

Closer Still to a Robust, All Digital, Empirical, Reproducible Sociolinguistic Methodology

Christopher Cieri, Stephanie Strassel

University of Pennsylvania, Linguistic Data Consortium

ccieri@ldc.upenn.edu, strassel@ldc.upenn.edu

www.ldc.upenn.edu/Papers

History

- ◆ 1963 Quantitative study of variation and change in the speech community has been intensively corpus based since inception
- ◆ 1971 Montreal Group began to create first computer based corpus for speech community study
- ◆ 1999 Gregory Guy convened a workshop on publicly available corpora, invited us to present on LDC corpora of potential use to sociolinguistics
- ◆ 2001 presented on corpus based sociolinguistics, our DASL project and the –t/d deletion study
- ◆ 2002 presented with William Labov on the SLx Corpus of classic sociolinguistic interviews and the DASLTrans
- ◆ 2003 organized Workshop at Penn of robust sociolinguistic methodology
- ◆ 2007 Malcah Yaeger-Dror convened workshop, invited Reva Schwartz, and MIT-LL and LDC to present on transcription practice and Phanotics project
- ◆ 2009 today we are very close to the realization of this ideal

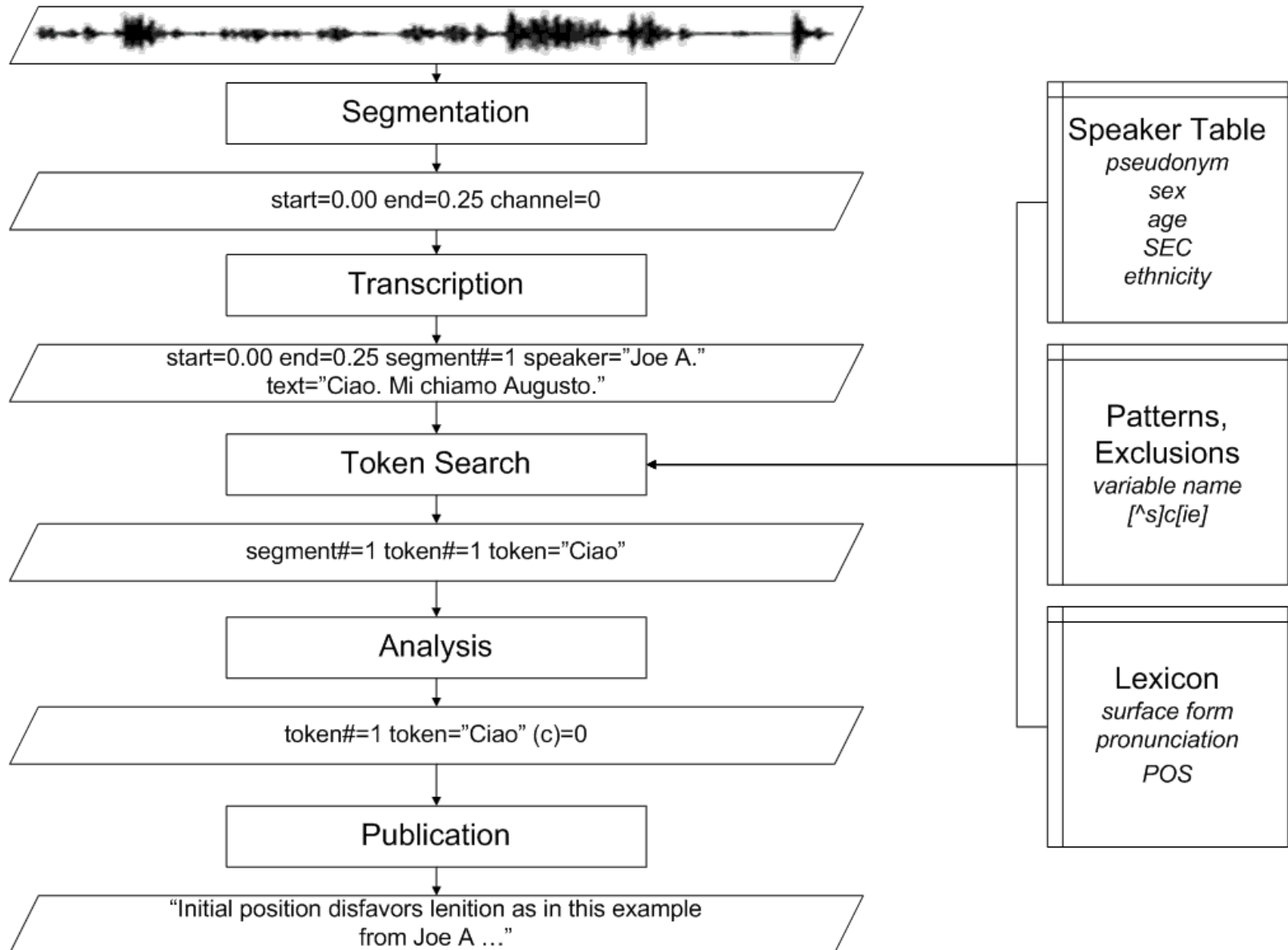
Vision

- ◆ raw data – text, audio, video – is digital as are annotations, specifications
- ◆ transcripts other annotations are linked back to the original, raw data
 - time stamped for speech, linked via word offsets for text
- ◆ raw data or transcript proxy is computer searched for target variables
 - lexicons, speaker tables, other data external to recordings consulted as needed
- ◆ coding decisions are still made by humans
 - though the potential for partial automation exists
- ◆ variables, coding practice described to permit replication by others on the same or comparable data
- ◆ coding strings, examples in a paper, dots on a scatter plot or tracked backed to original recordings
- ◆ ideally data also publicly accessible.

Model



Model



Segmentation

- ◆ Virtually divide digital audio stream into manageable units
- ◆ Greatly facilitates downstream transcription, token retrieval, coding, analysis
