



Language Technologies Institute



# Multimodal Affective Computing

# Lecture 5: Vocal Messages

# Louis-Philippe Morency Jeffrey Girard

Originally developed with help from Stefan Scherer and Tadas Baltrušaitis

# **Outline of this week's lectures**

- Multiple Layers of Vocal Messages
  - What we can convey with speech
- Fundamentals of speech production and hearing
  - Anatomy of the vocal tract and the physiology of hearing
  - Fundamental speech measures (direct vs. perceptual measures)
- Prosodic manipulation and its meaning
- Use and detection of varying voice quality
- Nonverbal vocal expressions
  - Laughter, pause filler (e.g. uh, um), and moans
- Practical tools for speech signal processing
- Automatic Techniques for visual processing



# **Upcoming Schedule**

- Week 5
  - Tuesday 2/12: Lecture on Vocal Messages
  - Thursday 2/14: Discussion (visual & vocal messages)
- Week 6:
  - Tuesday 2/19: Lecture on Verbal Messages
  - Thursday 2/21: Proposal presentations
  - Sunday 2/24: Due date for proposal reports
- Week 7:
  - Tuesday 2/26: Lecture on Statistical Analysis
  - Thursday 2/28: Discussion (verbal messages)



A story is told as much by silence as by speech. Susann Griffin

# **Multiple Layers of Vocal Messages**

#### Linguistic layer

- Carries the semantic information
- Language, grammar, phonemes, parts of prosody

#### Paralinguistic layer

- Non-linguistic and nonverbal layer
- Conveys information about current affective state, mood, attitude etc.
- Voice quality, prosody, nonverbal vocal expressions (e.g. laughter, moans, sighs)
- Main topic of this lecture!

#### Extralinguistic layer

 Identifies the speaker (e.g. age, gender, pitch range, habitual characteristics)



# Fundamentals of Speech



Language Technologies Institute



# Anatomy of the vocal tract

- Speech production involves multiple organs that shape the sound
  - Lungs
  - Vocal folds
  - Mouth/tongue
  - Nasal cavity
  - Lips





Pictures from Gray's anatomy 1918



# Anatomy of the vocal tract (2)

- Vocal folds vibrate and produce fundamental frequency (f0)
- Vibration is based on muscular tension and air pressure



Source: youtube.com



# Anatomy of the vocal tract (3)

- Vowels
- Tongue position influences the tone
- Tongue modulates the length of the cavities in the mouth





Source: wikipedia.org



# Anatomy of the vocal tract (4)

- Mouth opening, lip rounding and teeth influence the sound
- Other Sounds:
  - Nasal sounds (e.g. [m], [n])
  - Stops (e.g. [k],[t],...)
  - · · ·



Ladefoged, A course in phonetics, 2004



Language Technologies Institute

# Anatomy of the hearing organ

- Outer ear collects sound waves
  - Filters and allows directionality detection
- In the middle ear the sound waves are transferred via drum and three little bones called ossicles
- The inner ear transfers physical waves into nervous signals



Pictures from Gray's anatomy 1918



# Anatomy of the hearing organ (2)

- Transformation between sound waves and nervous signals sent to the brain is not linear
- There are perceptual differences:
  - Not every frequency is perceived with the same intensity or accuracy



Pictures from Gray's anatomy 1918



# Anatomy of the hearing organ (3)



![](_page_12_Picture_2.jpeg)

# Anatomy of the hearing organ (4)

#### Perceptual measures

- Loudness vs. intensity
- Pitch vs. fundamental frequency

![](_page_13_Figure_4.jpeg)

#### Try it out: http://tinyurl.com/px5g7cf

![](_page_13_Picture_6.jpeg)

# Representation Of Speech

![](_page_14_Picture_1.jpeg)

Language Technologies Institute

![](_page_14_Picture_4.jpeg)

# **Representations of speech**

Different domain representations:

- Time
- frequency
- Signal can be transferred from one domain to the other using Fourier transformation
  - Measures the amount of "presence" of a sine wave with a certain frequency in the original signal
- Two types: broadband and narrowband spectrogram
  - Bandwidth is defined by width of analysis window

![](_page_15_Figure_8.jpeg)

![](_page_15_Picture_9.jpeg)

# **Fundamental Frequency**

Fundamental frequency (f0) is the basic resonance of the vocal folds

 Harmonics are multiples of f0

![](_page_16_Figure_3.jpeg)

![](_page_16_Picture_4.jpeg)

![](_page_16_Picture_6.jpeg)

# **Formant Frequencies**

Formant frequencies are resonance frequencies dependent on the length of the vocal tract cavities

 The vocal tract is manipulated by the tongue etc.

