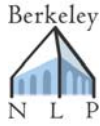


Natural Language Processing



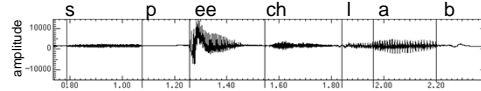
The Speech Signal

Dan Klein – UC Berkeley

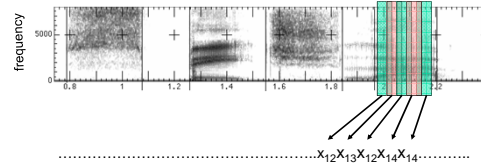


Speech in a Slide

- Frequency gives pitch; amplitude gives volume



- Frequencies at each time slice processed into observation vectors



Articulation



Articulatory System



- Nasal cavity
- Oral cavity
- Pharynx
- Vocal folds (in the larynx)
- Trachea
- Lungs

Sagittal section of the vocal tract (Techmer 1880)
Text from Ohala, Sept 2001, from Sharon Rose slide

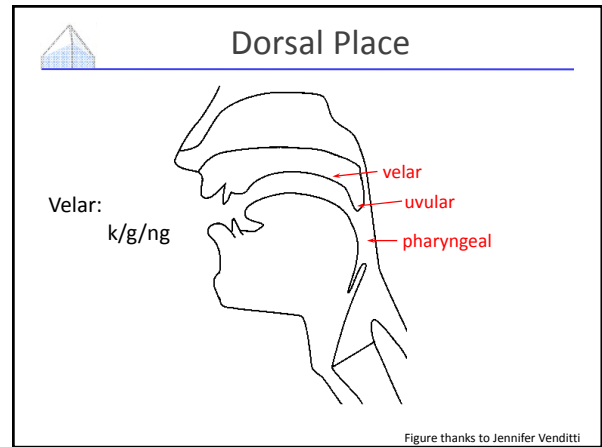
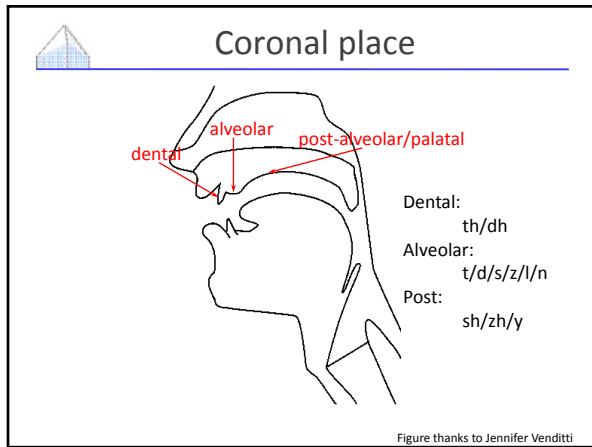
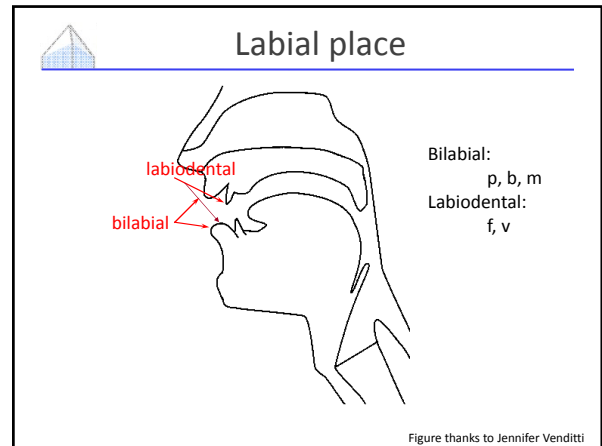
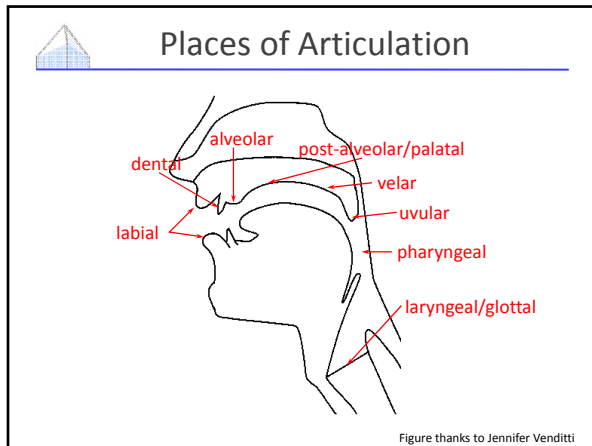


