Speech

CSC401/2511 – Natural Language Computing – Fall 2024 Lecture 8 University of Toronto

This lecture

- Speech signals
- Articulatory phonetics

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



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 ϕ 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'.
- Frequency: *n*. The number of cycles within a unit of time. e.g., **1 Hertz (Hz) = 1 oscillation/second**



The inner ear



- Time-variant waves enter the ear, vibrating the tympanic membrane.
 - This membrane causes tiny bones (the **ossicles**) 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.

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





Speech waveforms



Superposition of sinusoids

• Superposition: *n*. the adding of sinusoids together.





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





CSC401/2511 - Fall 2024

• 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: 10-30 ms
 - E.g. frame length: 25-40 ms
- 10-30 ms 25-40 ms



Window types



Extracting a spectrum



Any Colour You Like (track 8)



Extracting a spectrum in a window



Frequency (Hz)



Euler's formula

 Extracting sinusoids uses a relationship between natural exponent *e* and sinusoids expressed in Euler's formula:

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







The continuous Fourier transform



Input:

Continuous signal x(t).

Output: Spectrum X(F)

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



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)...



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

- *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 \leq 4 kHz,
 - ∴ At least 8 kHz sampling (16kHz the norm)





Discrete Fourier transform (DFT)

• Input: Windowed signal $x[0] \dots x[N-1].$ • Output: N complex numbers $X[k] \ (k \in \mathbb{Z})$ $X[k] = \sum_{n=0}^{N-1} x[n]e^{-i2\pi k\frac{n}{N}}$

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



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

SPECTROGRAM, R = 128 SPECTROGRAM, R = 512 3500 3500 3000 3000 2500 2500 kouenbeut Acuenbeur 1500 1500 1000 1000 500 500 C 0.45 0.05 0.4 0.05 35 0 0.5 n 0.4 0.45 **Narrow-band** Wide-band (better time (better frequency resolution) resolution)



Spectrograms





Articulatory phonetics



Sounds and transcriptions

- We are often interested in the meaning of an utterance
- In English, we often transcribe utterances as word tokens
 - We write: <How to recognize speech>
- Is this "what was said?"
 - We might write instead: <How to wreck a nice beach>
 - We can transcribe (or even adopt) foreign words
 - 沙发 = <sofa>, not <sandy hair>
 - We can even transcribe brand new words
 - <skibidi toilet>
- We can instead transcribe "speech sounds"



Phones and phonetics

- Phonetics is the study of speech sounds
- A phone is a unit of speech
 - Denoted with square braces: [t], [t^h], [u]
 - Language-independent
- Phones which are perceived "similarly" are grouped into phonemes
 - Denoted with slashes: /t/, /u/
 - $[t], [t^h] \mapsto /t/$
 - Language-dependent
- Transcriptions are often in-between:
 - $['t^hu] \mapsto [t^hu] \mapsto [tu] \mapsto /tu/$
- We will be very loose with the distinction



Phonetic transcription

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





Phonetic 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** phonetic alphabet will be used by **default** in this course.
 - Other popular alphabets include ARPAbet, Worldbet, and OGlbet, usually adding special cases.
 - E.g., [pcl] is the period of silence immediately before a [p].

TIMIT	IPA	e.g.
[iy]	[i ^y]	b <mark>ea</mark> t
[ih]	[I]	b <mark>i</mark> t
[eh]	[8]	b <u>e</u> t
[ae]	[æ]	b <u>a</u> t
[aa]	[a]	B <mark>o</mark> b
[ah]	[Λ]	b <u>u</u> t
[ao]	[כ]	b <u>ou</u> ght
[uh]	[ʊ]	b <u>oo</u> k
[uw]	[u]	b <u>oo</u> t
[ux]	[u]	s <u>ui</u> t
[ax]	[ə]	<u>a</u> bout



TIMITbet (incomplete)

