# Filtering and annotating web speech data

David Lutz
Parry Cadwallader
Mats Rooth

Cornell University

# Research with Web Speech Data

1. Locate

2. Collect

3. Filter

4. Annotate

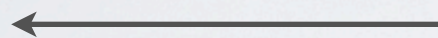5. Analyze

# Web Speech Data

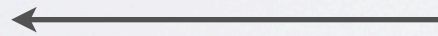1. Locate   &larr;

2. Collect   &larr;

3. Filter

4. Annotate

5. Analyze

Howell and Rooth 2009;
Howell 2012

# Web Speech Data

1. Locate

2. Collect

3. Filter ⟵——————————

                          Today: ezra

4. Annotate ⟵——————————

5. Analyze

# Locating and collecting tokens

- Content providers offer "media search"

- Many results, often with transcripts

- Results come from text search of ASR-generated transcripts

MLB Rule 5 draft and gives a recap of t

**in my mind** found at 7:30

**NFL Sunday breakdown: Pats off**

🔊Audio | Sun, 16 Dec 2012
Dale, Chris, Matt and Kevin take a look
that unit, including: Patrick Willis, Aldor
how the Patriots will attack this elite de
Rob Gronkowski.

**in my mind** found at 5:15

**Boomer Esiason on Patriots-49er
and the Sandy Hook shooting**

🔊Audio | Mon, 17 Dec 2012
Dennis and Callahan recap the Pats' nea
where things went so wrong early and s
discuss the tragedy in Newtown and Bo
Josh Brent on the sidelines after his DU