Space of Phonemes

	LABIAL		CORONAL				DORSAL			RADICAL		LABINGEAL
	Bilabial	Labio-dental	Dental	Alveolar	Palato-alveolar	Retroflex	Palatal	Velar	Uvular	Pharyngeal	Epi-glottal	Glottal
Nasal	m	ɱ		n	ɲ	ɳ	ɲ	ŋ	ɴ			
Plosive	p b	ɸ β	t d				k g	q ɢ				
Fricative	f v	θ ð	s z	ʃ ʒ	ʂ ʐ	ç ʝ	x ɣ	χ ʁ	ħ	ħ	h ɦ	
Approximant		ʋ	ɹ			ɻ	ɰ					
Trill	ʙ		ʀ						ʀ			
Tap, Flap		ɸ	ɾ									
Lateral fricative			ɬ ɮ			ɮ						
Lateral approximant			l			ʎ	ʟ					
Lateral flap			ɺ									

- Standard international phonetic alphabet (IPA) chart of consonants

Place



Space of Phonemes

	LABIAL		CORONAL				DORSAL			RADICAL		GLOTTAL
	Bilabial	Labio-dental	Dental	Alveolar	Palato-alveolar	Retroflex	Palatal	Velar	Uvular	Pharyngeal	Epi-glottal	Glottal
Nasal	m	ɱ		n	ɲ		ɲ	ŋ	ɴ			
Plosive	p b	ɸ β	t d	ʈ ɖ	ʈ ɖ	ʈ ɖ	c ɟ	k ɡ	q ɢ		ʔ	ʔ
Fricative		f v	θ ð	s z	ʃ ʒ	ʂ ʐ	ç ʝ	x ɣ	χ ʁ	ħ	ħ	h ɦ
Approximant		ʋ		ɹ			j	ɰ				
Trill	ʙ			ʀ					ʀ			
Tap, Flap		ɸ		ɾ								
Lateral fricative				ɬ ɮ			ɬ ɮ					
Lateral approximant				l			ʎ	ʟ				
Lateral flap				ɺ								

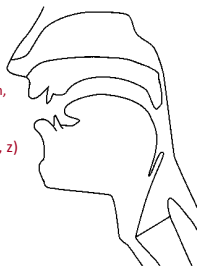
▪ Standard international phonetic alphabet (IPA) chart of consonants

Manner



Manner of Articulation

- In addition to varying by place, sounds vary by manner
- Stop: complete closure of articulators, no air escapes via mouth
 - Oral stop: palate is raised (p, t, k, b, d, g)
 - Nasal stop: oral closure, but palate is lowered (m, n, ŋ)
- Fricatives: substantial closure, turbulent: (f, v, s, z)
- Approximants: slight closure, sonorant: (l, r, w)
- Vowels: no closure, sonorant: (i, e, a)



