

L3: Organization of speech sounds

Phonemes, phones, and allophones

Taxonomies of phoneme classes

Articulatory phonetics

Acoustic phonetics

Speech perception

Prosody

Phonemes and phones

Phoneme

- The smallest meaningful contrastive unit in the phonology of a language
- Each language uses small set of phonemes, much smaller than the number of sounds than can be produced by a human
- The number of phonemes varies per language, with most languages having 20-40 phonemes
 - General American has ~40 phonemes (24 consonants, 16 vowels)
 - The Rotokas language (Paupa New Guinea) has ~11 phonemes
 - The Taa language (Botswana) has ~112 phonemes

Phonetic notation

- International Phonetic Alphabet (IPA): consists of about 75 consonants, 25 vowels, and 50 diacritics (to modify base phones)
- TIMIT corpus: uses 61 phones, represented as ASCII characters for machine readability
 - TIMIT only covers English, whereas IPA covers most languages

Phones in the TIMIT Database

TIMIT	IPA	Example	TIMIT	IPA	Example
pcl	p̚	(p closure)	bcl	b̚	(b closure)
tcl	t̚	(t closure)	dcl	d̚	(d closure)
kcl	k̚	(k closure)	gcl	g̚	(g closure)
p	p	pea	b	b	bee
t	t	tea	d	d	day
k	k	key	g	g	gay
q	ʔ	bat	dx	r	dirty
ch	tʃ	choke	jh	dʒ	joke
f	f	fish	v	v	vote
th	θ	thin	dh	ð	then
s	s	sound	z	z	zoo
sh	ʃ	shout	zh	ʒ	azure
m	m	moon	n	n	noon
em	m̩	bottom	en	n̩	button
ng	ŋ	sing	eng	ŋ	Washington
nx	ɹ̥	winner	el	l̥	bottle
l	l	like	r	r	right
w	w	wire	y	j	yes
hh	h	hay	hv	fi	ahead
er	ɝ	bird	axr	ɝ	butter
iy	i	beet	ih	I	bit
ey	e	bait	eh	ɛ	bet
ae	æ	bat	aa	ɑ	father
ao	ɔ	bought	ah	ʌ	but
ow	o	boat	uh	ʊ	book
uw	u	boot	ux	ü	toot
aw	ɑ ^w	about	ay	ɑ ^y	bite
oy	ɔ ^y	boy	ax-h	ɔ	suspect
ax	ə	about	ix	ɪ	debit
epi		(epenthetic sil.)	pau		(pause)
h#		(silence)			

[Gold & Morgan, 2000]

Phone

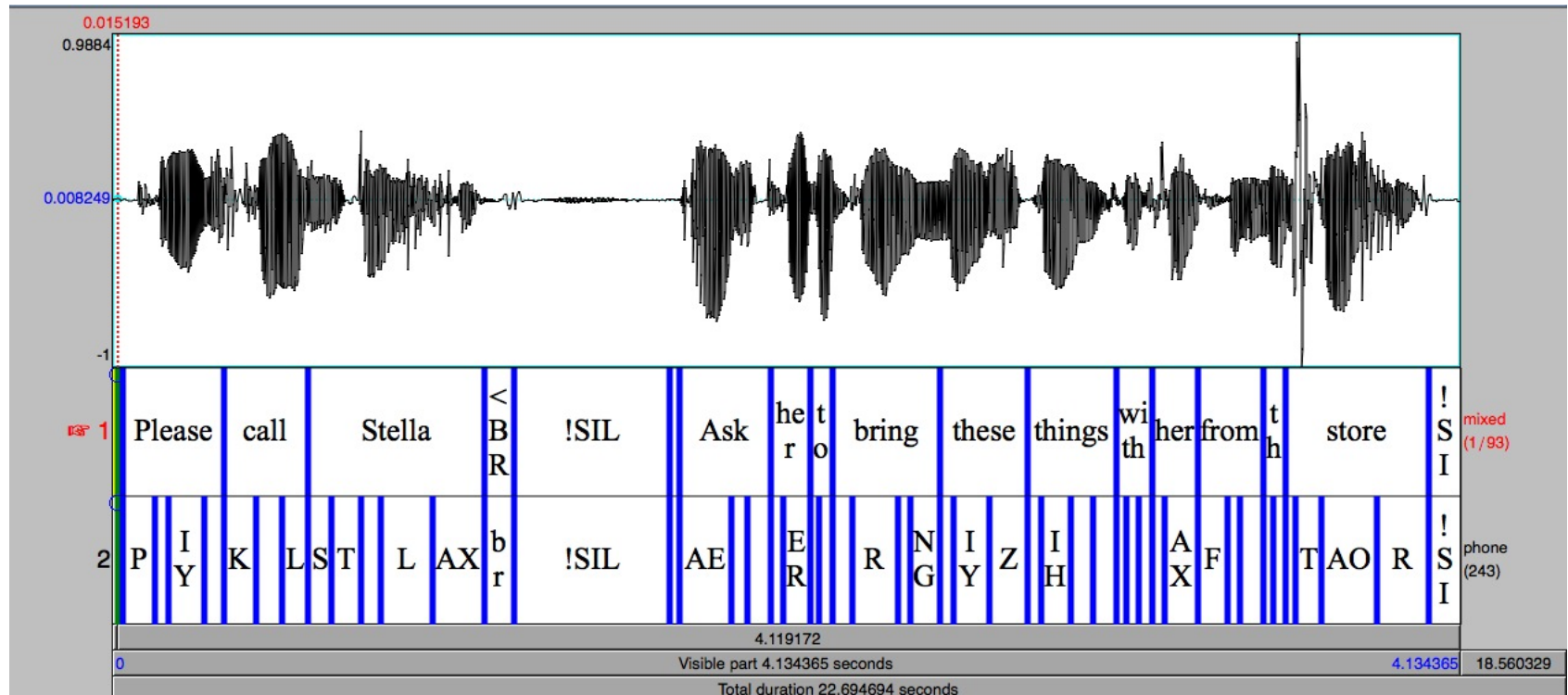
- The physical sound produced when a phoneme is articulated
- Since the vocal tract can have an infinite number of configurations, there is an infinite number of phones that correspond to a phoneme

