### speech

CSC401/2511 – Natural Language Computing – Spring 2019 Lecture 8 Frank Rudzicz University of Toronto

## **This lecture**

- Acoustics.
- Speech production.
- Speech perception.

 Some images from Gray's Anatomy, Jim Glass' course 6.345 (MIT), the Jurafsky & Martin textbook, Encyclopedia Britannica, the Rolling Stones, the Pink Floyds.





### acoustics

### What is sound?

- Sound is a time-variant pressure wave created by a vibration.
  - Air particles hit each other, setting others in motion.

    - High pressure  $\equiv$  **compressions** in the air (C).
    - Low pressure  $\equiv$  rarefactions within the air (R).





## What is sound?





Frequency F = 1/T



**phase**  $\phi$  is displacement of a signal in time. E.g., with  $\phi = \pi/2$ ,

 $\sin(x + \phi) = \cos(x)$ 



### What is sound?

• A single tone is a sinusoidal function of pressure and time.

- Amplitude: n. The degree of the displacement in the air. This is similar to 'loudness'. Often measured in Decibels (dB).
- Frequency: *n*. The number of cycles within a unit of time. e.g., **1 Hertz (Hz) = 1 oscillation/second**



## **Speech waveforms**



## **Superposition of sinusoids**

- Superposition: *n*. the adding of sinusoids together.
- Phase: *n*. The horizontal offset of a sinusoid (φ).





## **Extracting sinusoids from waveforms**

- As we will soon see, the relative amplitudes and frequencies of the sinusoids that combine in speech are often extremely indicative of the speech units being uttered.
  - If we could separate the waveform into its component sinusoids, it would help us classify the speech being uttered.
  - But the shape of the signal changes over time

(it's not a single repeating pattern)...





## **Short-time windowing**





• Speech waveforms change drastically over time.

- We *move* a short analysis window (assumed to be time-invariant) across the waveform in time.
  - E.g. frame shift: 5-10 ms
  - E.g. frame length: 10-25 ms
- 5-10 ms10-25 ms



## Window types



CSC401/2511 - Spring 2019

**TORONTO** 

#### **Extracting a spectrum**



Any Colour You Like (track 8)



### **Extracting a spectrum in a window**





## Aside – Euler's formula

 Extracting sinusoids is possible because of a relationship between *e* and sinusoids expressed in Euler's formula:

$$e^{ix} = \cos(x) + i\sin(x)$$







#### **The continuous Fourier transform**



Input:

Continuous signal x(t).

**Output**: Spectrum X(F)

$$X(F) = \int_{-\infty}^{\infty} x(t) e^{-i2\pi Ft} dt$$

(No need to memorize these )



• It's invertible, i.e.,  $x(t) = \int_{-\infty}^{\infty} X(F)e^{i2\pi Ft} dF$ . • It's linear, i.e., for  $a, b \in \mathbb{C}$ , if h(t) = ax(t) + by(t), then H(F) = aX(F) + bY(F)

Fun fact: Fourier instructed Champollion.

It needs **continuous** input x(t)... **uh oh?** 



### **Discrete signal representation**

- Sampling: vbg. measuring the amplitude of a signal at regular intervals.
  - e.g., 44.1 kHz (*CD*), 8 kHz (*telephone*).
  - These amplitudes are initially measured as continuous values at discrete time steps.



#### **Discrete signal representation**

• Nyquist rate:

*n.* the **minimum** sampling rate necessary to preserve a signal's **maximum** frequency.

- i.e., twice the maximum frequency, since we need ≥ 2 samples/cycle.
- Human speech is very informative ≤ 4 kHz,
  ∴ 8 kHz sampling.





## **Discrete Fourier transform (DFT)**



**Input**: Windowed signal  $x[0] \dots x[N-1].$ 

**Output**: *N* complex numbers X[k] ( $k \in \mathbb{Z}$ )

(No need to memorize these )

$$X[k] = \sum_{n=0}^{N-1} x[n] e^{-i2\pi k \frac{n}{N}}$$

 $M_{-1}$ 

• Algorithm(s): the Fast Fourier Transform (FFT) with complexity  $O(N \log N)$ .