- ◆ Can also indicate structural boundaries in recording
- ◆ Variable segment granularity to meet project needs
 - Maximum segment duration of 5-8 seconds makes downstream transcription and coding considerably more efficient
 - Sentence units (SU), breath/pause groups are convenient first-order units
 - Turns, discourse units, word, phones, etc. as optional second pass
- ◆ With right tools, SU or breath group segmentation can be performed in under 1.2x real time
 - Automatic segmentation, forced alignment with manual verification can also save time

Transcription

- ◆ Why a full transcription?
 - Index to speech, searchable
 - Provides stable basis for subsequent tasks
- ◆ Transcription specification to document conventions for orthographic representation
 - Use of standard orthography facilitates subsequent searching, retrieval of tokens, reanalysis
 - Specify treatment of common phenomena like disfluencies, non-standard forms, mispronunciations, transcriber uncertainty
- ◆ Transcription can be quite efficient given right tools combined with short audio segments

Comparison of Methods

| | <u>Quickest</u> | | | <u>Most Careful</u> |
|-----------------------------------|-----------------|---------------------------------|---------------------------------|------------------------------------|
| Segmentation | Automatic | Auto w/ verification | Manual | Manual w/ verification |
| Completeness | Content words | Add partial words, disfluencies | Add partial words, disfluencies | Add verification pass |
| Filled Pauses | Optional | Incomplete | Exhaustive | Exhaustive w/ verification |
| Disfluencies | None | Incomplete | Exhaustive | Exhaustive w/ verification |
| Transcriber Uncertainty | Flag and skip | Flag and best guess | Flag and best guess | Flagged best guess w/ verification |
| Feature Marking | None | Minimal | Full | Accurate, complete w/ correction |
| Speaker, Backgrnd Noise | None | Minimal | Exhaustive | Exhaustive w/ verification |
| Manual Passes | 1 | 1-2 | 2-3 | 4+ |
| Approx. Cost (x Real Time) | 5 x | 15 x | 25 x | 50 x |