Allophone

- A class of phones corresponding to a specific variant of a phoneme
 - Example: aspirated [p^h] and unaspirated [p] in the words *pit* and *spit*
 - Example: /t/ sounds in the words *tub*, *stub*, *but*, *butter*

Coarticulation

- The phenomenon whereby the articulatory configuration of a phoneme is affected by that of its neighboring phonemes
 - As a result, crisp boundaries between phonemes are hard to define
- Coarticulation is involved in the transformation of phonemes into allophones



http://groups.linguistics.northwestern.edu/documentation/images/praat_aligned.jpg

Branches of phonetics

Phonology vs. phonetics

- Phonology is concerned with the distribution and patterns of speech sounds in a particular language, or in languages in general
- Phonetics is concerned with the study of speech sounds and their production, classification, and transcription
 - In a nutshell, phonetics deals with the physical nature of speech sounds, and not with their relations to other speech sounds in particular languages

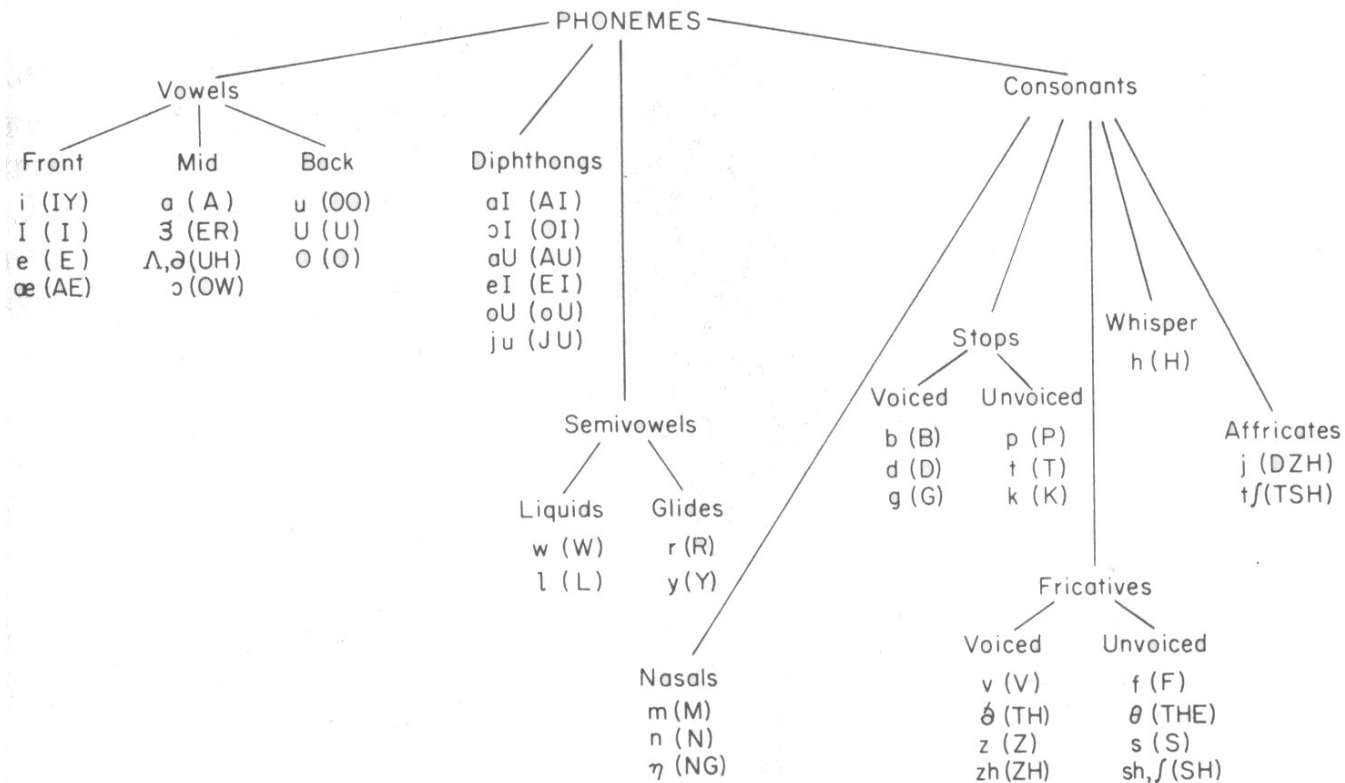
Three basic approaches to the study of phonetics