• (Aside) The **Cooley-Tukey algorithm** divides-and-conquers by breaking the DFT into smaller ones  $N = N_1 N_2$ .



## **Discrete Fourier transform (DFT)**

 Below is a 25 ms Hamming-windowed signal from /iy/ as in 'bull sh<u>ee</u>p', and its spectrum as computed by the DFT.



But this is all just for a small window ...



#### **Spectrograms**

• **Spectrogram**: *n.* a 3D plot of **amplitude** and **frequency** 

**Over time** (higher 'redness'  $\rightarrow$  higher amplitude).



## **Effect of window length**





#### **Spectrograms**





## **Aside – Filtering**

- Sometimes you only want part of a signal.
  - E.g., you have measurements of lip aperture over time you know that they can't move > 5-10 Hz.
  - E.g., you know there's some low-frequency Gaussian noise in either the environment or transmission medium.



 Low- and high-pass filters can be combined in series, yielding a band-pass filter.





# speech production

#### The vocal tract



- Many physical structures are co-ordinated in the production of speech.
- Generally, sound is generated by passing air through the vocal tract.
- Sound is modified by constricting airflow in particular ways.



## The neurological origins of speech

 Studying how systems break down can indicate how they work.



- **Reduced** hierarchical syntax.
- Anomia.
- Reduced "mirroring" between observation and execution.

- Normal intonation/rhythm.
- Meaningless words.
- 'Jumbled' syntax.
- Reduced comprehension.



## The neurological origins of speech

• Cranial nerves carry messages from the brain to the various **articulators**.



- Cranial nerves carry messages from the brain to the various articulators.
  - Damage to these nerves can result in neuro-motor disorders such as cerebral palsy.
  - These may be another example of the noisy channel.



## **Fundamental frequency**

 F<sub>0</sub>: n. (fundamental frequency), the rate of vibration of the glottis – often very indicative of the speaker.





## Prosody

- Sonorant: n. Any sustained sound in which the glottis is vibrating (i.e., the sound is 'voiced').
  - Includes some consonants (e.g., /w/, /m/).
- Prosody: n. the modification of speech acoustics in order to convey some extra-lexical meaning:
  - **Pitch**: Changing of  $F_0$  over time.
  - **Duration**: The length in time of sonorants.
  - Loudness: The amount of energy produced by the lungs.



## **Pitch prosody example**





## **Pitch can modify meaning**

e.g., I ask you
 "<u>who</u> is that?"





Pitch tends to rise when uttering novel or important information.



## **Pitch can modify meaning**

- <u>I</u> never said she stole my money. (Someone else said it)
- I <u>never</u> said she stole my money. (It never happened)
- I never <u>said</u> she stole my money. (I just hinted at it)
- I never said <u>she</u> stole my money. (Someone else stole it)
- I never said she <u>stole</u> my money. (She just borrowed it)
- I never said she stole <u>my</u> money. (She stole someone else's)
- I never said she stole my *money*. (She stole my heart).



## **Phonemes**

- Phoneme:
- Phonemes can be partitioned into **manners of articulation**:
  - Vowels:
  - Fricatives:
  - Stops/plosives:
  - Nasals:
  - Semivowels:
  - Affricates:

*n.* a distinctive unit of speech sound. open vocal tract, no nasal air. **noisy**, with air passing through a tight constriction (e.g., 'shift'). complete vocal tract constriction and burst of energy (e.g., 'papa'). air passes through the **nasal** cavity (e.g., '<u>m</u>a<u>m</u>a'). similar to vowels, but typically with more constriction (e.g., 'wall'). Alveolar stop followed by fricative.



