

# Statistical NLP Spring 2010

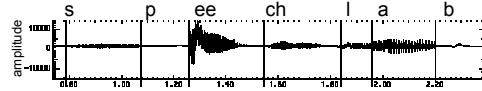


## Lecture 8: Speech Signal

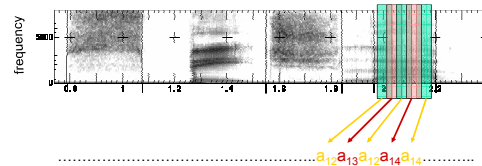
Dan Klein – UC Berkeley

### Speech in a Slide

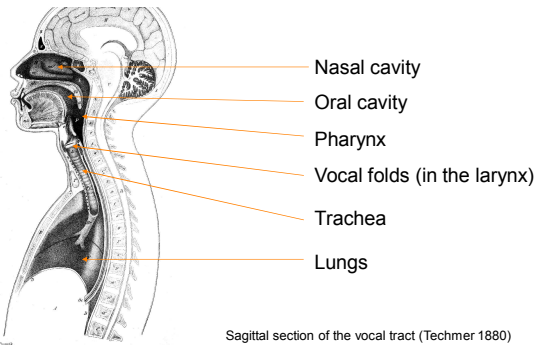
- Frequency gives pitch; amplitude gives volume



- Frequencies at each time slice processed into observation vectors



### Articulatory System



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

### Places of Articulation

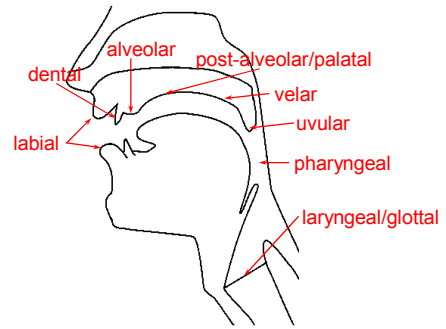
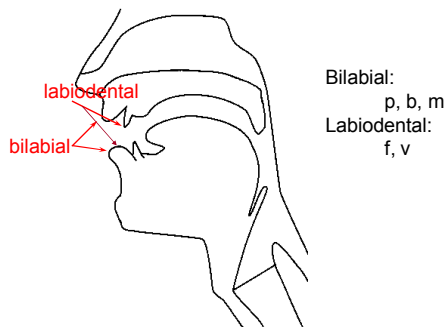


Figure thanks to Jennifer Venditti

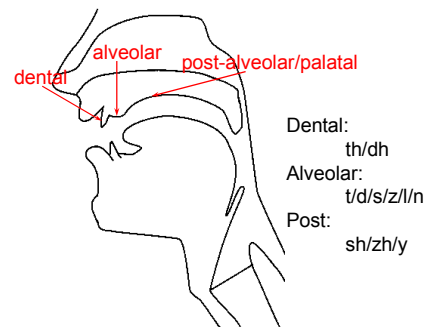
### Labial place



Bilabial:  
p, b, m  
Labiodental:  
f, v

Figure thanks to Jennifer Venditti

### Coronal place



Dental:  
th/dh  
Alveolar:  
t/d/s/z/l/n  
Post:  
sh/zh/y

Figure thanks to Jennifer Venditti

## Dorsal Place

Velar:  
k/g/ng

velar  
uvular  
pharyngeal

Figure thanks to Jennifer Venditti

## Space of Phonemes

	LABIAL		CORONAL				DORSAL			RADICAL		LARYNGAL
	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	ʈ ɖ	ʈ ɖ	c ɟ	k ɡ	q ɢ		ʔ	ʔ
Fricative	ɸ β	f v	θ ð	s z	ʃ ʒ	ʂ ʐ	ç ʝ	x ɣ	χ ʁ	ħ	ʕ	h ɦ
Approximant		ʋ		ɹ		ɻ	j	ɰ				
Trill	ʙ			ʀ					ʀ			
Tap, Flap		ɸ		ɾ		ɽ						
Lateral fricative				ɬ ɮ		ɮ	ɬ	ɮ				
Lateral approximant				l		ɭ	ʎ	ʟ				
Lateral flap				ɭ		ɮ						

- Standard international phonetic alphabet (IPA) chart of consonants

## 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			RADICAL		LARYNGAL
	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	ʈ ɖ	ʈ ɖ	c ɟ	k ɡ	q ɢ		ʔ	ʔ
Fricative	ɸ β	f v	θ ð	s z	ʃ ʒ	ʂ ʐ	ç ʝ	x ɣ	χ ʁ	ħ	ʕ	h ɦ
Approximant		ʋ		ɹ		ɻ	j	ɰ				
Trill	ʙ			ʀ					ʀ			
Tap, Flap		ɸ		ɾ		ɽ						
Lateral fricative				ɬ ɮ		ɮ	ɬ	ɮ				
Lateral approximant				l		ɭ	ʎ	ʟ				
Lateral flap				ɭ		ɮ						

- Standard international phonetic alphabet (IPA) chart of consonants

## Vowel Space

Front: Near front Central Nearback Back

Close i y ɨ ɥ ɯ ʉ

Near close ɪ ʏ ɘ ɚ ɝ

Close mid e ø ɘ ɚ ɝ

Mid ɛ œ ə ɜ ɞ

Open mid ɛ œ ə ɜ ɞ

Near open æ ɶ ɛ ɜ ɞ

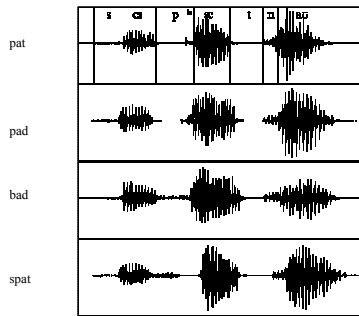
Open a ɶ ɛ ɜ ɞ

Vowels at right & left of bullets are rounded & unrounded.