Space of Phonemes

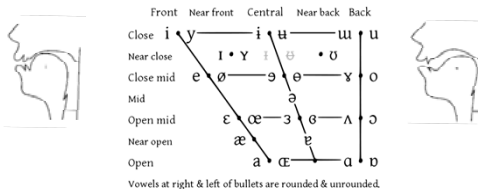
	LABIAL		CORONAL				DORSAL			RACIAL		LARYNGAL
	Bilabial	Labio-dental	Dental	Alveolar	Palato-alveolar	Retroflex	Palatal	Velar	Uvular	Pharyngeal	Epi-glottal	Glottal
Nasal	m	ɱ		n	ɲ	ɳ	ɰ	ŋ	ɴ			
Plosive	p b	ɸ β		t d	ʈ ɖ	ʈ ɖ	c ɟ	k ɡ	q ɢ		ʔ	ʔ
Fricative		f v	θ ð	s z	ʃ ʒ	ʂ ʐ	ç ʝ	x ɣ	χ ʁ	ħ	ħ	h ɦ
Approximant			ɹ	ɻ		ɻ	j	ɰ				
Trill	ʙ		ʀ						ʀ			ʀ
Tap, flap		ɸ		ɾ	ɽ							
Lateral fricative			ɬ ɮ		ɮ	ɮ	ɮ	ɮ				
Lateral approximant			l	ɭ	ɭ	ɭ	ɭ	ɭ				
Lateral flap			ɺ	ɻ	ɻ							