![](_page_17_Picture_3.jpeg)

Source: wikipedia.org

![](_page_17_Picture_5.jpeg)

# **Example Use of Formant Frequencies**

Vowel space

 Gender and age define vowel space size.

What could be a reason for reduced vowel space?

![](_page_18_Figure_4.jpeg)

![](_page_18_Picture_5.jpeg)

# **Example Use of Formant Frequencies**

#### Vowel space

 Medical relevance: Parkinson's Disease ALS Depression

![](_page_19_Figure_3.jpeg)

![](_page_19_Picture_4.jpeg)

![](_page_19_Picture_5.jpeg)

# Mel frequency cepstral coefficients (MFCC)

# Mel frequency cepstral coefficients (MFCC)

- Compact representation of the spectrum
- Emulates human hearing
- Popular for speech recognition

![](_page_20_Figure_5.jpeg)

![](_page_20_Picture_6.jpeg)

# Mel frequency cepstral coefficients (MFCC)

![](_page_21_Figure_1.jpeg)

![](_page_21_Picture_2.jpeg)

# Mel frequency cepstral coefficients (MFCC)

![](_page_22_Figure_1.jpeg)

Language Technologies Institute

23

# Prosody

![](_page_23_Picture_1.jpeg)

Language Technologies Institute

![](_page_23_Picture_4.jpeg)

#### **Prosody can strongly influence the meaning**

![](_page_24_Picture_1.jpeg)

![](_page_24_Picture_2.jpeg)

# **Elements of prosody**

- Prosody, the suprasegmental envelope of an utterance, is composed by:
  - Syllable length
  - Loudness
  - Pitch
  - Pauses
- Prosody influences and defines:
  - Prosodic boundaries
  - Question or statement
  - Sarcasm
  - Emotional state
  - Meaning of words (e.g. in Chinese)

![](_page_25_Picture_12.jpeg)

![](_page_25_Picture_13.jpeg)

![](_page_25_Picture_14.jpeg)

### **Prosodic Boundaries**

Can you hear the difference?

![](_page_26_Picture_2.jpeg)

I met Mary and Elena's mother at the mall yesterday.

Sally saw % the man with the binoculars. Sally saw the man % with the binoculars.

When Madonna sings % the song is a hit. When Madonna sings the song % it's a hit.

![](_page_26_Picture_6.jpeg)

# Same 'tune', different alignment

#### What is a good source of vitamins?

![](_page_27_Figure_2.jpeg)

![](_page_27_Picture_3.jpeg)

# Same 'tune', different alignment

Are legumes a source of vitamins?

![](_page_28_Figure_2.jpeg)

![](_page_28_Picture_3.jpeg)

# Same 'tune', different alignment

What are legumes a good source of?

![](_page_29_Figure_2.jpeg)

![](_page_29_Picture_3.jpeg)

# **Other Uses of Pitch Contour Analysis**

- Rising pitch contour towards the end of an utterance indicates a yes-no question
- WH-questions have a falling pitch contour
- Signaling doubt or uncertainty can be expressed by rising contour in the end of an utterance

![](_page_30_Picture_4.jpeg)

![](_page_30_Picture_6.jpeg)

#### **Yes-No question tune**

![](_page_31_Figure_1.jpeg)

Rise from the main accent to the end of the sentence.

![](_page_31_Picture_3.jpeg)

#### **Yes-No question tune**

![](_page_32_Figure_1.jpeg)

Rise from the main accent to the end of the sentence.

![](_page_32_Picture_3.jpeg)

#### **Yes-No question tune**

![](_page_33_Figure_1.jpeg)

Rise from the main accent to the end of the sentence.

![](_page_33_Picture_3.jpeg)

# **Rising statements**

![](_page_34_Figure_1.jpeg)

![](_page_34_Picture_2.jpeg)

# **Emotional state**

#### Prosody can be used to express emotion

Acoustic Parameters	Arousal/Stress	Happiness	Anger	Sadness	Fear	Boredom
Speech rate and fluency						
Number of syllables	>	>=	$\diamond$	<	>	<
Syllable duration	<	<=	$\diamond$	>	<	>
Number/duration of pauses	<	<	<	>	<>	>
f <sub>0</sub> and prosody						
f <sub>0</sub> mean	>	>	>	<	>	<=
f <sub>0</sub> deviation	>	>	>	<	>	<
f <sub>0</sub> range	>	>	>	<	<>	<=
Gradient of f <sub>0</sub>	>	>	>	<	<>	<=
Vocal Effort/phonation						
Intensity (dB) mean	>	>=	>	<=		<=
Intensity (dB) deviation	>	>	>	<		<
Jitter		>=	>=		>	=
Shimmer		>=	>=		>	=

