**Engineering**
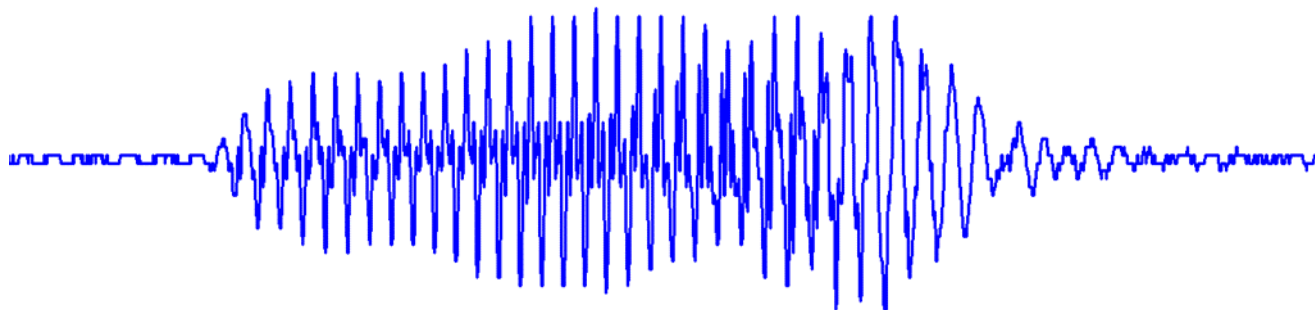
# Paralinguistic Speech Analysis

Ruben Bloom
ECE4095
30th June 2015

# Significant Contributions

- My supervisors, Professor Tom Drummond and Dr. Wai Ho Li provided invaluable guidance and planning for this project.

# Poster

# It's not just what you say, but how you say it.

# Intro to Computational Paralinguistics

- WHAT *are* underline paralinguistics?

- All aspects of speech beyond the words!

  - Pitch
  - Pauses
  - Speech Rate
  - Pitch Patterns
  - Umms and Ahhs

- Can be extracted from speech waveform automatically!

- Automatic Speech Recognition for non-verbals!

# The Goal: Speaking Style

- **IDENTIFY**

- **EXTRACT**

- **ANALYSE**

*for SCIENCE!!*
*for ENGINEERING!!!*

# System Overview

MONASH University

# Pitch

# Pauses and Utterances

- Why. They. Matter.
- whytheymatter . . .

# Pauses and Utterances

# VAD: Detecting Pauses and Utterances

- VAD

- (Can do with ASR, but didn't)

- Three attempts

# Finality Patterns

- Pitch movement at the end of utterance.

- English: Up, Down, Neutral.

- Up for statements: High Rising Terminal/Uptalk

# Uptalk: A Hot Topic

**Psychology Today** Find a Therapist ▾ Topics ▾

## The Uptalk Epidemic
Can you say something without turning it into a question?
Post published by Hank Davis on Oct 06, 2010 in Caveman Logic

**f SHARE** **y TWEET** **8+ SHARE** **✉ EMAIL**

**LEX·I·CON VALLEY** | A BLOG ABOUT LANGUAGE | DEC. 16 2014 3:48 PM

## Young Women Shouldn't Have to Talk Like Men to Be Taken Seriously

By Marybeth Seitz-Brown

13.5k   292   609

**Daily Mail AUSTRALIA**

Home | U.K. | U.S. | News | Sport | TV&Showbiz | Femail | Health | **Science** |
Latest Headlines | *Science* | Pictures

### Want a promotion? Don't speak like an AUSSIE: Rising in pitch at the end of sentences make you sound 'insecure'

**Well**

**MIND**

## Overturning the Myth of Valley Girl Speak
By JAN HOFFMAN    DECEMBER 23, 2013 4:12 PM    🗩 422 Comments

# Finality Patterns

# System Overview

# Speaking Style Comparisons

Three groups

1. Oral Presentations by 2<sup>nd</sup> Year Psychology Undergraduates (n = 7)

2. TED Talks with a minimum of 1M views each (n = 7)

3. Ceremonial wedding speeches (n = 6)

# Results: Group Differences

- Differences in expected direction from previous studies! (Rosenburg and Hirschburg, Stangert, Strangert and Gustafson)

- TED Speakers have higher mean pitch and pitch variation, longer and more pauses, and higher percentage of utterances with Falling Intonation.

- Australian Psychology Undergrads have high levels of High Rising Terminal.

|  | Student Presentations | TED Speakers | Wedding Speakers |
|---|---|---|---|
| Length (s) | 261 | 300 | 300 |
| f0 Mean (Hz) | 178 | 222 | 166 |
| f0 Std (Hz) | 43 | 53 | 30 |
| Mean Pause Length (s) | 0.689 | 0.745 | 0.980 |
| Pause Length Std (s) | 0.462 | 0.432 | 0.702 |
| Mean Utterance Length (s) | 2.837 | 2.491 | 2.348 |
| Utterance Length Std (s) | 5.325 | 3.779 | 1.951 |
| Speech/Pause Percentage (%/%) | 80/20 | 76/24 | 71/29 |
| HRT Percentage (%/100) | 0.16 | 0.06 | 0.00 |
| FI Percentage (%/100) | 0.10 | 0.26 | 0.15 |

*Table 2 Mean of each speaker group for each variable listed.*

# Results: Individual Differences

- Dramatic differences!

- Occur across group boundaries: it is possible to be charismatic in many ways!

- Need more sophisticated measures to differentiate style for charisma.

| | Range of Top 5 | Range of Bottom 5 |
|---|---|---|
| **Mean Pitch** | 220-260Hz | 155-161Hz[1] |
| **Pitch Standard Deviation[2]** | 73-46Hz | 30-42Hz[1] |
| **Mean Pause Duration** | 900-1600ms | 500-600ms |
| **Pause Duration Std[3,4]** | 900-1500ms | 100-160ms |
| **Pause Percentage** | 28-42% | 11-18% |
| **HRT Percentage** | 13-42% | 0-0% |
| **FI Percentage** | 20-62% | 2-5% |

*Table 3 Range of scores for highest scoring and lowest scoring five individuals on each measure.*

# Limitations

- Inaccuracies in segmentation have downstream effects.

- Small samples.

- Limited testing of feature extraction.

- No control of audio recording environment and resultant quality.

# Outcomes

Overall Goal: Use computational paralinguistics to develop tools useful for the scientific and engineering analysis of speaking style.

- Speech features identified!

- Automatic extraction achieved!

- Tools for quantitatively identifying paralinguistic differences between speakers achieved!

- Sample analysis achieved!

Project changed from original specifications, but majority of requirements still met!

# Future Directions

- Rigorous testing of extracted features.

- More speech features:
    - Speech Rate
    - Filled Pauses
    - Energy Variation (emphasis)
    - Whole pitch contours
- Improved Segmentation

- Better analysis of speech/pause rhythm, e.g. "frequency analysis"

- More speech types: interviews, political speeches, conversations

- Speaker Diarisation

- PCA on analysed data

- Combination with linguistic features from ASR for comprehensive speaker profile.

# Graphical Output Demo & Spreadsheet