## **Place of articulation**

- The **location** of the *primary constriction* can be:
  - **Alveolar**: constriction near the alveolar ridge (e.g., /t/)
  - **Bilabial**: touching of the lips together (e.g., /m/, /p/)
  - **Dental**: constriction of/at the teeth (e.g., /th/)
  - Labiodental: constriction between lip and teeth (e.g., /f/)
  - Velar: constriction at or near the velum (e.g., /k/).





## **Phonemic alphabets**

- There are several alphabets that categorize the sounds of speech.
  - The International Phonetic Alphabet (IPA) is popular, but it uses non-ASCII symbols.
  - The TIMIT phonemic alphabet will be used by default in this course.
  - Other popular alphabets include ARPAbet, Worldbet, and OGIbet, usually adding special cases.
    - E.g., /pcl/ is the period of silence immediately before a /p/.

TIMIT	IPA	e.g.
/iy/	/i <sup>y</sup> /	b <mark>ea</mark> t
/ih/	/1/	b <u>i</u> t
/eh/	/ɛ/	b <u>e</u> t
/ae/	/æ/	b <mark>a</mark> t
/aa/	/a/	B <mark>o</mark> b
/ah/	/_/	b <u>u</u> t
/ao/	/၁/	b <u>ou</u> ght
/uh/	/ʊ/	b <u>oo</u> k
/uw/	/u/	b <u>oo</u> t
/ux/	/ʉ/	s <u>ui</u> t
/ax/	/ə/	<u>a</u> bout



## TIMIT Phonemic alphabet (incomplete)

Vow	el e.g.	st	ор	e.g.		fricative	e.g.
/iy,	' b <u>ea</u> t	/	b/	<u>B</u> il <u>b</u> o		/s/	<u>S</u> ea
/ih,	′ b <u>i</u> t	/	d/	<u>d</u> a <u>d</u> a		/f/	<u>F</u> rank
/eh	/ b <u>e</u> t	/	g/	<u>G</u> aga		/z/	<mark>∠</mark> appa
/ae	/ b <u>a</u> t	/	p/	<u>P</u> ippin		/th/	<u>th</u> is
/aa	/ B <u>o</u> b	/	′t/	<u>T</u> oo <u>t</u> s		/sh/	<u>Sh</u> ip
/ah	/ b <u>u</u> t	/	k/	<u>k</u> i <u>ck</u>		/zh/	a <mark>z</mark> ure
/ao	/ b <u>ou</u> ght	-				/v/	Vogon
/uh	/ b <u>oo</u> k	na	asal	e.g.		/dh/	then
/uw	/ b <u>oo</u> t	/	m/	<u>M</u> a <u>m</u> a	Ľ		
/ux	/ s <u>ui</u> t	/	n/	<u>n</u> oo <u>n</u>		(Incom	olete)
/ax	l <u>a</u> bout	/r	ng/	thi <mark>ng</mark>			



#### **Phoneme sequences**

- Often, we assume that a spoken utterance can be partitioned into a sequence of non-overlapping phonemes.
  - Demarking the periods during which certain phonemes are being uttered is called transcription or annotation (\*).
  - This approach has problems (e.g., when *exactly* does one phoneme end and another begin?), but it's useful for **classification**.



What are some characteristics of the six manners of articulation?



## Vowels (1/6)

- There are approximately 19 vowels in Canadian English, including diphthongs in which the articulators move over time.
- Vowels are distinguished primarily by their formants. (?)

other	e.g.
/er/	B <u>er</u> t
/axr/	b <u>u</u> tter

diphthong	e.g.
/ey/	b <u>ai</u> t
/ow/	b <mark>oa</mark> t
/ay/	b <u>i</u> te
/oy/	b <u>oy</u>
/aw/	b <u>ou</u> t
/ux/	s <u>ui</u> t

Mono- phthong	e.g.
/iy/	b <mark>ea</mark> t
/ih/	b <u>i</u> t
/eh/	b <u>e</u> t
/ae/	b <mark>a</mark> t
/aa/	B <u>o</u> b
/ao/	b <u>ou</u> ght
/ah/	b <u>u</u> t
/uh/	b <u>oo</u> k
/uw/	b <u>oo</u> t
/ax/	<u>a</u> bout
/ix/	ros <u>e</u> s
(	<u>.</u>

## The uniform tube



 The positions of the tongue, jaw, and lips change the shape and cross-sectional area of the vocal tract.