Comparison of Methods

| | <u>Quickest</u> | | | <u>Most Careful</u> |
|-----------------------------------|-----------------|---------------------------------|---------------------------------|------------------------------------|
| Segmentation | Automatic | Auto w/ verification | Manual | Manual w/ verification |
| Completeness | Content words | Add partial words, disfluencies | Add partial words, disfluencies | Add verification pass |
| Filled Pauses | Optional | Incomplete | Exhaustive | Exhaustive w/ verification |
| Disfluencies | None | Incomplete | Exhaustive | Exhaustive w/ verification |
| Transcriber Uncertainty | Flag and skip | Flag and best guess | Flag and best guess | Flagged best guess w/ verification |
| Feature Marking | None | Minimal | Full | Accurate, complete w/ correction |
| Speaker, Backgrnd Noise | None | Minimal | Exhaustive | Exhaustive w/ verification |
| Manual Passes | 1 | 1-2 | 2-3 | 4+ |
| Approx. Cost (x Real Time) | 5 x | 15 x | 25 x | 50 x |

Approximately 10x including segmentation

Quick Transcription Example

XTrans: 571959_a.tdf

File Edit View Tools

Respondent: A good story is that uh
Respondent: when I was in high school that integrated
Respondent: the -- the schools
Respondent: I was taking algebra and there's a lot of stories about how I was mistreated and stuff.
Respondent: But I was having difficulty with -- with algebra.
Respondent: And I was sitting at the kitchen table trying to do my homework.
Respondent: And I said -- I got frustrated and said I just can't figure this out, I'm just --
Respondent: So my father said what's the problem, he came by, he said what's the problem?
Respondent: And I said ((it's this)) algebra, and he said well let me look at it. I said Dad this is algebra.
Respondent: They didn't even have algebra in your day.
Respondent: And uh -- and I went to sleep, I went to bed.
Respondent: And around four o'clock that morning he woke me up, he said
Respondent: come on son get up.
Respondent: And I said what -- what's wrong. He said let's talk about this algebra.
Respondent: He sat me at the kitchen table.
Respondent: And he -- we went over algebra. He taught me algebra.
Respondent: What he had done is sit up all night and read the algebra book.
Respondent: And then he explained the problems to me,
Respondent: so I could do them, and understand them.
Respondent: And to this day I live my life trying to be half the man my father was.
Respondent: Just half the man.
Respondent: And uh
Respondent: I would be a success if my children loved me half as much as I loved my father.

Interviewer (F)
Respondent (M)

VOS VAS SRT CLR
NSI ESIq MRG VRGq
LRS LAS LAG

2:38.1585 2:20.6893 2:25.3747 4.6855

2:05.0 2:10.0 2:15.0 2:20.0 2:25.0 2:30.0 2:35.0 2:40.0 2:45.0 2:50.0 2:55.0 3:00.0

1.00

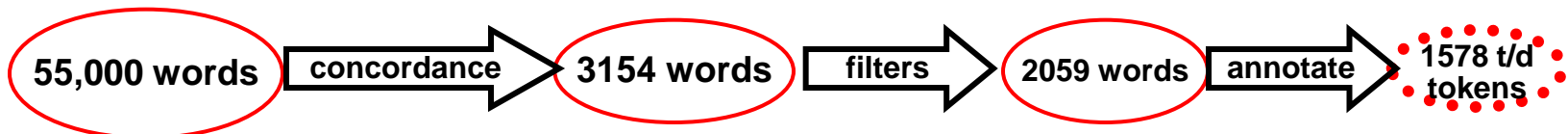
ready

<http://www ldc.upenn.edu/XTrans>

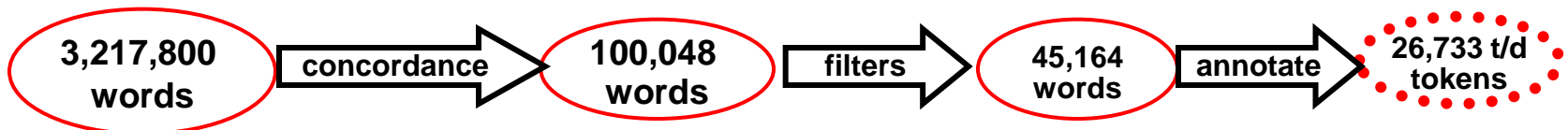
Token Selection

- ◆ Selection of tokens for analysis can be automated to large extent
 - Concordance to identify tokens of interest
 - Using string matching, regular expression queries
 - Filters to remove additional non-tokens
- ◆ More robust than manual selection, which might miss or overlook tokens
- ◆ Implemented in DASL t/d study

