

# Massive Auditory Lexical Decision: Going Big in the Auditory Domain

Benjamin V. Tucker and Daniel Brenner

Department of Linguistics, University of Alberta  
{bvtucker, brenner}@ualberta.ca

Mental Lexicon 2016: Ottawa, 20 Oct. 2016

# Acknowledgements

- R. Harald Baayen, D. Kyle Danielson, Danielle Fonseca, Catherine Ford, Matt Kelly, Pearl Lorentzen, Filip Nenadic, Katelynn Pawlenchuk, Michelle Sims
- With funding provided by:



Conseil de recherches  
en sciences humaines  
du Canada

Social Sciences and  
Humanities Research  
Council of Canada

Canada



UNIVERSITY OF  
ALBERTA



## 'Megastudies' have several important advantages:

- statistical power
- minimization of strategic effects
- comprehensiveness
- multi-functionality
- complementing traditional small factorial experiments
- model development and testing

## Visual Lexical Decision

For example:

- the English Lexicon Project (Balota et al., 2007): 40,000 words and non-words
- the French Lexicon Project (Ferrand et al., 2010): 38,000 words and non-words
- the Dutch Lexicon Project (Keuleers et al., 2010): 14,000 words and non-words
- the British English Lexicon Project (Keuleers et al., 2012): 28,700 words and non-words

**The only Auditory Lexical Decision megastudy:** BALDEY (Ernestus Cutler, 2015):

- 5,541 words and 5,541 pseudo-words
- 10 female and 10 male listeners

**The only Auditory Lexical Decision megastudy** BALDEY (Ernestus Cutler, 2015):

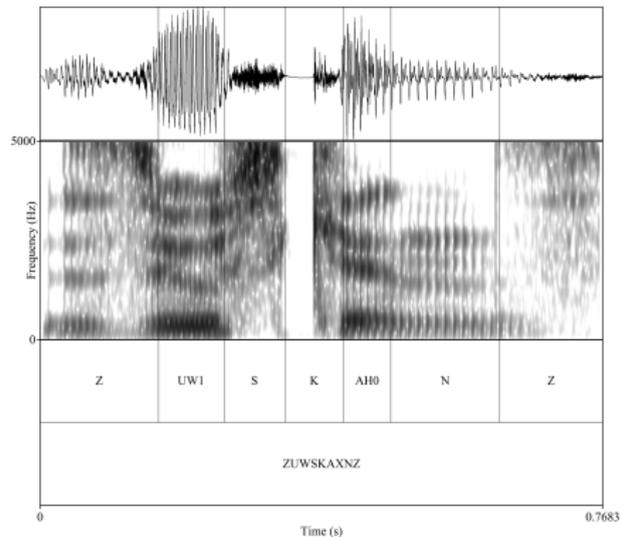
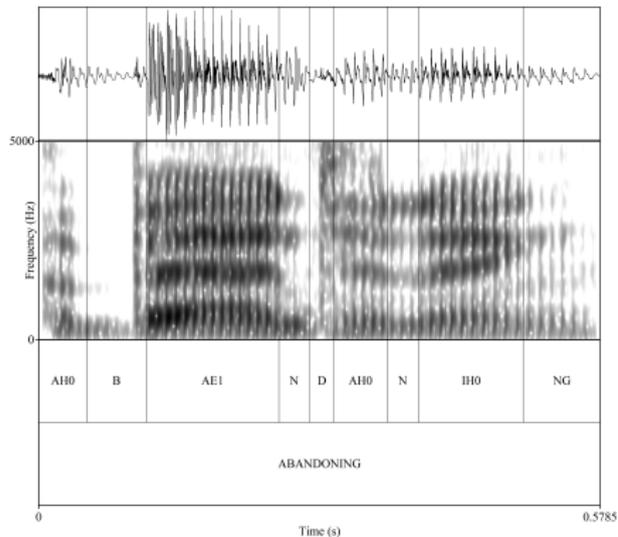
- 5,541 words and 5,541 pseudo-words
- 10 female and 10 male listeners

And now:

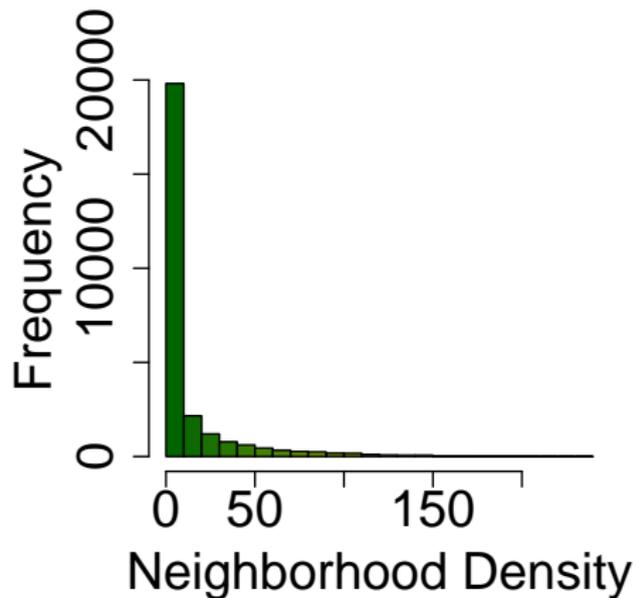
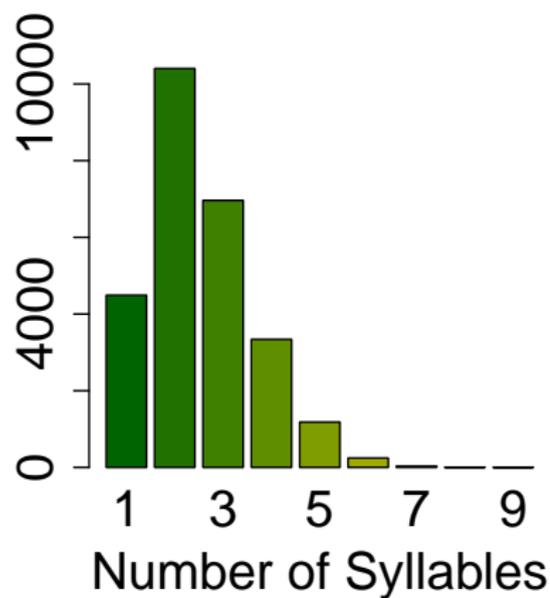
- MALD: Massive Auditory Lexical Decision

- Male Western Canadian English speaker (age 32) recorded in a sound attenuated booth
  - 28,511 words
  - 11,400 non-words (wuggy, Keuleers Brysbaert, 2010b)
    - Words and non-words are morphologically complex
    - 1000 compound words and non-words
- About 2000 words/day or 800 non-words/day
- Items were extracted and mispronunciations removed leaving:
  - 26,800 English words
  - 9,600 pseudo-words
- All items provided with segmental level mark-up (p2fa, Yuan Liberman, 2008)