- Articulatory phonetics is concerned with the position, shape and movements of speech articulators
- Acoustic phonetics is concerned with the spectro-temporal properties of the speech sound waves
- Auditory phonetics is concerned with the perception, categorization, and recognition of speech sounds and the role of the auditory system

General organization of sounds

Four general classes of sounds in American English

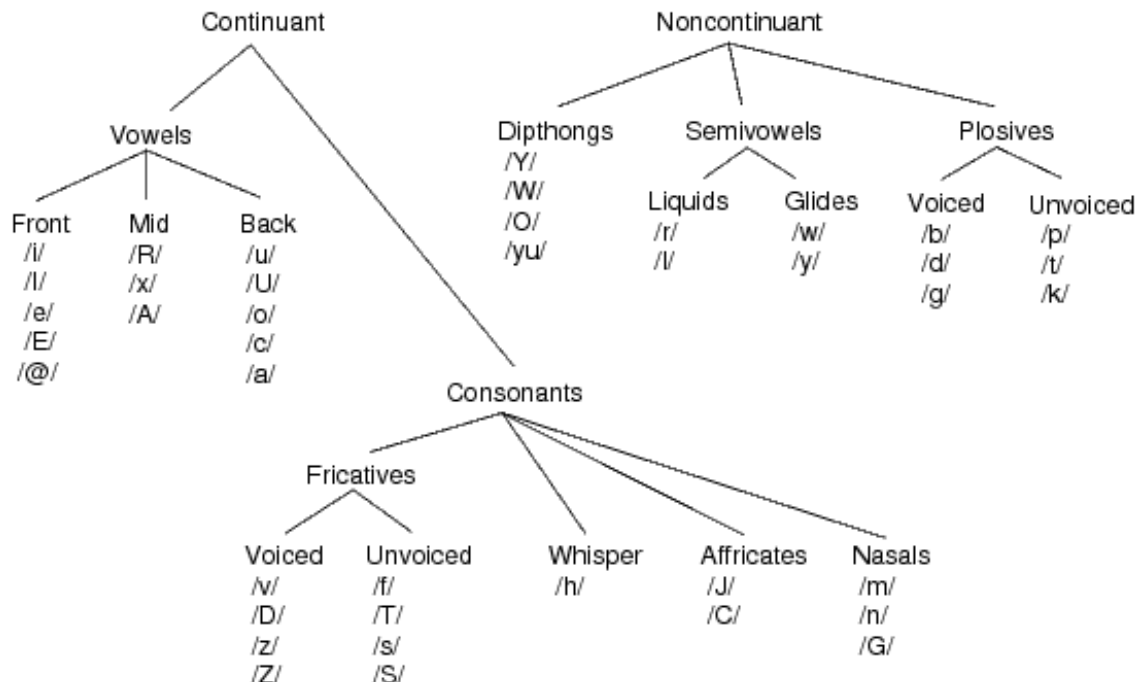
- Vowels, diphthongs, semivowels, and consonants
- Each can be further divided according to articulators (manner, place)



[Rabiner & Schafer, 1978]

Alternatively, phoneme classes can be divided into

- Continuant: produced by a fixed vocal tract configuration
 - Includes vowels, fricatives, and nasals
- Non-continuant: vocal tract configuration changes over time
 - Diphthongs, semivowels, stops and affricatives



<http://cnx.org/content/m18086/latest/phoneme.png>

Articulatory phonetics (consonants)

In terms of articulators, consonants can be described by

- Place of articulation: Defines the place of contact between an active articulator (i.e. tongue) and a passive articulator (i.e. palate)
- Manner of articulation: Concerned with airflow, the path it takes and the degree to which it is impeded
- Voicing: Determined by the behavior of the vocal folds (vibrating vs. open)

		MANNER	VOICING	PLACE					
				Bilabial	Labiodental	Interdental	Alveolar	Palatal	Velar
OBSTRUENTS	Stop	Voiceless	p			t		k	ʔ
		Voiced	b			d		g	
	Fricative	Voiceless		f	θ	s	ʃ		h
		Voiced		v	ð	z	ʒ		
	Affricate	Voiceless					tʃ		
		Voiced					dʒ		
SONORANTS	Nasal	Voiced	m			n		ŋ	
	LIQUID	Lateral	Voiced				l		
		Rhotic	Voiced					r	
	Glide	Voiced	w				j	w	