- Standard international phonetic alphabet (IPA) chart of consonants

Vowels



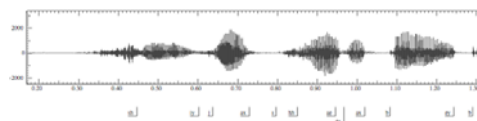
Vowel Space



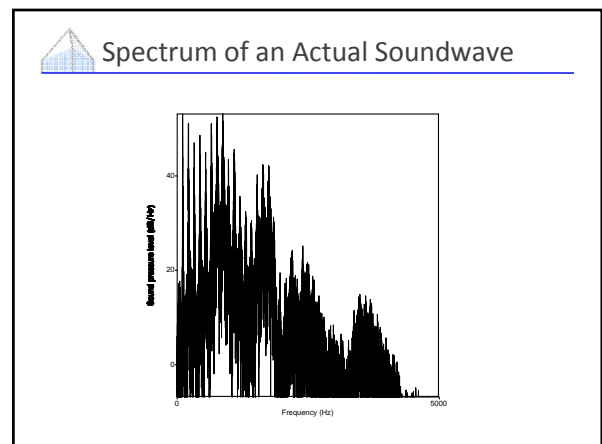
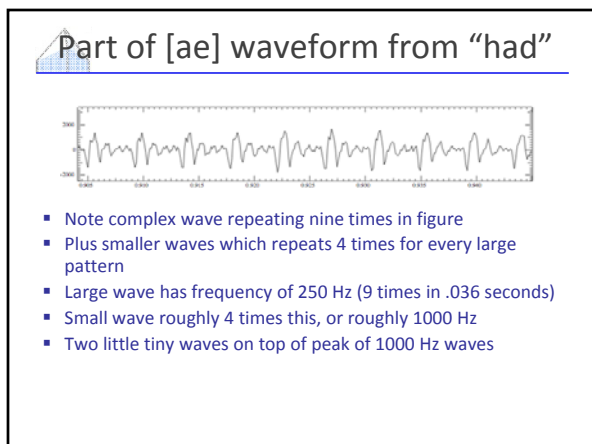
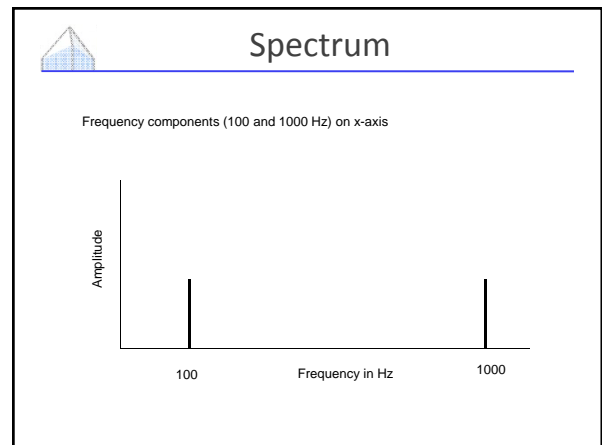
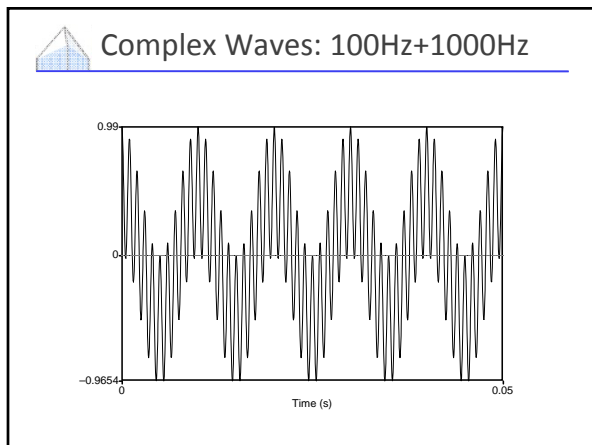
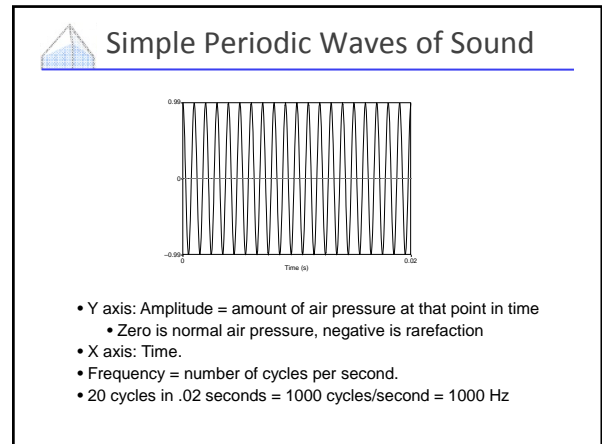
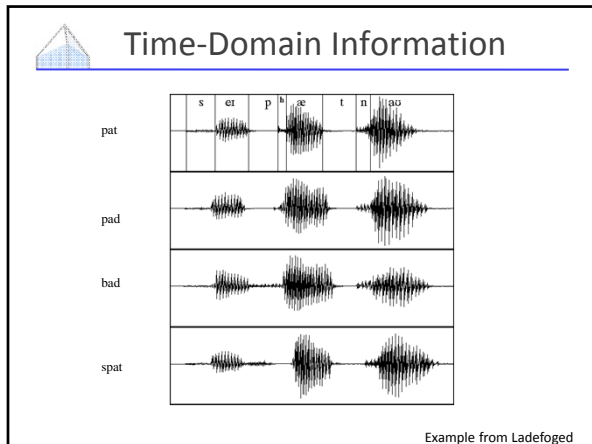
Acoustics



"She just had a baby"



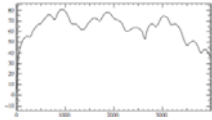
- What can we learn from a waveform?
 - No gaps between words (!)
 - Vowels are voiced, long, loud
 - Length in time = length in space in waveform picture
 - Voicing: regular peaks in amplitude
 - When stops closed: no peaks, silence
 - Peaks = voicing: .46 to .58 (vowel [ɪ]), from second .65 to .74 (vowel [æ]) and so on
 - Silence of stop closure (1.06 to 1.08 for first [b], or 1.26 to 1.28 for second [b])
 - Fricatives like [ʃh]: intense irregular pattern; see .33 to .46





Back to Spectra

- Spectrum represents these freq components
- Computed by Fourier transform, algorithm which separates out each frequency component of wave.

