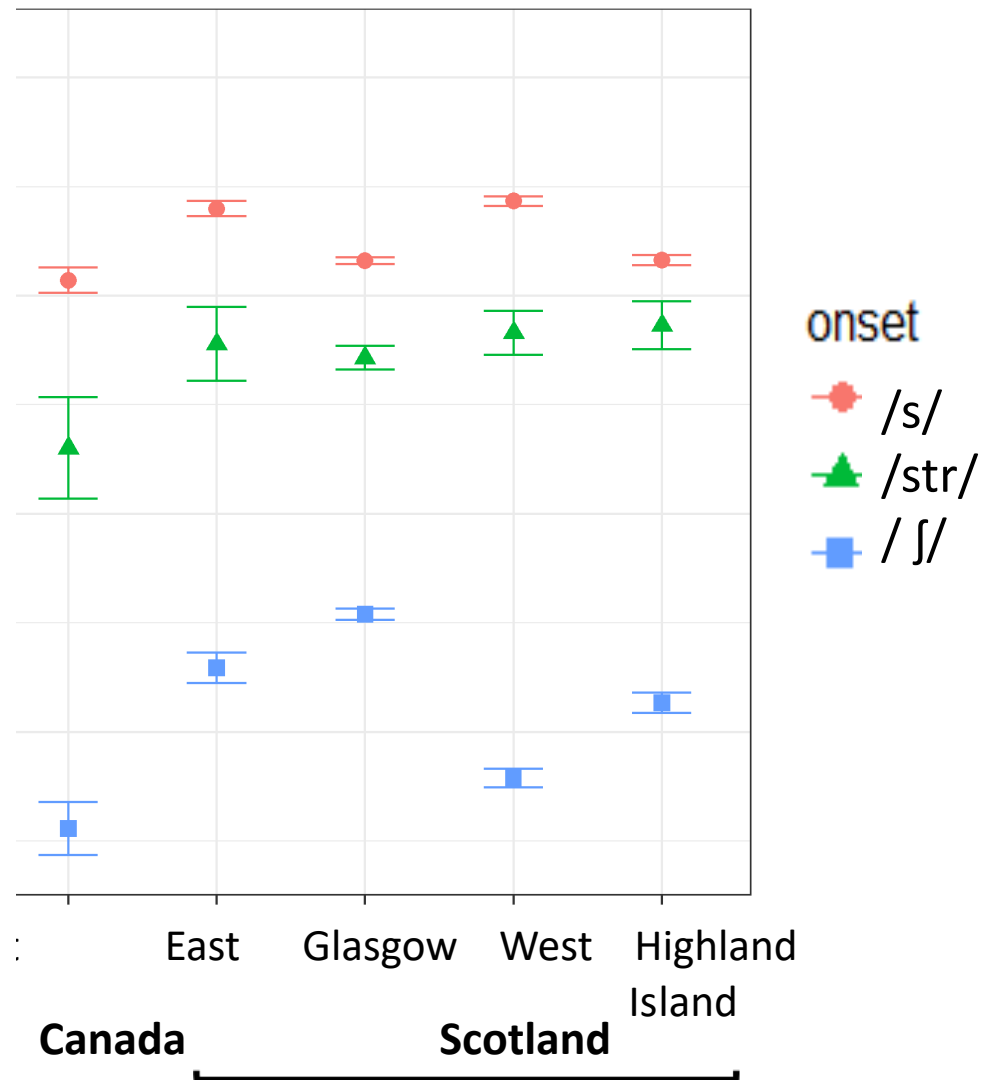
onset

- /s/
- ▲ /str/
- /ʃ/

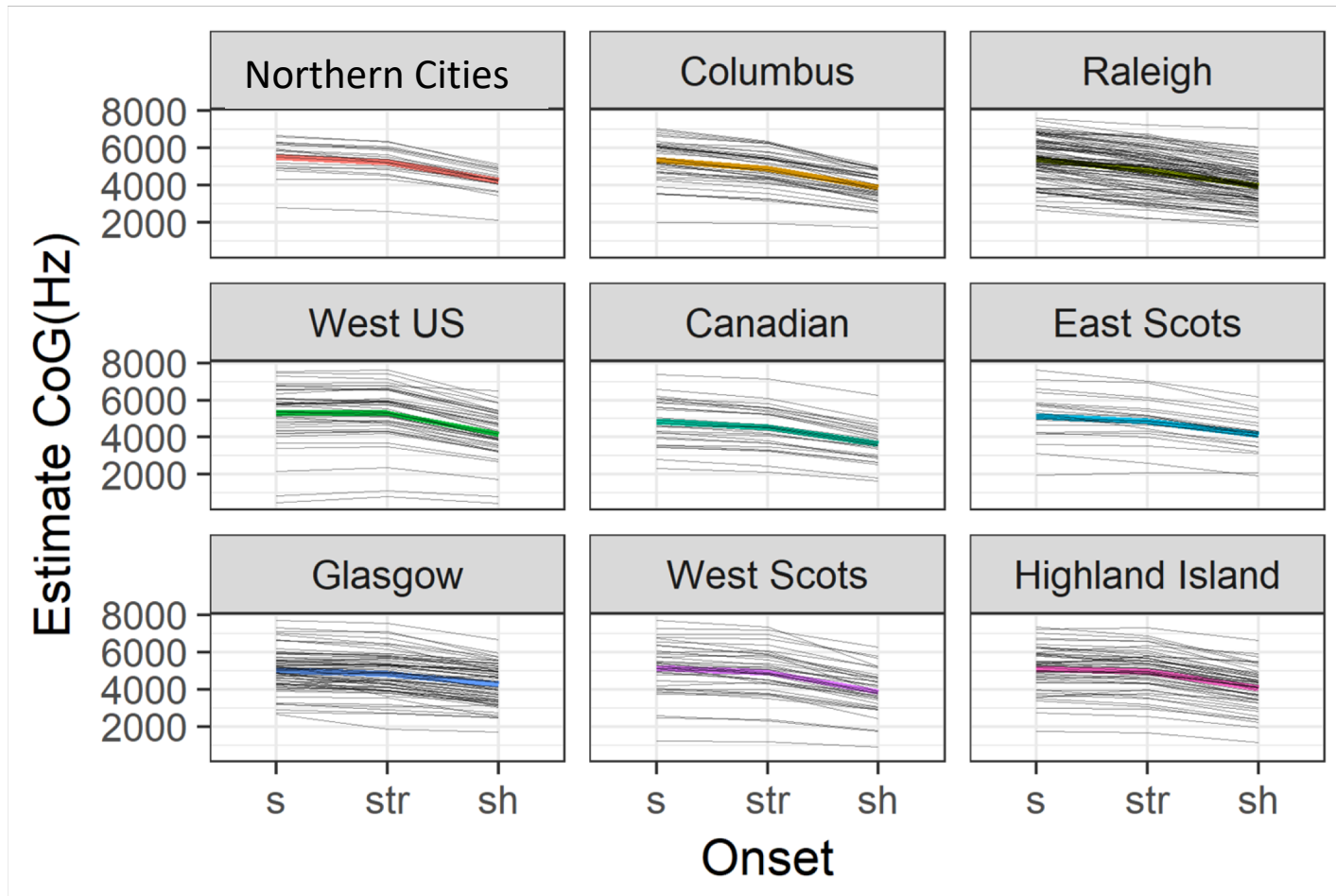


In **Scottish** and **Canadian** dialects, smaller differences between /str/ and /s/

/s/ is lower in frequency overall.