## **Uniform tubes in practice**

- Many musical instruments are based on the idea of uniform (or, in many cases, bent) tubes.
- Longer tubes produce 'deeper' sounds (lower frequencies).
  - A tube ½ the length of another will be 1 octave higher.





#### **Vowels as concatenated tubes**

• The vocal tract can be modelled as the concatenation of dozens, hundreds, or thousands of tubes.



#### Aside – waves in concatenated tubes

• We model the **volume velocity**  $U_k$  and the **pressure variation**  $p_k$  at position x in the  $k^{th}$  lossless tube (whose area is  $A_k$ ) at time t



#### Waves in concatenated tubes

 Because of partial wave reflections that occur at tube boundaries, we can generate spectra with particular resonances.



#### **Formants and vowels**

• Formant: *n*. A concentration of energy within a frequency band. Ordered from low to high bands (e.g.,  $F_1$ ,  $F_2$ ,  $F_3$ ).



### The vowel trapezoid



CSC401/2511 - Spring 2019

TORONTO

#### **Tongues and formants**



#### Fricatives (2/6)

 Fricatives are caused by acoustic turbulence at a narrow constriction whose position determines the sound.





#### **Fricatives**

- Fricatives have four places of articulation.
- Each place of articulation has a voiced fricative (i.e., the glottis can be vibrating), and an unvoiced fricative.

	Unvo	biced	Voiced		
Labial	/f/	<mark>f</mark> ee	/v/	<u> </u>	
Dental	/th/	<u>th</u> ief	/dh/	<u>Th</u> ee	
Alveolar	/s/	<u>s</u> ee	/z/	<u>Z</u> ardo <u>z</u>	
Palatal	/sh/	<u>sh</u> e	/zh/	<u>Zh</u> a- <u>zh</u> a	



#### **Unvoiced fricatives**





## **Plosives (3/6)**

- Plosives build pressure behind a complete closure in the vocal tract.
- A sudden release of this constriction results in brief noise.



#### **Plosives**

• **Plosives** have three places of articulation:

	Unvo	biced	Voiced		
Labial	/p/	<mark>p</mark> or <mark>p</mark> oise	/b/	<u></u> bab₀oon	
Alveolar	/t/	<u>t</u> or <u>t</u>	/d/	<u>d</u> o <u>d</u> o	
Velar	/k/	<u>k</u> i <u>ck</u>	/g/	<u><b>G</b></u> oo <mark>g</mark> le	

- Voiced stops are usually characterized by a "voice bar" during closure, indicating the vibrating glottis.
- Formant transitions are very informative in classification.



## **Voicing in plosives**



#### **Formant transitions in plosives**



• Despite a **common** vowel, the **motion** of  $F_2$  and  $F_3$  into (and out of) the vowel helps identify the plosive.

## Nasals (4/6)

- Nasals involve lowering the velum so that air passes through the nasal cavity.
- Closures in the oral cavity (at same positions as plosives) change the resonant characteristics of the nasal sonorant.



### **Formant transitions among nasals**



• Despite a common vowel, the motion of  $F_2$  and  $F_3$  before and after each nasal helps to identify it.

#### Semivowels (5/6)

- Semivowels act as consonants in syllables and involve constriction in the vocal tract, but there is less turbulence.
  - They also involve slower articulatory motion.
- Laterals involve airflow around the sides of the tongue.



## Semivowels

• Semivowels are often sub-classified as glides or liquids.

	Semiv	Nearest vowel	
Glides	/w/	<u>W</u> o <u>w</u>	/uw/
	/y/	<mark>у</mark> о <u>у</u> о	/iy/
Liquids	/r/	<u>r</u> ea <u>r</u>	/er/
	/1/	<u></u> Lu <u>l</u> u	/ow/

- Semivowels are more constricted versions of corresponding vowels.
  - Similar formants, though generally weaker.