# Item markup



# Item information

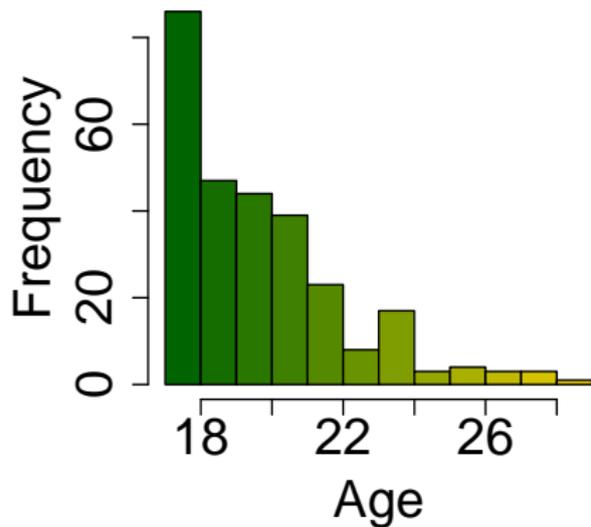


# Procedure

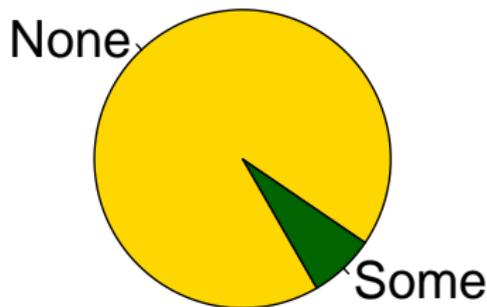
- Hearing screening
- Auditory Lexical Decision task
- Session lasts approximately 25min
  - Goal: At least 4 responses per word (400 words/400 pseudowords per experiment)
- Participants could participate in up to three sessions
  
- 232 monolingual Canadian English participants
- 285 total experimental sessions

**228,000 total button presses**

## Subject Ages

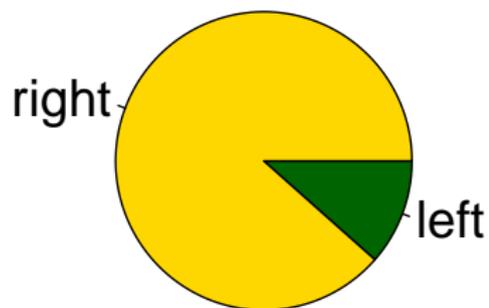


## Hearing Loss

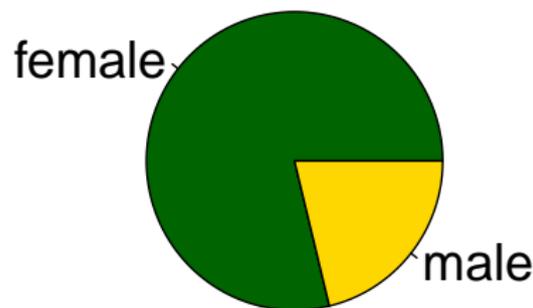


# Participants

## Subject Handedness



## Subject Sex



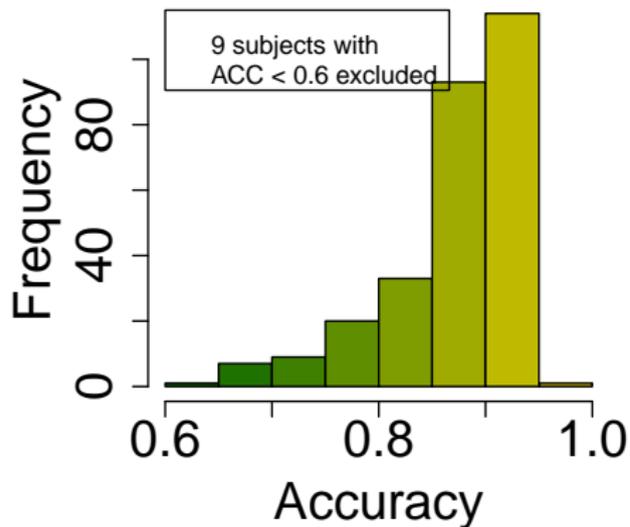
## Independent Variables (already in the data file):

- Word Duration
- Education
- Neighborhood Density
- Frequency (COCA, COCA Spoken, Google nGram)
- Non-word characteristics (e.g. phonotactic probability, Phonological Neighborhood Density)
- Age
- Sex
- Handedness
- Word Run Length

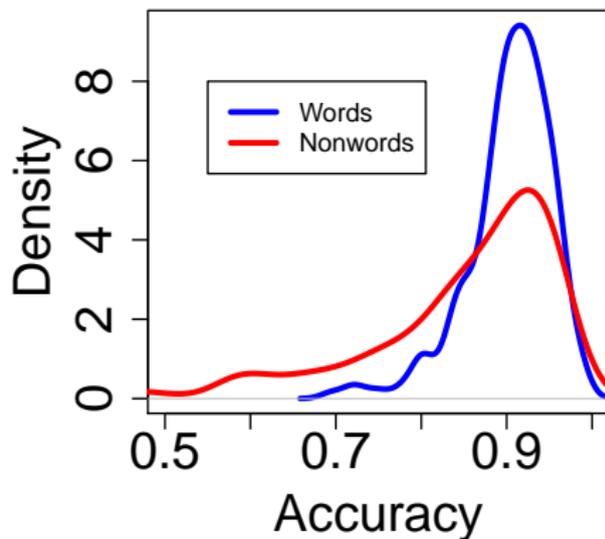
## Dependent variables:

- Acoustic characteristics
- Response Latency
- Accuracy

## Subject Accuracy

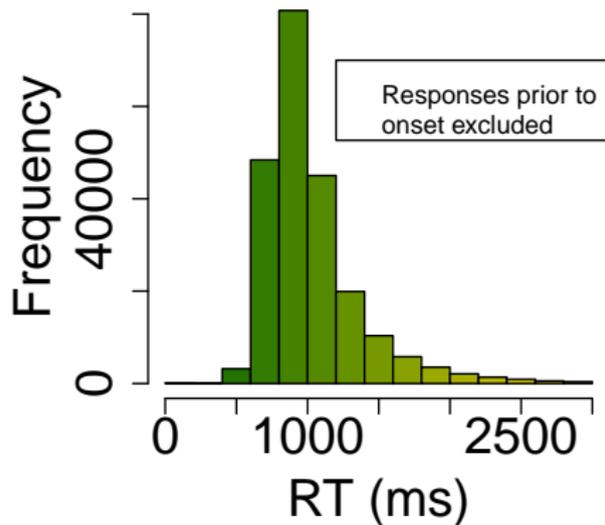


## Word vs. Non-word

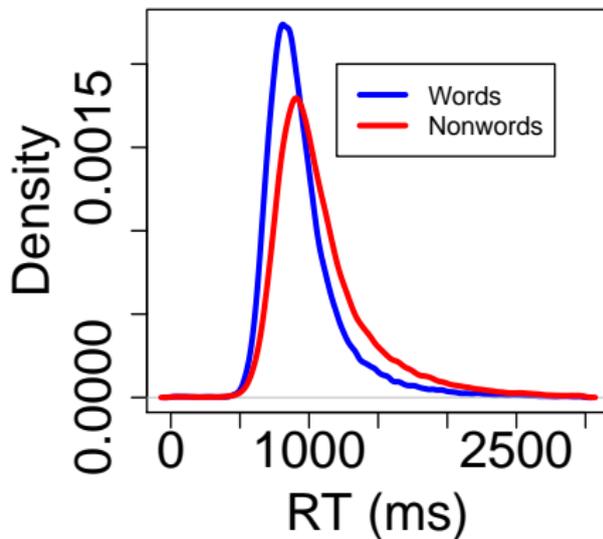


# Response Latency

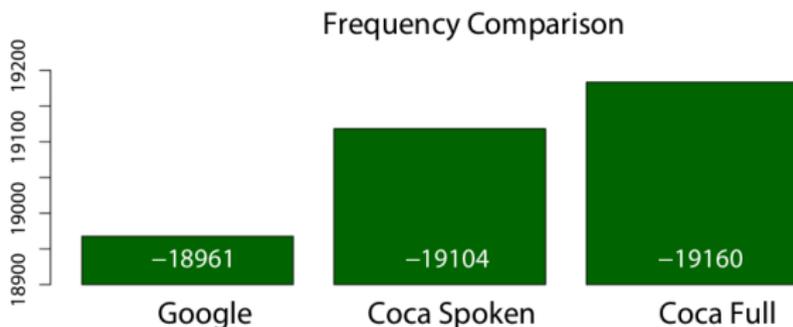
## All Latencies



## Word vs. Nonword



## Linear mixed effects regression (R core team, 2013; Bates et al., 2014)



- Counts based on all genres in COCA are better than spoken language
- Counts based on the Google Unigram corpus are less predictive than COCA

Thank you!

Stay tuned for the public release of the database:  
<https://aphl.artsrn.ualberta.ca/>