TIMIT Corpus (LDC93S1)



Switchboard Corpus (LDC97S62)



Coding Spec Challenges

- ◆ Difficulty of achieving perfectly explicit guidelines
 - Even when working on well-studied variable
- ◆ In DASL t/d deletion study, goal was to investigate comparability of corpus-based approaches with previous studies involving sociolinguistic interview data
- ◆ But previous t/d coding specs not typically published
 - Had to resort to personal communication with authors, detective work, reverse engineering from results
- ◆ Variation in coding for some factor groups inhibits direct comparison of results
 - Morphological factors, e.g. passives ("I was frightened")
- ◆ Some categories unmentioned - how were these coded?
 - Nasal flaps? Glottalized segments? What constitutes a pause?

Coding Spec Best Practices

- ◆ Formal annotation/coding specifications promote coder reliability and direct comparison of results
- ◆ Developed iteratively over several rounds of pilot labeling including analysis of inter-coder reliability, via (double-blind) dual coding
 - Consider removal, merging of rules/categories with low consistency
- ◆ Written guidelines include
 - Title, date, version number
 - Introduction with framing/contextual info and general description of rule syntax
 - Screenshots of annotation/coding interface
 - Multiple examples for each rule
 - Including some difficult cases as well as counter-examples
 - Embedded sound files to illustrate application & non-application of rule
 - Appendix, glossary
 - Rules of thumb to promote consistent labeling
 - Can't tell, difficult decision flags
- ◆ (Link to) guidelines published along with results

Coding

- ◆ Careful data preparation (segmentation, transcription) and pre-selection of all candidate tokens enables efficient coding
- ◆ "Regions of interest" already identified
- ◆ Attention directed at a single task: how is this variable realized in this batch of tokens
- ◆ Some customization of coding tools can increase efficiency further still

DASL t/d Coding Tool

DASL - Project: t/d Deletion - Netscape
File Edit View Go Communicator Help

Welcome:
ccieri

Jump to:
Next Page
DASL Home
t/d Deletion Page

Data and Annotations for SocioLinguistics

| | | | | | |
|--|---------------------------------------|--|---------------|--------------------|-----------------------|
| Independent Variable File: /Shared/TDdeletion.tag | Token File: /Shared/TDdeletion.tok | Annotation File: /ccieri/TDdeletion.ann | Page: 1/83 | Tokens/Page: 25 | Total Tokens: 2059 |
|--|---------------------------------------|--|---------------|--------------------|-----------------------|

1. ... loved to chew on the old rag doll.
2055, Male, New York City, 25, White, Bachelor's Degree

| | |
|----------------|--|
| t/d: | <input type="radio"/> Untouched <input type="radio"/> Deleted <input checked="" type="radio"/> Retained <input type="radio"/> Unsure <input type="radio"/> NA |
| Morphological: | <input checked="" type="radio"/> Monomorpheme <input type="radio"/> Irregular_Past <input type="radio"/> Regular_Past |
| Preceding: | <input type="radio"/> Stop <input checked="" type="radio"/> Lateral <input type="radio"/> Rhotic <input type="radio"/> Alveolar_Nasal <input type="radio"/> Other_Nasal <input type="radio"/> Alveolar_Fricative <input type="radio"/> Other_Fricative |
| Following: | <input type="radio"/> Obstruent <input type="radio"/> Lateral <input checked="" type="radio"/> Rhotic <input type="radio"/> Clustering_Glide <input type="radio"/> Other_Glide <input type="radio"/> Vowel <input type="radio"/> Pause |
| comments: | vocalized l |