Vowel	e.g.	stop	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 <mark>i</mark> t	[d]	<u>d</u> a <u>d</u> a	[f]	<u>F</u> rank
[eh]	b <mark>e</mark> t	[g]	<u>G</u> a <u>g</u> a	[z]	<mark>_</mark> appa
[ae]	B <u>a</u> t	[p]	<u>P</u> i <u>pp</u> in	[th]	<u>th</u> is
[aa]	B <mark>o</mark> 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]	∨ ogon
[uh]	b <u>oo</u> k	nasal	e.g.	[dh]	then
[uw]	b <u>oo</u> t	[m]	<u>M</u> a <u>m</u> a		
[ux]	s <u>ui</u> t	[n]	<u>n</u> oo <u>n</u>	(Incomplete)	
[ax]	<u>a</u> bout	[ng]	thi <u>ng</u>		



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.
- We can classify phones by how they are **produced**

A taxonomy of phones

- Phones fall into two broad categories
- Vowels are
 - Always periodic
 - Produced with relatively unobstructed airflow
 - Use tongue, lips, and jaw to produce resonances in vocal tract, in turn generating formants
- Consonants are
 - Mostly noisy (not nasals, semivowels)
 - Produced by obstructing airflow
 - Classified by the place and manner of primary obstruction, as well as voicing

Voicing and fundamental frequency

- Voiced phones are produced with vibrating vocal folds
 - The space between the folds is the **glottis**
- All vowels are voiced; consonants can be **unvoiced**
- **F**₀: *n*. (fundamental frequency), the rate of vibration (Hz)
 - Very indicative of speaker

	Avg F_0 (Hz)	Min F_0 (Hz)	Max F_0 (Hz)
Male	125	80	200
Female	225	150	350
Children	300	200	500

Vowels

- 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

diphthonge.g.[ey]bait[ow]boat[ow]boat[ay]bite[oy]boy[aw]bout[ux]suit

Mono- phthong	e.g.
[iy]	b <u>ea</u> t
[ih]	b <u>i</u> t
[eh]	b <u>e</u> t
[ae]	b <u>a</u> t
[aa]	B <mark>o</mark> 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

Uniform tubes

- Formants and resonances can be approximated with tubes
- 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.

The uniform tube

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

Vowels as concatenated tubes

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

Waves in concatenated tubes

Reflections at tube boundaries produce resonances which amplify certain frequencies

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).

Tongues, lips, and formants

The vowel trapezoid

TORONTO

Manner of articulation

- Consonants are classified by place and manner of obstruction
- For manner:
 - Fricatives:
 - Stops/plosives:
 - Nasals:
 - Semivowels:
 - Affricates:
 - Taps:

noisy, with air passing through a tight constriction (e.g., '<u>sh</u>ift').

- **complete** vocal tract constriction and burst of energy (e.g., '*papa*'). air passes through the **nasal** cavity (e.g., '*mama*').
- similar to vowels, but typically with more constriction (e.g., '<u>w</u>a<u>ll</u>'). Alveolar stop followed by fricative.
- Quick collision of articulators ('butter')

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]).
 - **Glottal:** constriction of the glottis ([q])

Fricatives

• 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 folds can be vibrating), and an unvoiced fricative.

Labia-dontal				Voiced		
Labio-dental	[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:

	Unve	oiced	Voiced		
Labial	[p]	<mark>p</mark> or <mark>p</mark> oise	[b]	<u>b</u> a <u>b</u> oon	
Alveolar	[t]	<u>t</u> or <u>t</u>	[<i>d</i>]	<u>d</u> o <u>d</u> o	
Velar	[k]	<u>k</u> i <u>ck</u>	[g]	<u>G</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 F₂ and F₃ into (and out of) the vowel helps identify the plosive.

Nasals

- 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 F₂ and F₃ 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	
Clidad	[w]	<u>W</u> o <u>w</u>	[uw]
Glides	[y]	<mark>у</mark> оуо	[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

- There are two affricates: [jh] (voiced; e.g., judge) and [ch] (unvoiced; e.g., <u>ch</u>ur<u>ch</u>).
 - These involve an **palatal 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

UNIVERSITY OF TORONTO

Other topics in phonetics

- The grouping of phones into syllables
 - Consisting of a vowel (nucleus), and optionally preceding (onset) and succeeding (coda) consonants
 - Only certain sequences are permissible in English
 - Syllables may be made more prominent via pitch, duration, or loudness
- The prosody, or intonation and rhythm, of an utterance
 - Prominence can also indicate phrase boundaries
 - Gradual F0 movement (tune) can indicate a question or statement
- These are especially important to text-to-speech synthesis