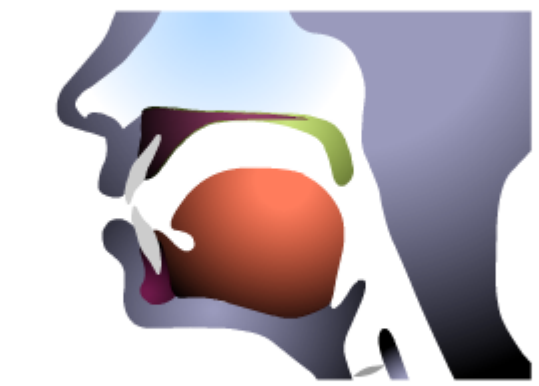
<http://www.speech-language-therapy.com/pvm.pdf>

Phonetics: The Sounds of American English

consonants — manner — place — voice — vowels — monophthongs — diphthongs
voiced — voiceless

Voiced

- | | | | |
|------|---|-----------|--------|
| /b/ | <input checked="" type="checkbox"/> /d/ | Stop | |
| /g/ | | | |
| /v/ | /ð/ | Fricative | |
| /z/ | /ʒ/ | | |
| /dʒ/ | | Affricate | |
| /m/ | /n/ | /ŋ/ | Nasal |
| /l/ | /r/ | | Liquid |
| /w/ | /j/ | | Glide |



/d/ play

animation with sound step-by-step description

fonetiks anatomy feedback

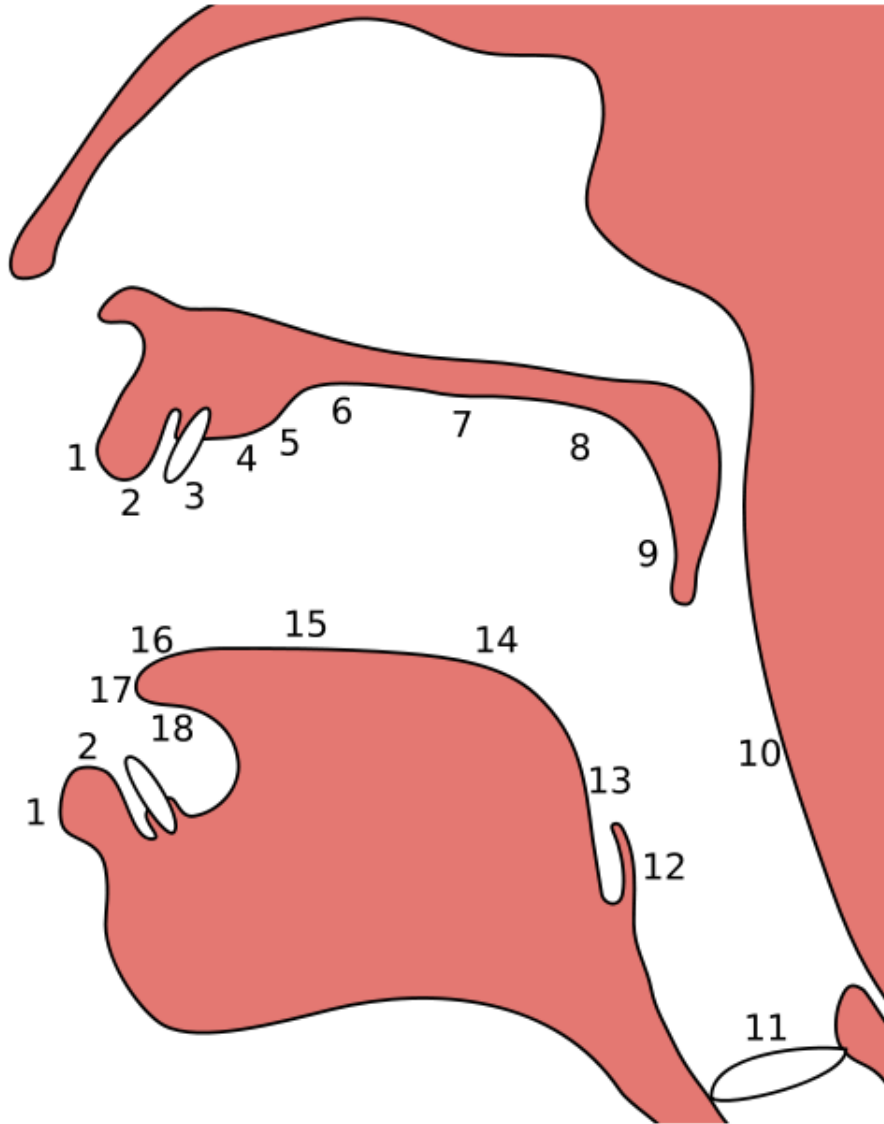


/d/

- deer
- radar
- bleed

Place of articulation

- **Bilabial**: constriction at the lips: [b], [m]
- **Labiodental**: Lower lips against upper teeth: [f], [v]
- **Interdental**: constriction between the teeth: [θ] *thing*, [ð] *that*
- **Alveolar**: constriction is at the alveolar ridge: [t], [n], [z]
- **Palatal-alveolar**: constriction slightly behind alveolar ridge: [ʃ] *sherry*, [ʒ] *measure*
- **Palatal**: constriction in the hard palate: [j] *joke*
- **Velar**: constriction closer to the soft palate: [k], [ŋ] *sing*
- **Labiovelar**: constriction both at lips and velum: [w]
- **Glottal**: when closure occurs as far back as the glottis: glottal stop [ʔ], as in the negative utterance uh-uh
- **Uvular**: constriction in the uvula; none in English, French /r/ in *rouge*