# /s/-retraction in individuals within dialects

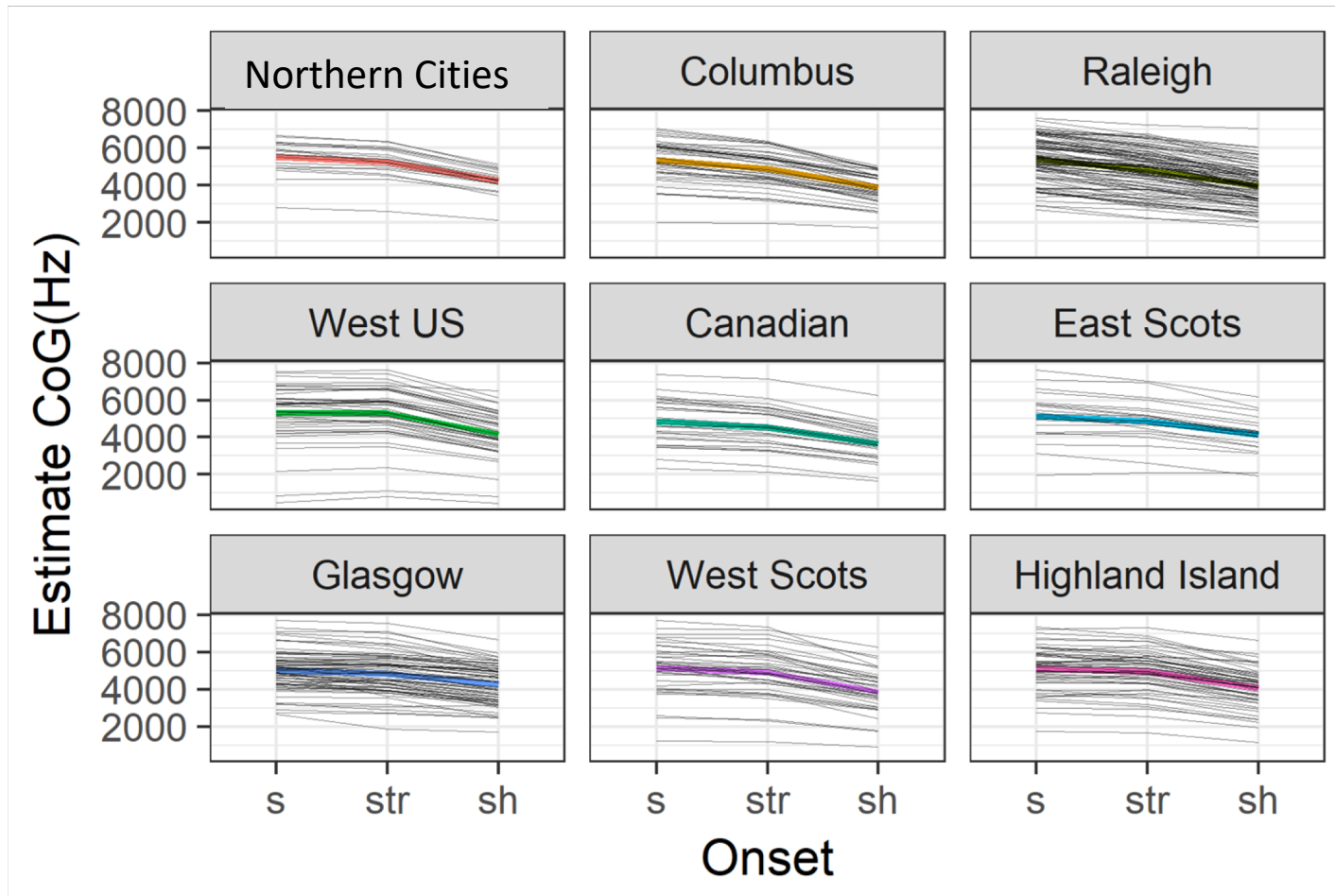


N = 76,440

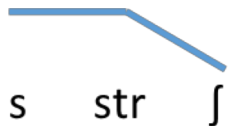
s str ʃ

retracting pattern shared by many Raleigh and Glasgow speakers.