![](_page_35_Picture_3.jpeg)

# **Emotional state: Neutral**

![](_page_36_Figure_1.jpeg)

- Low pitch
- Low pitch variation
- Moderate loudness

![](_page_36_Picture_5.jpeg)

![](_page_36_Picture_6.jpeg)

# **Emotional state: Angry**

![](_page_37_Figure_1.jpeg)

- High pitch
- High pitch variation
- High intensity

![](_page_37_Picture_5.jpeg)

![](_page_37_Picture_6.jpeg)

# **Emotional state: Fear**

![](_page_38_Figure_1.jpeg)

- Average pitch
- Average pitch variation
- Average intensity

![](_page_38_Picture_5.jpeg)

![](_page_38_Picture_6.jpeg)

Language Technologies Institute

# **Emotional state (5)**

 QUIZ (Who knows the Germans best?)

![](_page_39_Figure_2.jpeg)

- Allowed answers:
  - Happiness
  - Anger
  - Sadness

- Disgust
- Fear
- Neutral

![](_page_39_Picture_10.jpeg)

# Voice Quality

![](_page_40_Picture_1.jpeg)

Language Technologies Institute

![](_page_40_Picture_4.jpeg)

# **Voice Quality**

- Does not refer to a term concerning fidelity or "goodness"
- Refers to the timbre or coloring of a voice
- Functions:
  - Meaning or disambiguation of words (e.g. in Gujarati) (<u>http://www.phonetics.ucla.edu/vowels/chapter12/gujarati.html</u>)
  - Paralinguistic signal
    - Attitude
    - Mood
    - Social factors (e.g. standing, (*inter*-)personality)
    - Affective state
    - Turn-management (e.g. in Finnish)

![](_page_41_Picture_11.jpeg)

# **Voice Quality**

- Can be seen as the signal residue after removing effects of the vocal tract filter
- The phonation and manner of the vocal folds vibrate play a major role (exceptions: e.g. whisper)

![](_page_42_Figure_3.jpeg)

![](_page_42_Figure_4.jpeg)

![](_page_42_Picture_5.jpeg)

# **Vocal Folds: Phonation Gestures**

![](_page_43_Figure_1.jpeg)

#### Adductive tension: Movement toward the mi dline of the vocal fold.

- Medial compression: adductive force on vocal processes
- Longitudinal pressure: tension of vocal folds

![](_page_43_Picture_5.jpeg)

![](_page_44_Picture_1.jpeg)

- "Neutral" mode
- Muscular adjustments moderate
- Vibration of vocal folds periodic, full closing of glottis, no audible friction
- Frequency of vibration and loudness in low to mid range for conversational speech

![](_page_44_Picture_6.jpeg)

![](_page_44_Picture_8.jpeg)

# **Harsh/Tense Voice**

![](_page_45_Picture_1.jpeg)

 Very strong tension of vocal folds, very high tension in vocal tract

![](_page_45_Figure_3.jpeg)

![](_page_45_Picture_4.jpeg)

![](_page_45_Picture_5.jpeg)

# **Whispery Voice**

![](_page_46_Figure_1.jpeg)

![](_page_46_Picture_2.jpeg)

- Little or no vocal fold vibration
- Very low adductive tension
- Medial compression moderately high
- Longitudinal tension moderately high
  - Turbulence generated by friction of air in and above larynx, with vocal folds not vibrating

![](_page_46_Picture_8.jpeg)

# **Creaky Voice (vocal Fry)**

![](_page_47_Picture_1.jpeg)

- Vocal fold vibration at low frequency, irregular
- Low tension
- The vocal folds strongly adducted
- Longitudinal tension weak
- Moderately high medial compression

![](_page_47_Figure_7.jpeg)

![](_page_47_Picture_8.jpeg)

C<mark>arnegie Mellon University</mark>

#### Language Technologies Institute

lilule

### **Breathy Voice**

- Tension low
  - Minimal adductive tension
  - Weak medial compression
- Medium longitudinal vocal fold tension
- Vocal folds do not come together completely, leading to frication

![](_page_48_Figure_8.jpeg)

![](_page_48_Picture_9.jpeg)

![](_page_48_Picture_10.jpeg)

![](_page_49_Picture_0.jpeg)

# **Creaky voice (Vocal fry)**

![](_page_49_Figure_2.jpeg)

![](_page_49_Picture_3.jpeg)

### **Examples of Voice Quality Measures**

- Open Quotient (OQ)
- Normalized Amplitude Quotient (NAQ)
- Peak Slope

![](_page_50_Picture_4.jpeg)