2. ... those who te
2055, Male, New Y

WaveView 1.1 - Netscape
File Edit View Go Communicator Help

WaveView version 1.1 Corpus: timit, Filename: /train/dr6/mabc0/sx331.wav
Time: 0.0sec D: 0.32717142sec L: 1.42765714sec R: 1.75482857sec

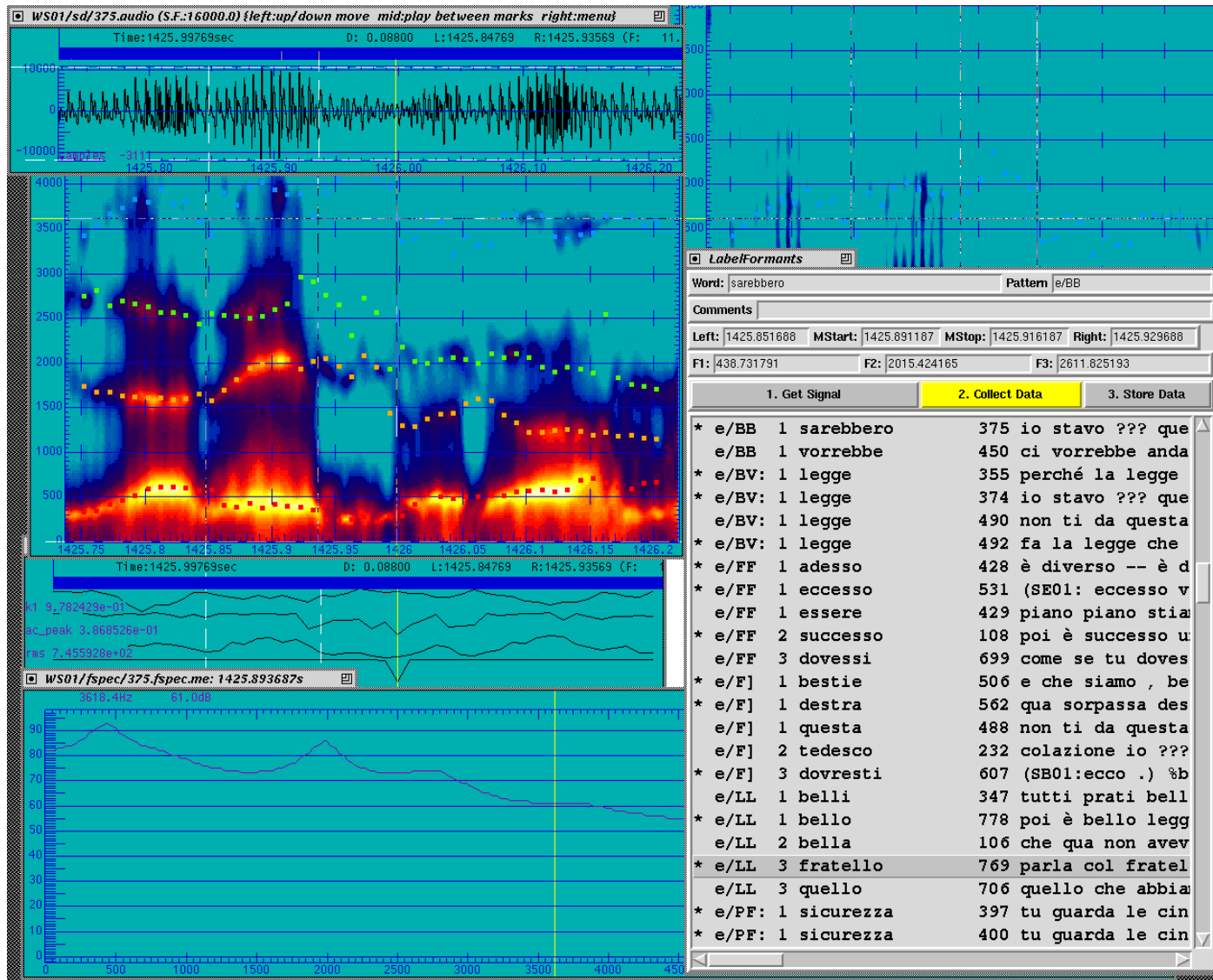
Zoom In Zoom Out Zoom Full Out Bracket Mark Window Forward Window Backward
Stop Play Play Mark Play Window Play All