## “She just had a baby”

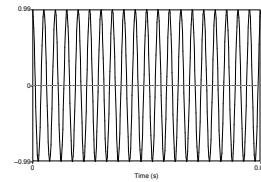
- What can we learn from a wavefile?
  - No gaps between words (!)
  - Vowels are voiced, long, loud
  - Length in time = length in space in waveform picture
  - Voicing: regular peaks in amplitude
  - When stops closed: no peaks, silence
  - Peaks = voicing: .46 to .58 (vowel [ɪ]), from second .65 to .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

## Non-Local Cues



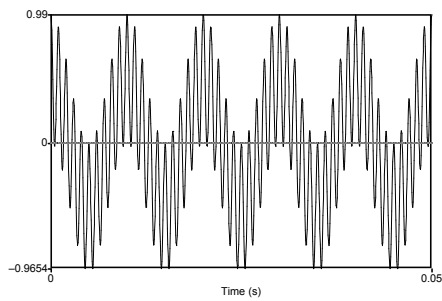
Example from Ladefoged

## Simple Periodic Waves of Sound



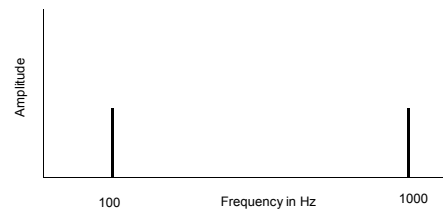
- Y axis: Amplitude = amount of air pressure at that point in time
  - Zero is normal air pressure, negative is rarefaction
- X axis: Time.
- Frequency = number of cycles per second.
- 20 cycles in .02 seconds = 1000 cycles/second = 1000 Hz

## Complex Waves: 100Hz+1000Hz

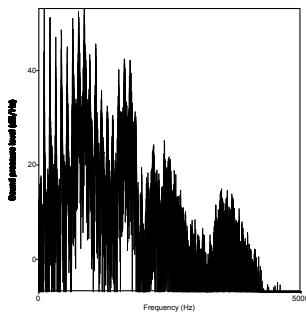


## Spectrum

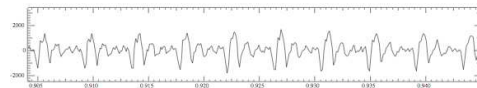
Frequency components (100 and 1000 Hz) on x-axis



## Spectrum of an Actual Soundwave



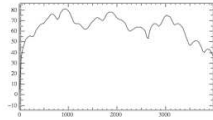
## Part of [æ] waveform from "had"



- Note complex wave repeating nine times in figure
- Plus smaller waves which repeats 4 times for every large pattern
- Large wave has frequency of 250 Hz (9 times in .036 seconds)
- Small wave roughly 4 times this, or roughly 1000 Hz
- Two little tiny waves on top of peak of 1000 Hz waves

## Back to Spectra

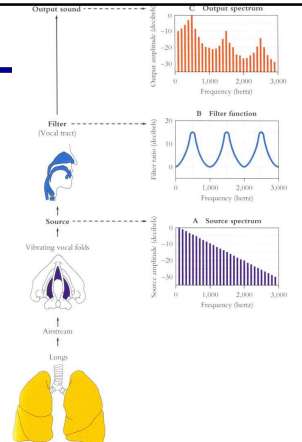
- Spectrum represents these freq components
- Computed by Fourier transform, algorithm which separates out each frequency component of wave.



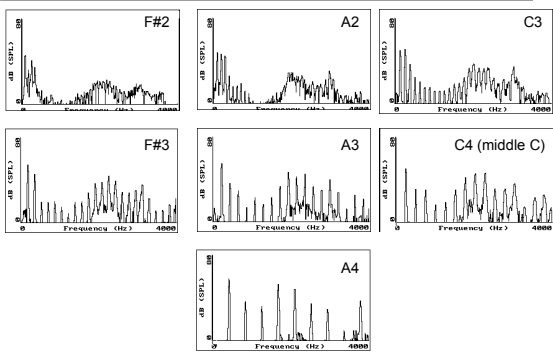
- x-axis shows frequency, y-axis shows magnitude (in decibels, a log measure of amplitude)
- Peaks at 930 Hz, 1860 Hz, and 3020 Hz.

## Why these Peaks?

- **Articulator 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] sung at successively higher 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.

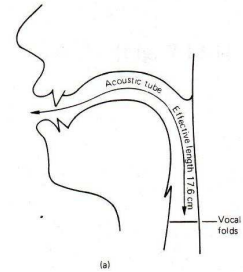
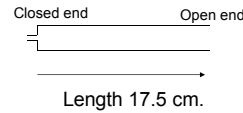
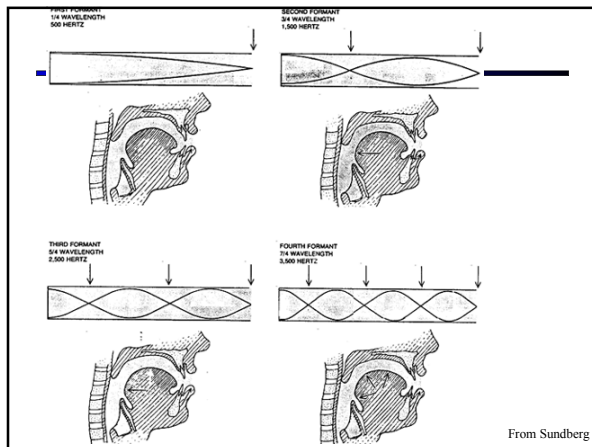


Figure from W. Barry



From Sundberg

## Computing the 3 Formants of Schwa

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