![](_page_50_Picture_5.jpeg)

# Suicide Prevention

[ICASSP 2013]

- Nonverbal indicators of suicidal ideations
- Dataset: 30 suicidal adolescents/30 non-suicidal adolescents
- Suicidal teenagers use more breathy tones

![](_page_51_Figure_5.jpeg)

![](_page_51_Picture_6.jpeg)

# Nonverbal Vocal Expressions

![](_page_52_Picture_1.jpeg)

Language Technologies Institute

![](_page_52_Picture_4.jpeg)

# Nonverbal vocal expressions

- Nonverbal vocal expressions are paralinguistic utterances
  - Backchannel (e.g. uhu, hm, yeh)
  - Laughter
  - Moans/Sighs
  - Pause fillers (e.g. um, uh)
- Varying functions and possible interpretations
- Distinct prosodic and vocal characteristics
- Often accompanied with distinct facial expressions (e.g. laughing)

![](_page_53_Picture_9.jpeg)

# **Backchannels**

- Backchannels are important fragments of speech
- Functions:
  - Turn-taking management
  - Signals agreement and attention

•••

![](_page_54_Figure_6.jpeg)

![](_page_54_Picture_7.jpeg)

- Important social communicational expression
- Understood by all cultures
- Varying meanings:
  - Humorous laughter
  - Signals agreement
  - Uncertainty (e.g. social laughter, nervous laughter)
- Multimodal!

![](_page_55_Picture_8.jpeg)

*"Laughter* [...] *is performed almost exclusively during social encounters; solitary laughter seldom occurs except in response to media, a source of vicarious social stimulation."* – Provine and Yong

![](_page_55_Picture_10.jpeg)

![](_page_56_Picture_1.jpeg)

![](_page_56_Picture_2.jpeg)

0

![](_page_57_Figure_1.jpeg)

![](_page_57_Figure_2.jpeg)

Time [s]

![](_page_57_Picture_3.jpeg)

Language Technologies Institute

![](_page_57_Picture_6.jpeg)

- Acoustically very distinct but variable
  - Snort-like
  - Inhaled
  - Exhaled
- Often arranged in *bouts* consisting of single *calls*
  - Often about 200ms in length, but quite variable in length

![](_page_58_Figure_7.jpeg)

![](_page_58_Picture_8.jpeg)

# **Hesitations/Pause fillers**

- Hesitations (e.g. pause fillers: um, uh)
- Multiple functions
  - Signal anxiety or nervousness
  - Proficiency in speaking
- Characterized by prolonged vowels and static spectrum

![](_page_59_Picture_6.jpeg)

![](_page_59_Picture_7.jpeg)

# Practical Tools for Speech Processing

![](_page_60_Picture_1.jpeg)

Language Technologies Institute

### **Useful Software**

COVAREP - A Cooperative Voice Analysis Repository for Speech Technologies (Matlab and Octave) <u>https://github.com/covarep/covarep/</u>

![](_page_61_Picture_2.jpeg)

![](_page_61_Picture_3.jpeg)

Degottex, G., Kane, J., Drugman, T., Raitio, T., and Scherer, S., COVAREP - A collaborative voice analysis repository for speech technologies, under review at International Conference on Acoustics, Speech, and Signal Processing (ICASSP).

![](_page_61_Picture_5.jpeg)

# openSMILE - Speech & Music Interpretation by Large Space Extraction

 openSMILE is a fast, real-time (audio) feature extraction utility (C++)

http://opensmile.sourceforge.net/

![](_page_62_Picture_4.jpeg)

![](_page_62_Picture_5.jpeg)

![](_page_62_Picture_6.jpeg)

# **Useful Software**

- Praat
- http://www.fon.hum.uva.nl/praat/
- Useful for feature extraction/visualization and more!

![](_page_63_Picture_4.jpeg)

![](_page_63_Picture_5.jpeg)

![](_page_63_Picture_6.jpeg)

# **Upcoming Conferences**

- Affective Computing & Intelligent Interaction (ACII)
  - http://acii-conf.org/2019/
  - Deadline: April 12, 2019
- International Conference on Multimodal Interaction (ICMI)
  - https://icmi.acm.org/2019/
  - Deadline: May 7<sup>th</sup>, 2019
- Affective Computing Pre-Conference at ISRE (International Society of Research on Emotion)
  - <u>https://www.isre2019.org/program/pre-conferences/affective-computing</u>
  - Deadline: May 7<sup>th</sup>, 2019
- Empirical Methods in Natural Language Processing
  - https://www.emnlp-ijcnlp2019.org/
  - Deadline: May 21th, 2019

![](_page_64_Picture_13.jpeg)