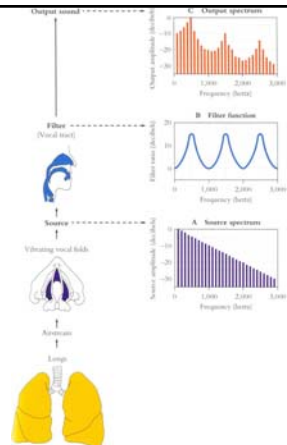


- x-axis shows frequency, y-axis shows magnitude (in decibels, a log measure of amplitude)
- Peaks at 930 Hz, 1860 Hz, and 3020 Hz.

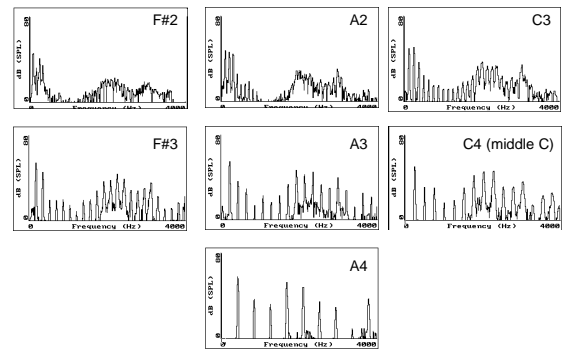
Source / Channel

Why these Peaks?

- **Articulation process:**
 - The vocal cord vibrations create harmonics
 - The mouth is an amplifier
 - Depending on shape of mouth, some harmonics are amplified more than others



Vowel [i] at increasing pitches

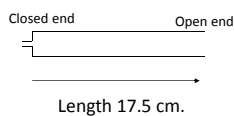


Figures from Ratre Wayland



Resonances of the Vocal Tract

- The human vocal tract as an open tube:



- Air in a tube of a given length will tend to vibrate at resonance frequency of tube.
- Constraint: Pressure differential should be maximal at (closed) glottal end and minimal at (open) lip end.

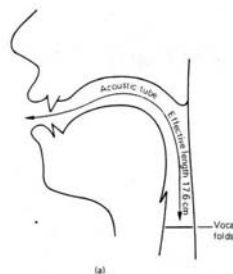
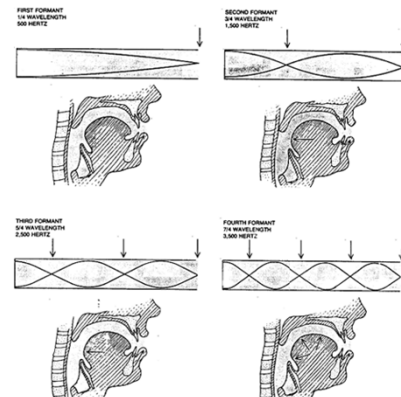