Places of articulation

1. Exo-labial
2. Endo-labial
3. Dental
4. Alveolar
5. Post-alveolar
6. Pre-palatal
7. Palatal
8. Velar
9. Uvular
10. Pharyngeal
11. Glottal
12. Epiglottal
13. Radical
14. Postero-dorsal
15. Antero-dorsal
16. Laminal
17. Apical
18. Sub-apical

http://en.wikipedia.org/wiki/Place_of_articulation

Manner of articulation

- **Stops:** produced by complete stoppage of the airstream: [p]
- **Fricatives:** tongue comes very close to a full closure: [f], [sh] *sherry*
- **Affricatives:** combination b/w stops and fricatives: *cherry*
- **Nasals:** closed oral passage (as in stops), open nasal cavity: [n], [ng] *sing*
- **Approximants** : halfway between consonants and vowels
 - **Liquids:** [l], [r]
 - **Glides:** [y], [w]

Voicing

- When the vocal folds vibrate, it is voiced; otherwise it is voiceless
- Examples of voiceless/voiced: *Sue* vs. *zoo*, *pat* vs. *bat*

Articulatory phonetics (vowels)

Vowels can be described in a similar way

- Manner of articulation, just considered to be “vowel”
- Place of articulation is generally described with three major parameters: frontness, height, and roundness

Frontness (or backness)

- Provides a general indication of the greatest place of constriction, and correlates with F2
- Three positions in English
 - Front: [iy] *beat*, [ih] *bit*, [eh] *bet*, [ae] *bat*
 - Central: “schwa” [ax] *about*
 - Back: [uw] *boot*, [ao] *bought*, [ah] *but*, [aa] *father*

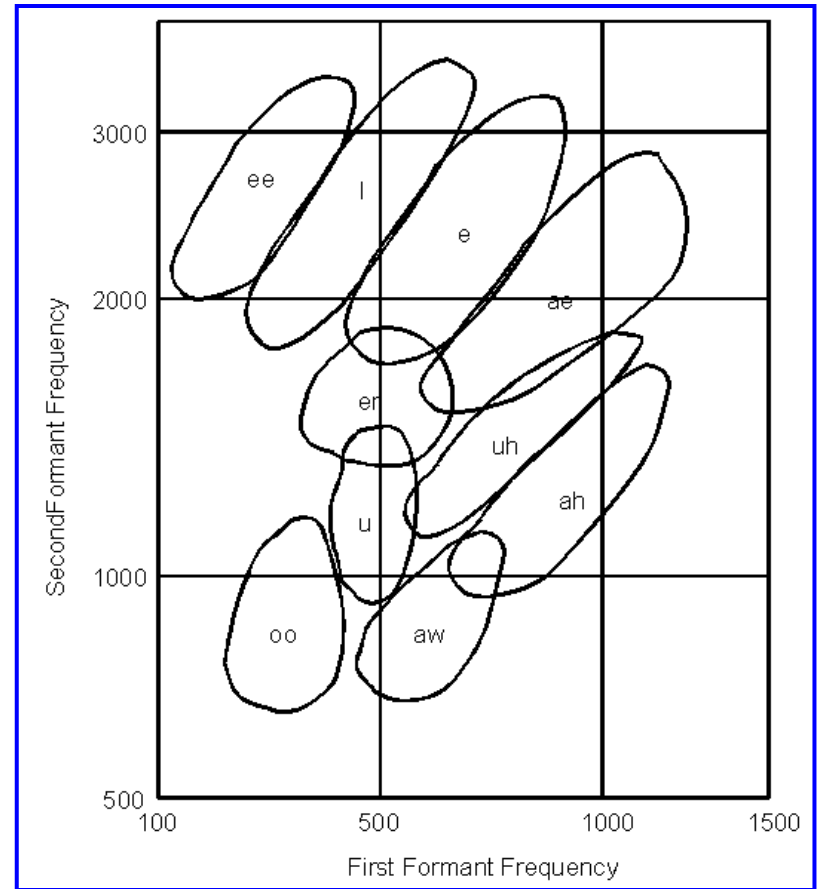
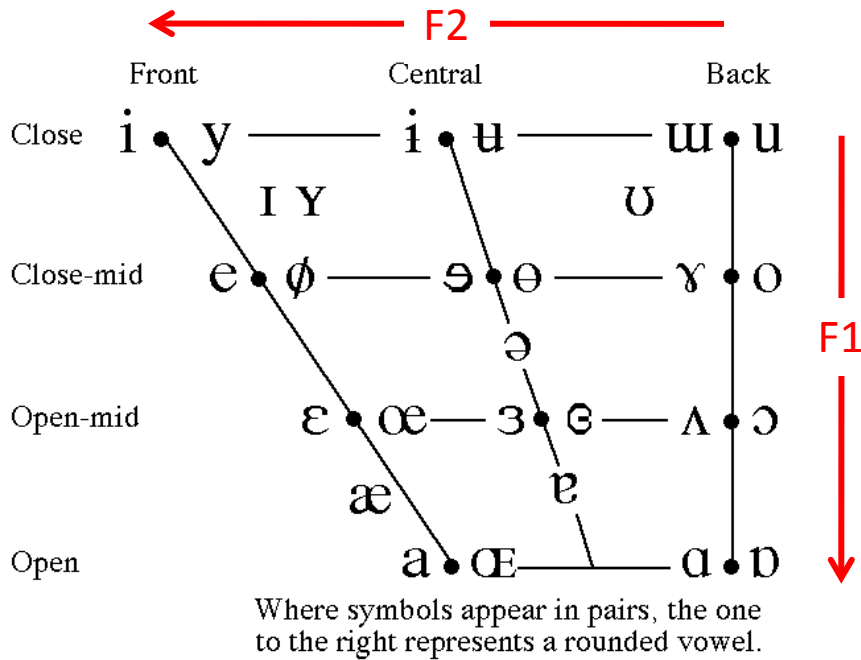
<http://www.utexas.edu/courses/linguistics/resources/phonetics/vowelmap/index.html>