**in my mind** found at 12:21

# Locating and Collecting tokens

- Howell and Rooth: tools to search and download results

- Yielded hundreds of hits for target phrases

- Many of these hits were false positives

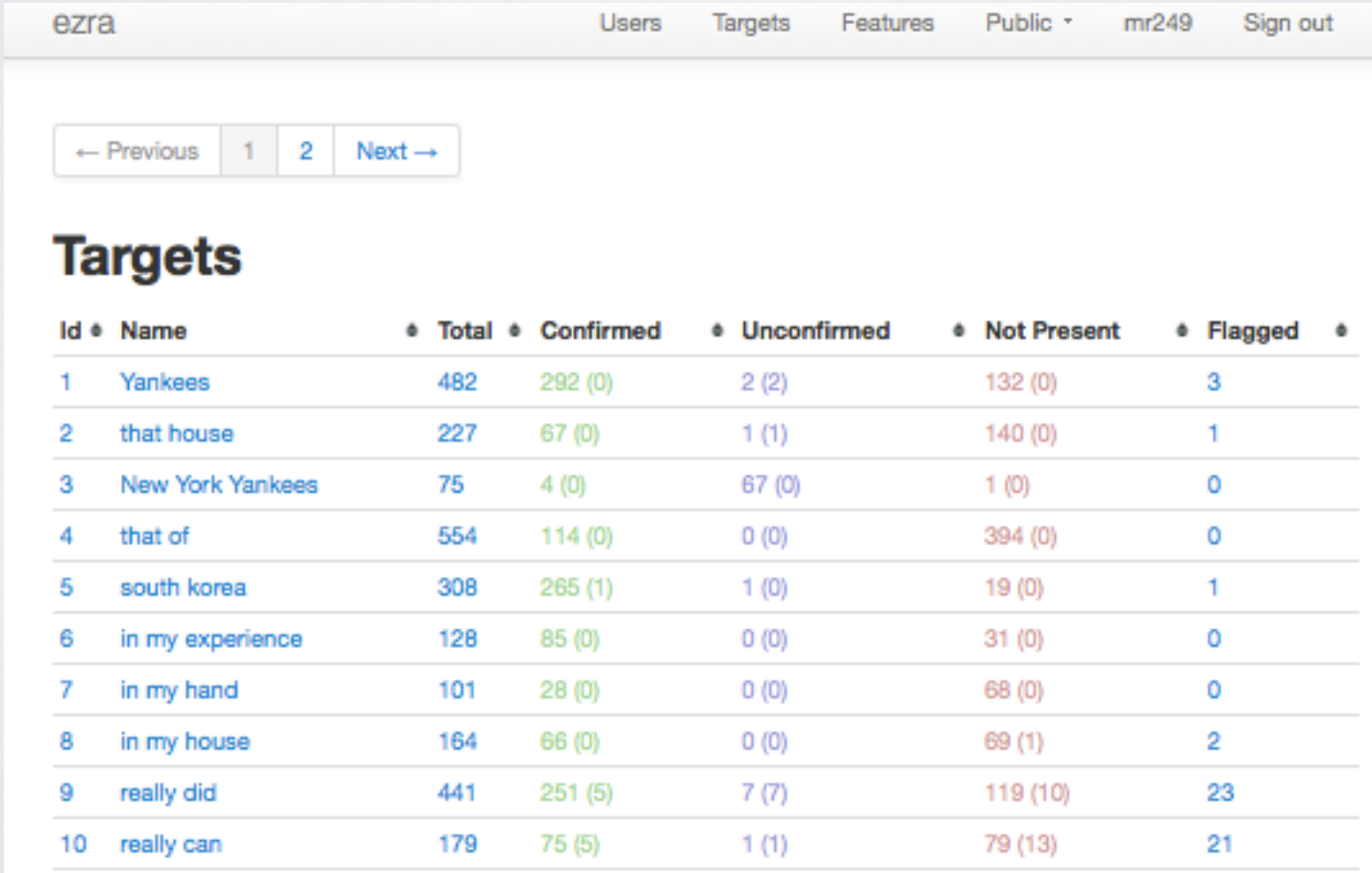- Those that weren't needed annotation

# Filtering and Annotation

1. Separate the true tokens from the false positives / duplicates

2. Mark the token location and reasonable clip boundaries

3. Add or correct the clip transcript

4. Target-specific annotation tasks

Much time was spent repeating the same mindless tasks

# ezra

Work is organized around **targets** and **features**.

**targets** are words or phrases of interest to the researcher.

# ezra

Work is organized around **targets** and **features**.

**features** are annotation tasks for a particular target.

## 18. some people

Total: 448
Confirmed: 383
Unconfirmed: 0
Not Present: 56

### Features

| Name | Created by | Number of targets | Instructions |
|------|-----------|-------------------|--------------|
| focus (Edit) | 4 | 3 | Is the target word in the ngram (e.g. "South" in "South |
| Phonological phrase position (Edit) | 9 | 2 | PPhrase-initial should be marked if the phrase is at the beginning |

### Hits

| ← Previous | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | ... | 14 | 15 | Next → |

| Hit Number | Status | Flagged |
|------------|--------|---------|
| ✏ 4303 | Not Present | |
| ✏ 4304 | Confirmed | |

# Features

- Features are annotation tasks

- Associated with targets

- Annotators complete these tasks for each hit

- Feature values are stored with each hit

# ezra: filtering and annotation

# Workflow

1. Define a target and import data

2. Define features and associate them with target(s)

3. Filter and annotate hits

4. Export for analysis

# Two user classes

- Supervisors

  - Create targets and features

  - Import/export data

  - Monitor user activity

  - Supervise annotation

- Annotators

  - More limited privileges

  - Focus on filtering and annotation

# Web application

- Accessible from anywhere

    - Users need a modern browser and internet connection

    - Team members can be remote

- Standard technology

    - Ruby on Rails

    - HTML5 / Javascript

    - SQLite

- Works in modern browsers

# Ezra: benefits

- Web application

- Built on standard technology

- Two user classes

- Multiple targets/projects

- Flexible feature definition

- Simple interface

- Public site

- Efficiency

https://github.com/del82/ezra

# Ezra: future work

- More complete user auditing

- Integrate locating and collecting tokens: plugins

- Improve administrator interface

- Automatic duplicate detection

- Partnerships with content creators

- Crowdsourcing?

https://github.com/del82/ezra

# Thanks

Neil Ashton

Jonathan Howell

Michael Wagner

Ross Kettleson

Michael Schramm

Lauren Garfinkle