SPAAT (Super Phonetic Annotation & Analysis Tool)

- ◆ One variable, one ROI at a time
- ◆ Average of 250 judgments/hour, up to 400+ for experienced labelers



Formant Analysis



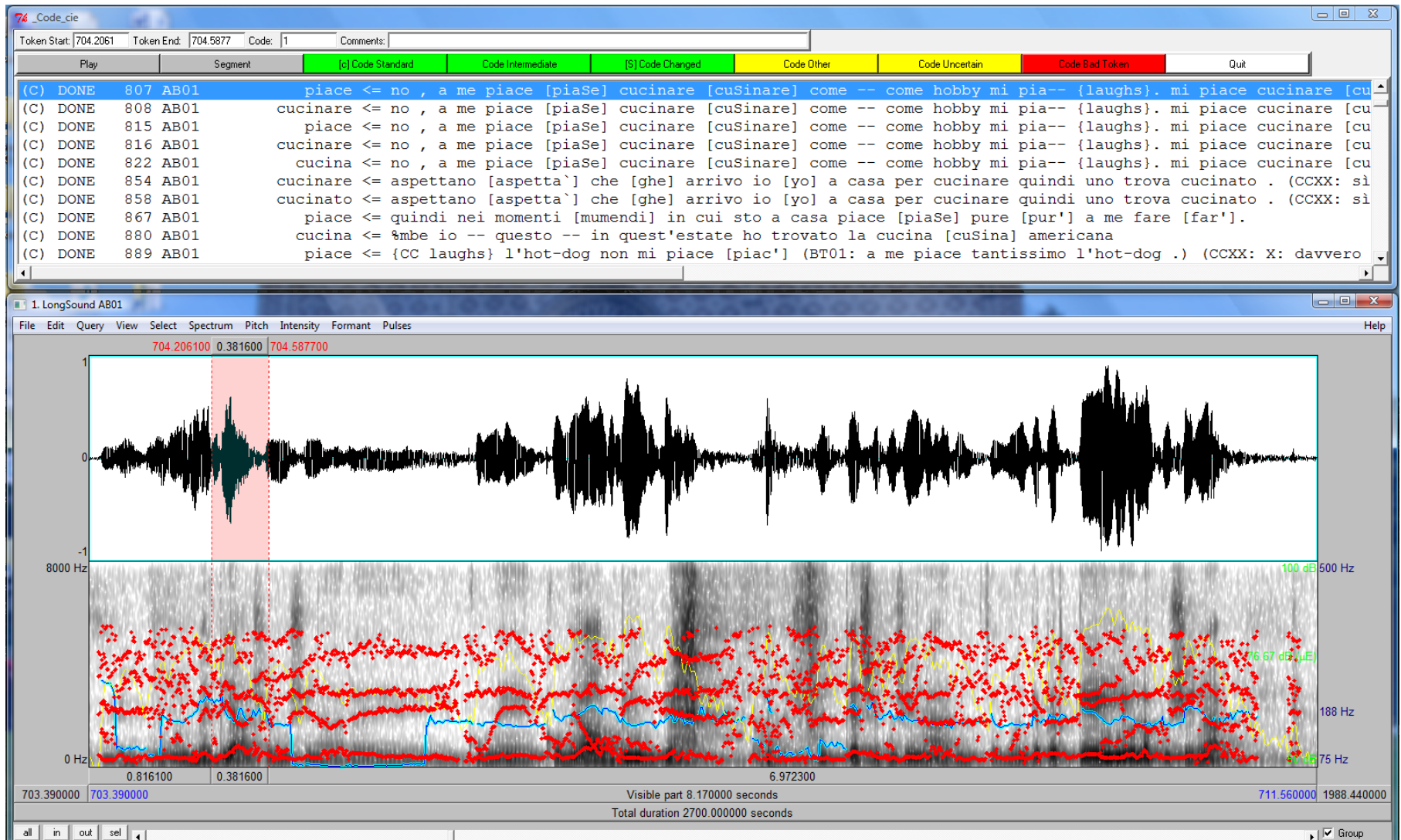
Token Selection

Vowel
Segmentation

Identification of
central tendency
of word stressed
vowel

Hand checking
of formant
tracker values
for F1 and F2

Data Management



Vision

- ☑ raw data – text, audio, video – is digital as are annotations, specifications
- ☑ transcripts other annotations are linked back to the original, raw data
 - Xtrans, Praat, various Concordancers
- ☑ raw data or transcript proxy is computer searched for target variables
 - Ottawa Workshop, Montreal Project, SPAAT
- ☑ coding decisions are still made by humans
 - though the potential for partial automation exists
 - Yuan's Forced Aligner, Evanini's formant extractor
 - Other HLTs: ASR, Universal Phonetic Decoders, Energy Detectors, POS Taggers
- ☑ variables, coding practice described to permit replication by others on the same or comparable data
 - DASL Project, SLx,
- ☑ coding strings, examples, points on a graph tracked to original recordings
 - HTML <a> tags, Stefan Dollinger's Bank of Canadian English, Tom Veatch's 1993 dissertation
- ☑ ideally data also publicly accessible
 - Michelle Minnick-Fox, Nationwide Speech Project, NECTE Corpus

Journal of Experimental Linguistics

[eLanguage](#) [Home](#) [About](#) [Log In](#) [Register](#) [Search](#) [Current](#) [Archives](#)

[Announcements](#)

[Home](#) > [Journal of Experimental Linguistics](#)

Journal of Experimental Linguistics

JEL is an interdisciplinary journal of reproducible research on topics related to speech and language. Regular publication will begin towards the end of 2009. For further details, see the announcements below.

The *Journal of Experimental Linguistics* is part of the Linguistic Society of America's eLanguage initiative. Like the rest of eLanguage, JEL is an Open Access online journal.