Height

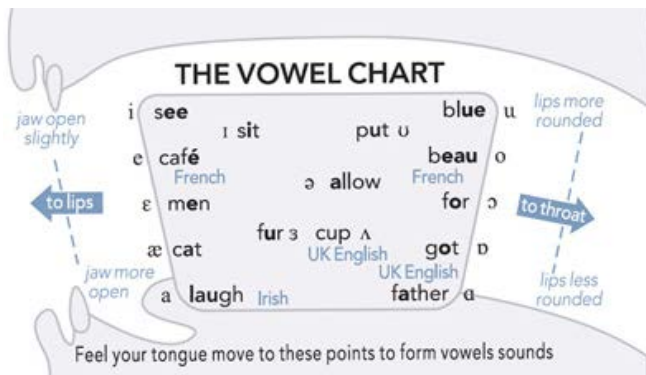
- Refers to how far lower jaw is from upper jaw when making the vowel
 - High vowels have lower and upper jaw close: [iy], [uw]
 - Low vowels have a more open oral cavity: [ae], [aa]
- Correlates with F1 (high vowel: low F1; low vowel: high F1)

Roundness

- Refers to whether the lips have been rounded as opposed to spread
- In English, front vowels are unrounded whereas back vowels are rounded: *bit* vs. *boot*