Figure from W. Barry

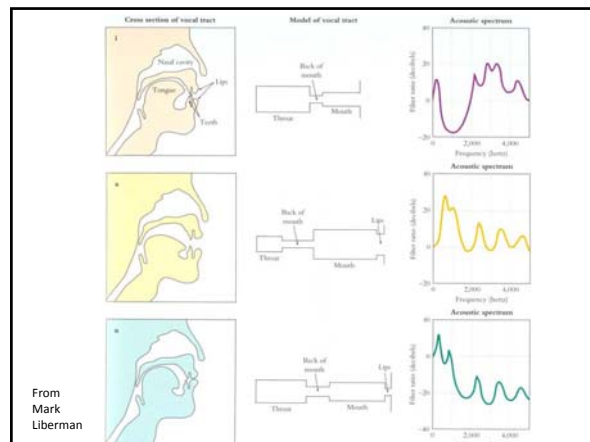


From Sundberg

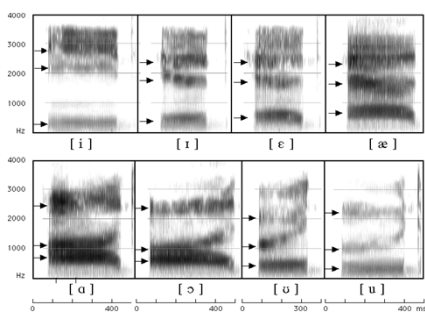


Computing the 3 Formants of Schwa

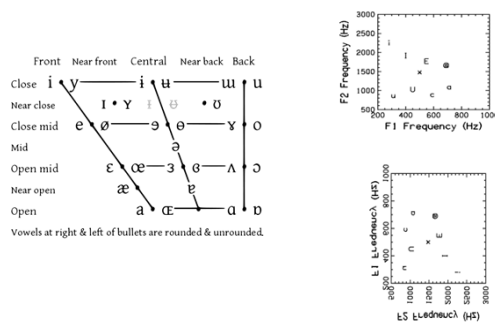
- Let the length of the tube be L
 - $F_1 = c/\lambda_1 = c/(4L) = 35,000/4 * 17.5 = 500\text{Hz}$
 - $F_2 = c/\lambda_2 = c/(4/3L) = 3c/4L = 3 * 35,000/4 * 17.5 = 1500\text{Hz}$
 - $F_3 = c/\lambda_3 = c/(4/5L) = 5c/4L = 5 * 35,000/4 * 17.5 = 2500\text{Hz}$
- So we expect a neutral vowel to have 3 resonances at 500, 1500, and 2500 Hz
- These vowel resonances are called **formants**



Seeing Formants: the Spectrogram



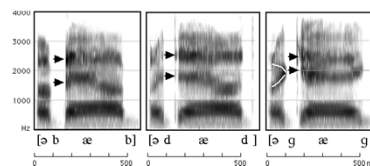
Vowel Space



Spectrograms



How to Read Spectrograms



- [bab]: closure of lips lowers all formants: so rapid increase in all formants at beginning of "bab"
- [dad]: first formant increases, but F2 and F3 slight fall
- [gag]: F2 and F3 come together: this is a characteristic of velars. Formant transitions take longer in velars than in alveolars or labials

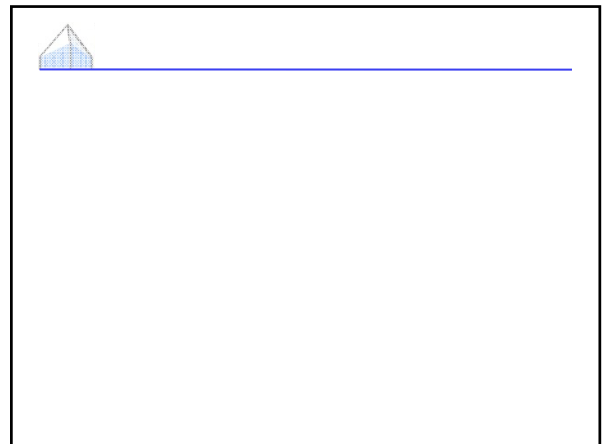
From Ladefoged "A Course in Phonetics"

"She came back and started again"

1. lots of high-freq energy

3.
4.
5.
6.
7.
8.
9.

From Ladefoged "A Course in Phonetics"



Deriving Schwa

- Reminder of basic facts about sound waves
 - $f = c/\lambda$
 - c = speed of sound (approx 35,000 cm/sec)
 - A sound with $\lambda=10$ meters: $f = 35$ Hz (35,000/1000)
 - A sound with $\lambda=2$ centimeters: $f = 17,500$ Hz (35,000/2)

American English Vowel Space

Figures from Jennifer Venditti, H. T. Bunnell

Dialect Issues

- Speech varies from dialect to dialect (examples are American vs. British English)
 - Syntactic ("I could" vs. "I could do")
 - Lexical ("elevator" vs. "lift")
 - Phonological
 - Phonetic
- Mismatch between training and testing dialects can cause a large increase in error rate

	American	British
all		
old		