# /s/-retraction in individuals within dialects

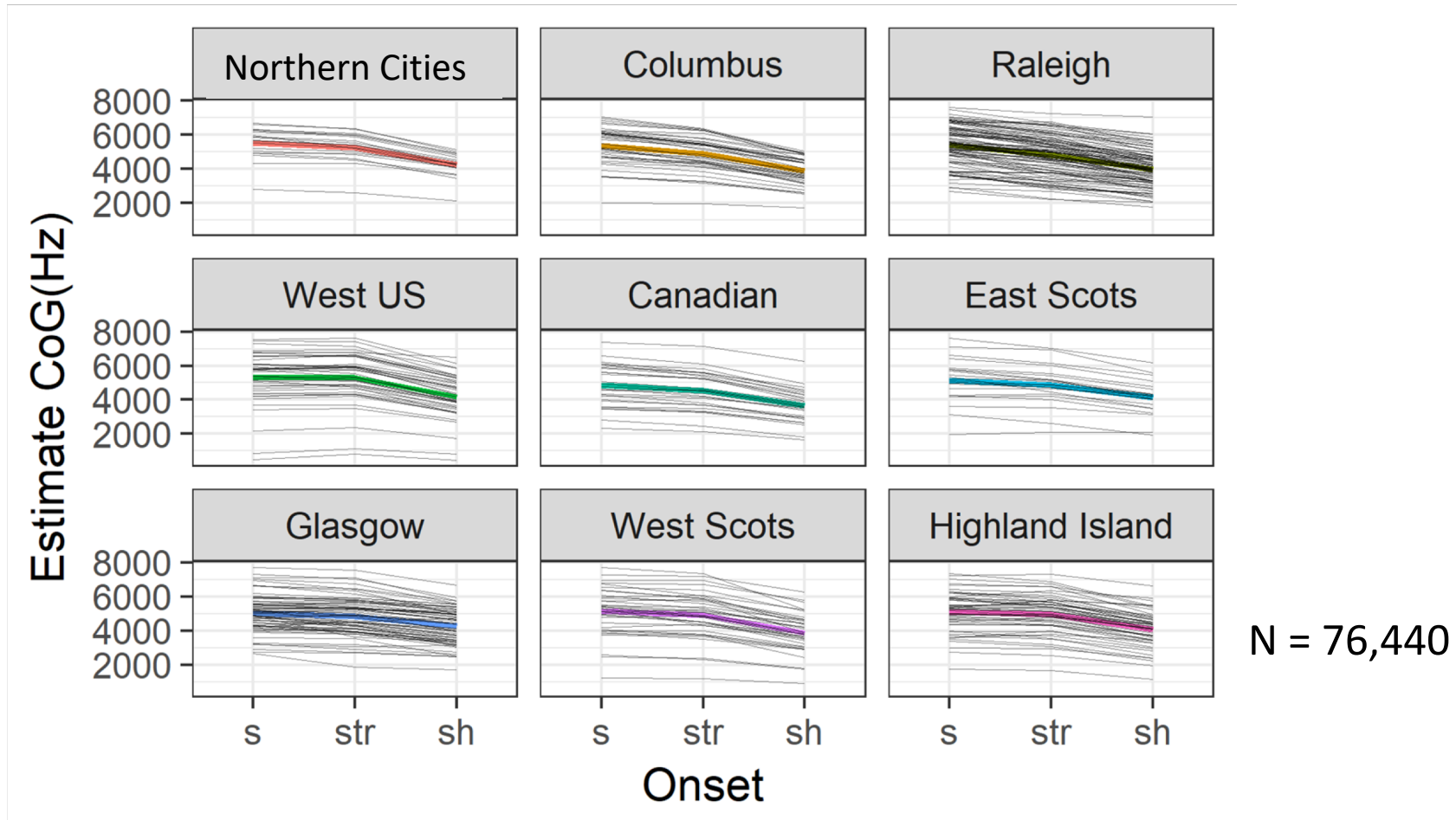


N = 76,440



non-retracting pattern in individual speakers  
in West US and Northern Cities dialects

# /s/-retraction in individuals within dialects



Both retracting and non-retracting patterns seen in some individuals in all dialects.

# Discussion

- Q1: what is evidence for /s/ retraction across English?
  - Some dialects show retraction of /str/
  - Large differences by dialect and by country
  - Impression of “/s/-retraction” depends on which dialects are considered
- Q2: continuum vs. dichotomy in /s/-retraction
  - ?
- Scaling up analysis across dialects, with consistent measures, allows identification of new patterns

# Study 2: vowel formants

Mielke et al. *Proc.  
ICPhS 2019*

- Influential hypothesis from sociolinguistics:  
(Labov, 1994)
  - **Intraspeaker variation** in vowel production ~ same axis as **diachronic change** in community
- Intuitively plausible, but unchecked
- (go to [poster..](#))

# Thanks!

- SPADE Team, especially
  - Jane Stuart-Smith, Michael McAuliffe, James Tanner, Vanna Willerton, Jeff Mielke
- MCQLL lab RAs
  - Michael Goodale, Arlie Coles, Elias Stengel-Eskin
- Funding
  - Digging Into Data, SSHRC, NSERC

# Questions