<http://www.singwise.com/images/CardinalVowelChart.gif>



http://www.geofex.com/Article_Folders/sing-wah/sing-wah.htm

http://www.thedialectcoach.com/images/content/vowel_chart.jpg

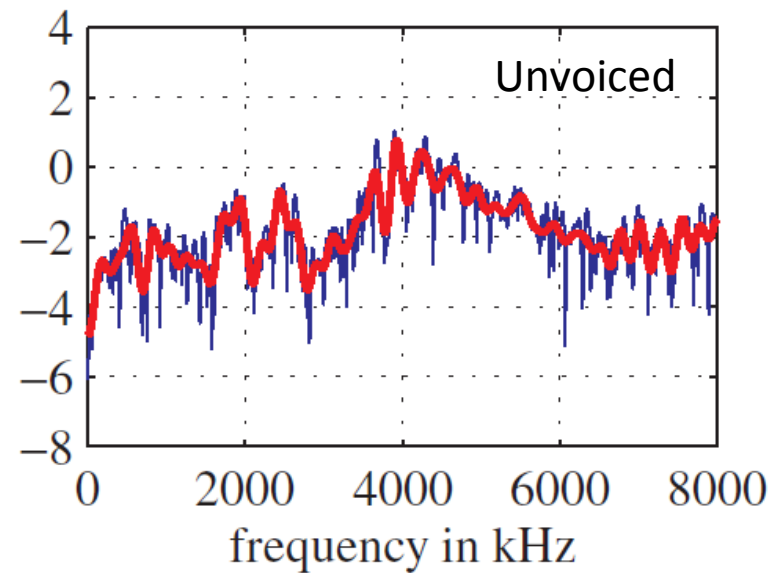
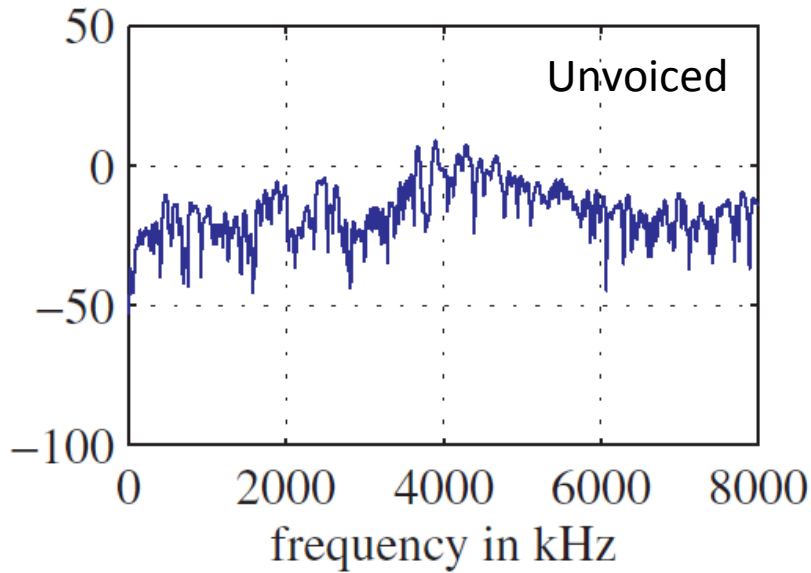
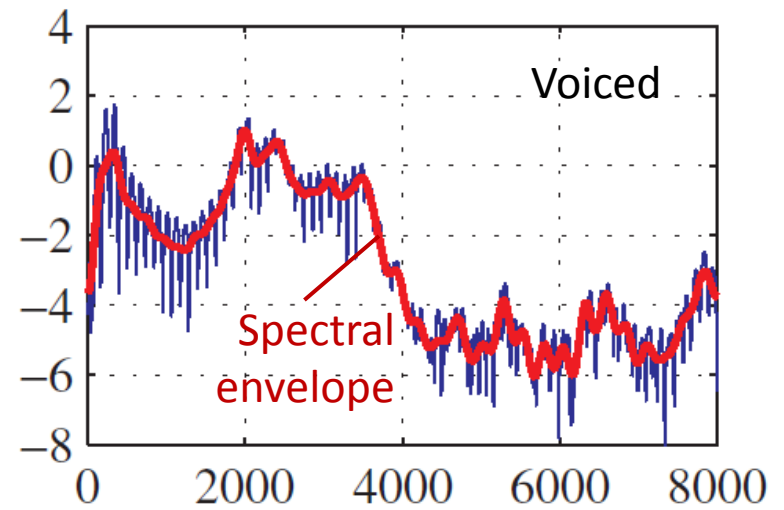
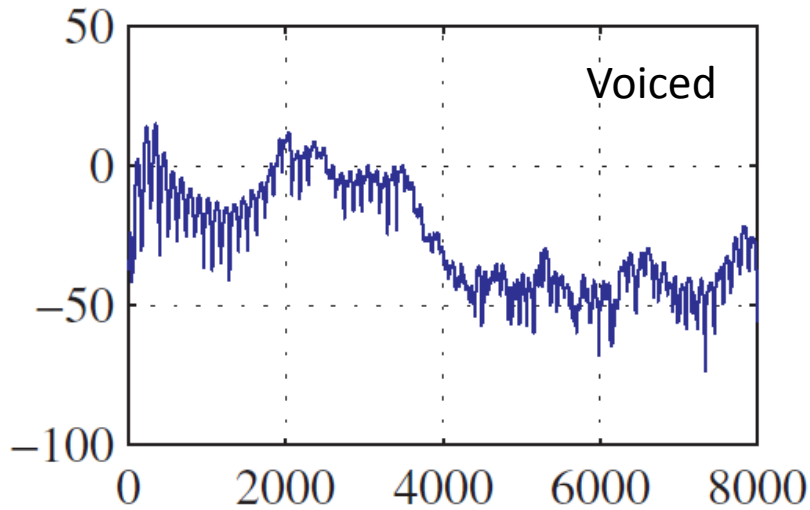
Acoustic phonetics

Acoustic phonetics is concerned with

- Time domain waveform of the speech signal, and
- Its time-varying spectral characteristics

Visualizations of speech waveforms

- Time-domain waveforms are rarely studied directly
 - This is because phase differences can significantly affect its shape but are rarely relevant for speech perception
- Instead, frequency-domain signals are commonly used
- The spectrum (log-magnitude) of a voiced phone shows two types of information
 - A comb-like structure, which represent the harmonics of F_0 (the source),
 - A broader envelope, which represents the resonances (formants) of the vocal tract filter
- Various techniques exist to separate the two sources of information
 - Linear prediction, homomorphic (cepstral) analysis ...



