LanguageARC:

Using Citizen Science to Augment Sociolinguistic Data Collection and Coding Christopher Cieri, Jonathan Wright, James Fiumara, Alex Shelmire and Mark Liberman University of Pennsylvania, Linguistic Data Consortium

Corresponding: {ccieri, jfiumara}@ldc.upenn.edu, www.ldc.upenn.edu

Introduction

- Discipline-specific advances, exploiting improvements in computing, allow quantitative sociolinguistics to analyze more subjects, variables and situations.
- Forced alignment [1, 2], vowel formant extraction [3] and statistical tools [4] allow researchers to undertake investigation and re-analysis until recently impossible [5]. Classifiers for other variables are emerging [6, 7] and search-guided coding [8] alone provides faster access to tokens where no classifiers exist.
- Unfortunately, we continue to suffer from a paucity of recorded observations of languagerelated behavior and of judgements which cannot be accelerated by purely technical mechanisms.
- Recording new data for under-represented communities and situations, transcription [9] for most languages [10, 11], evaluative judgements of the kind used, e.g., in perceptual dialectology [12] and matched guise tests require human input which is relatively costly and time consuming to obtain.

Novel Incentives and Workflows

- NSF sponsored NIEUW (NSF CRI CI-NEW #1730377) augments existing data sources by offering novel incentives -- opportunities to learn, compete, promote a linguistic variety and make real contributions to science -- on a significantly larger scale with less effort.
- Similar approaches proven successful in other fields such as the > 443 million judgments submitted by > 1.9 million Zooniverse [13] volunteers
- Our own language ID game, NameThatLanguage.org has collected 382,014 HITs from 31,368 players in < 1 year.

Citizen Linguists

- LanguageARC is NIEUW's Citizen Linguist portal designed specifically to collect linguistic data and judgements supporting multiple projects and tasks.
- Underlying toolkit already used in > 100 collection and coding tasks that yielded > 1,000,000 judgements
- Simplified for researcher use
- Tasks short enough to complete during breaks, commuting, etc. and many require nothing more than a smart phone.
- Project Builder allows researchers to deploy new tasks <1 hour, given a design and appropriately formatted data.
- Example Uses: eliciting word, sentence, passage reading; DiaPix and Map Task speech; silent movie narration and judgements such as language, speaker and dialect identification, transcription, grammaticality, and variable coding

Structure





| ID | Text |
|----|-------------|
| 1 | Rainbow.txt |
| 2 | NWS.txt |
| 3 | Stella.txt |
| 4 | Wolf.txt |
| 5 | Arthur.txt |
| 6 | Comma.txt |

Comma.txt Grandfather.txt

Project Builder

7

Ne

Ne

Ne

То

Allow

| N | Project | | | |
|--------|---|---------------------|--------------------------|--|
| | Name | Title | Subtitle | |
| | Pitch | News / Blog URL | Project Image | |
| | | | | |
| | Forums | Announcements | General Discussion | |
| | | Questions | Help & Technical Support | |
| | | | | |
| | Research Team | Name | Title | |
| | | Description | Photo | |
| | Partners | Name | LIRI | |
| | 1 artifers | Badao | ORE | |
| | | Dauge | | |
| w Task | | | | |
| | Name | Title | Image | |
| | Description | Tutorial | Reference Guide | |
| | | | | |
| | Assignment Order | In Order | Random | |
| | Coverage | Within Contributors | Across Contributors? | |
| | Forum | General Discussion | | |
| | Detect | | | |
| N | Dataset | Description | | |
| | Name | Description | | |
| | Manifest File Data Files | | | |
| | Randomize Manifest Item Order? Yes No | | | |
| ы | | | | |
| | Instructions | | | |
| | Media Type: Text Audio, Image Video, None, Media Column | | | |
| | Language Limit Language Selection=>Languages | | Languages | |
| | | | | |
| | Item ID | | | |
| | Item Specific Text | Primary Text/Label | Secondary Text/Label | |
| | Response Audio | Level Test? Yes No | Level Meter? Yes No | |
| | Response Text Yes No Text/Label | | | |
| | Judgment Buttons | Multiple Choice | | |

Skip? Yes No

Report Bad Item Yes No

Passage Name

Rainbow

NWS

Stella

Wolf

Arthu

Comma

Grandfathe



Projects/Tasks



[14] Z. Boyd, Z. Elliott, J. Fruehwald and L. Hall-Lew, "An Evaluation of Sociolinguistic Elicitation Methods (2015)," in Proceedings of ICPhS, 2015.

Soft Launch:

https://www.languagearc.org