JEL is a linguistic "journal of reproducible research", that is, a journal of reproducible computational experiments on topics related to speech and language. These experiments may involve the analysis of previously published corpus data, or of experiment-specific data that is published for the occasion. Other relevant categories include computational simulations, implementations of diagnostic techniques or task scoring methods, methodological tutorials, and reviews of relevant new publications (including new data and software).

In all cases, JEL articles will be accompanied by executable recipes for recreating all figures, tables, numbers and other results. These recipes will be in the form of source code that runs in some generally--available computational environment.

Although JEL is centered in linguistics, we aim to publish research from the widest possible range of disciplines that engage speech and language experimentally, from electrical engineering and computer science to education, psychology, biology, and speech pathology. In this interdisciplinary context, "reproducible research" is especially useful in helping experimental and analytical techniques to cross over from one subfield to another.

Publication is in online digital form only, with articles appearing as they complete the review process. A rigorous but rapid process of peer review, designed to take no more than 4-6 weeks from submission to publication, will be supplemented by a vigorously-promoted system for adding moderated remarks and replies after publication.

The editorial board, in alphabetical order, is Alan Black, Steven Bird, Harald Baayen, Paul Boersma, Tim Bunnell, Khalid Choukri, Christopher Cieri, John Coleman, Eric Fosler-Lussier, John Goldsmith, Jen Hay, Stephen Isard, Greg Kochanski, Lori Levin, Mark Liberman, Brian MacWhinney, Ani Nenkova, James Pennebaker, Stuart Shieber, Chinlin Shih, David Talkin, Betty Tuller, and Jiahong Yuan. Mark Liberman is the editor in chief.