## **Semivowels**



 Note the drastic formant transitions which are more typical of semivowels.



## Affricates and aspirants (6/6)

- There are two affricates: /jh/ (voiced; e.g., judge) and /ch/ (unvoiced; e.g., <u>ch</u>ur<u>ch</u>).
  - These involve an alveolar stop followed by a fricative.
  - Voicing in /jh/ is normally indicated by voice bars, as with plosives.
- There's only one aspirant in Canadian English: /h/ (e.g., <u>h</u>at)
  - This involves turbulence generated at the glottis,
  - In Canadian English, there is **no** constriction in the vocal tract.



## **Affricates and aspirants**





## **Alternative pronunciations**

- **Pronunciations** of words can vary significantly, but with observable **frequencies**.
  - The Switchboard corpus is a phonetically annotated database of speech recorded in telephone conversations.

because				ab	out		
ARPAbet	%	ARPAbet	%	ARPAbet	%	ARPAbet	%
b iy k ah z	27%	k s	2%	ax b aw	32%	b ae	3%
b ix k ah z	14%	k ix z	2%	ax b aw t	16%	b aw t	3%
k ah z	7%	k ih z	2%	b aw	9%	ax b aw dx	3%
k ax z	5%	b iy k ah zh	2%	ix b aw	8%	ax b ae	3%
b ix k ax z	4%	b iy k ah s	2%	ix b aw t	5%	b aa	3%
b ih k ah z	3%	b iy k ah	2%	ix b ae	4%	b ae dx	3%
b ax k ah z	3%	b iy k aa z	2%	ax b ae dx	3%	ix b aw dx	2%
k uh z	2%	ax z	2%	b aw dx	3%	ix b aa t	2%



## **Known effects of pronunciation**

- Speakers tend to drop or change pronunciations in predictable ways in order to reduce the effort required to co-ordinate the various articulators.
  - Palatalization generally refers to a conflation of phonemes closer to the frontal palate than they 'should' be.
  - Final t/d deletion is simply the omission of alveolar plosives from the ends of words.

Palatalization			Final t/d Deletion			
Phrase	Lexical	Reduced	Phrase	Lexical	Reduced	
set your	s eh t y ow r	s eh ch er	find him	f ay n d h ih m	f ay n ix m	
not yet	n aa t y eh t	n aa ch eh t	and we	ae n d w iy	eh n w iy	
did you	d ih d y uw	d ih jh y ah	draft the	d r ae f t dh iy	d r ae f dh iy	



## Variation at syllable boundaries





## **Phonological variation**

The acoustics of a phoneme depend strongly on the context in which that phoneme occurs.



That must make **recognizing** phonemes hard, right? How do humans do it?



## The inner ear



- Time-variant waves enter the ear, vibrating the tympanic membrane.
- This membrane causes tiny bones (incl. **malleus**) to vibrate.
- These bones in turn vibrate a structure within a shellshaped bony structure called the cochlea.



## The cochlea and basilar membrane





- The basilar membrane is covered with tiny hair-like nerves – some near the base, some near the apex.
- High frequencies are picked up near the base, low frequencies near the apex.
- These nerves fire when activated, and communicate to the brain.



#### **The Mel-scale**

- Human hearing is not equally sensitive to all frequencies.
  - We are **less** sensitive to frequencies > 1 kHz.
- A **mel** is a unit of pitch. Pairs of sounds which are **perceptually** equidistant in pitch are separated by an equal number of **mels**.

$$Mel(f) = 2595 \log_{10} \left( 1 + \frac{f}{700} \right)$$



(No need to memorize this either)



## **Aside – Challenges of perception**

 Cochlear implants replace the basilar membrane and stimulate the auditory nerve directly.









#### Next...

- How the Mel scale is used in ASR.
- Automatic speech recognition.