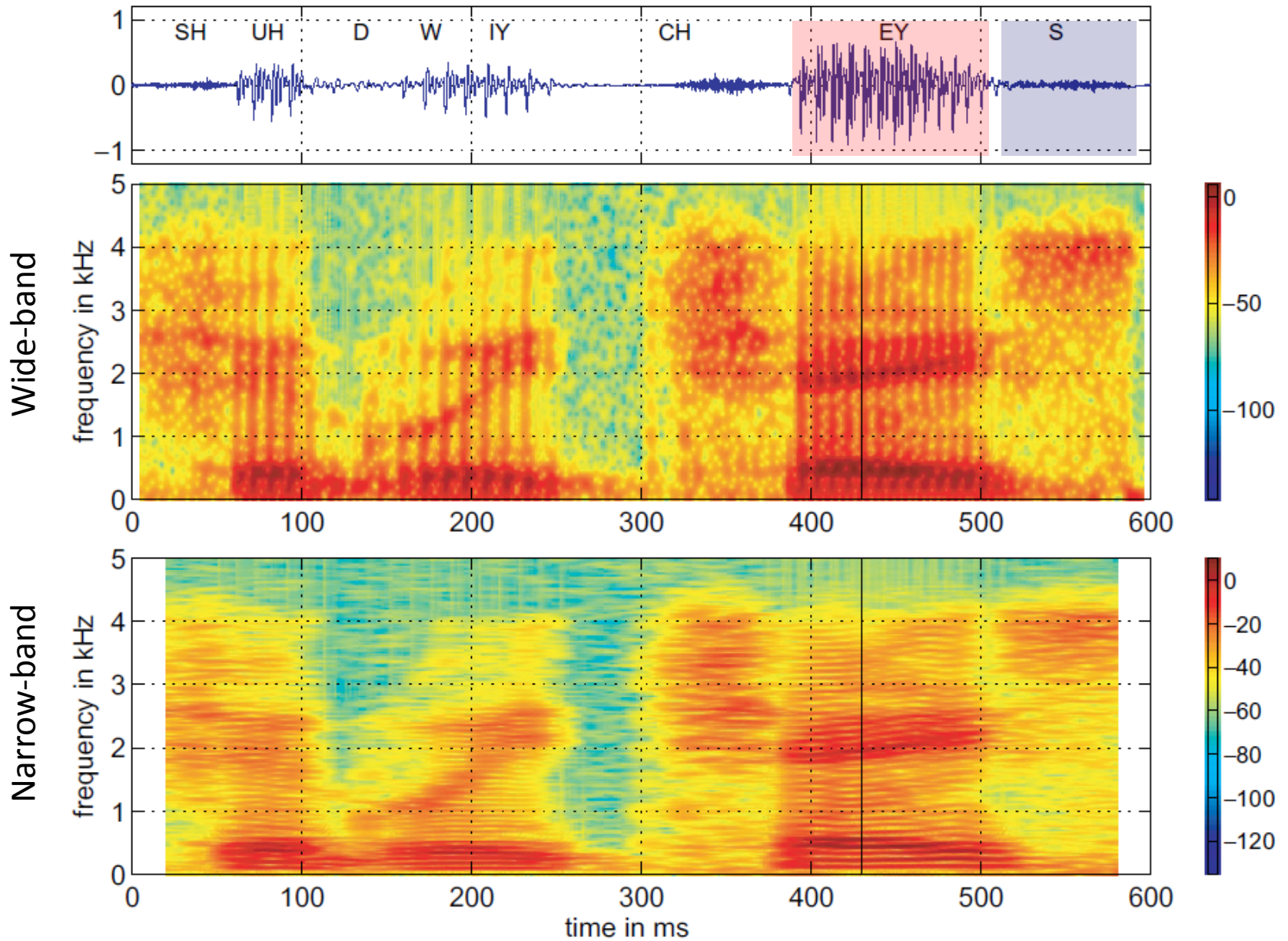
The spectrogram

- Can be thought of as a moving spectrum over time
- Typically represented as a 3D graphic where
 - Horizontal dimension represents time,
 - Vertical dimension represents frequency, and
 - Color represents magnitude (typically in log scale)

Two general types of spectrograms

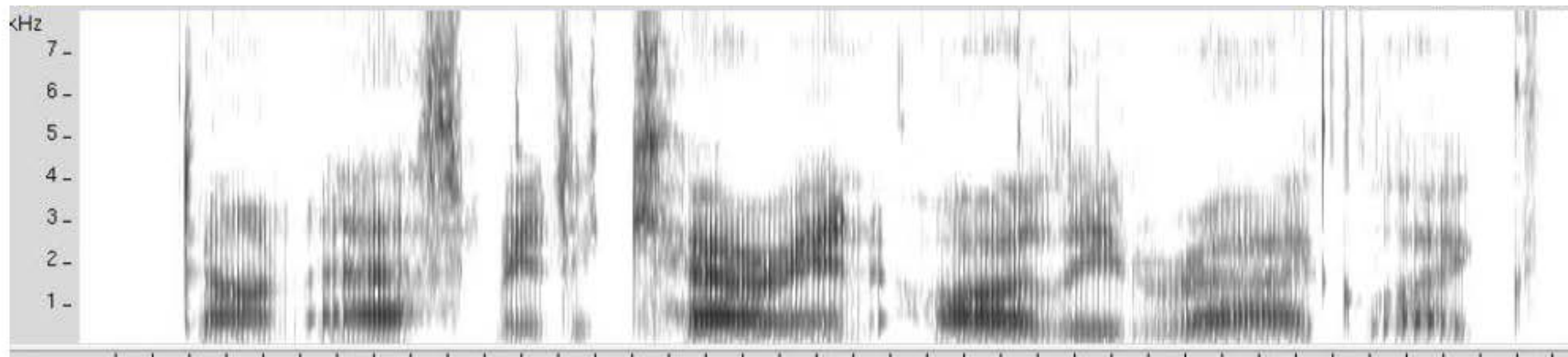
- Wideband
 - Computed over a short window of time (e.g., 5 ms)
 - High temporal resolution, poor frequency resolution
 - Vertical striations represent individual pitch periods (for voiced phones)
- Narrowband
 - Computed over a relatively large window of time (e.g., 25 ms)
 - High frequency resolution, poor temporal resolution

"Should we chase." Voiced Unvoiced

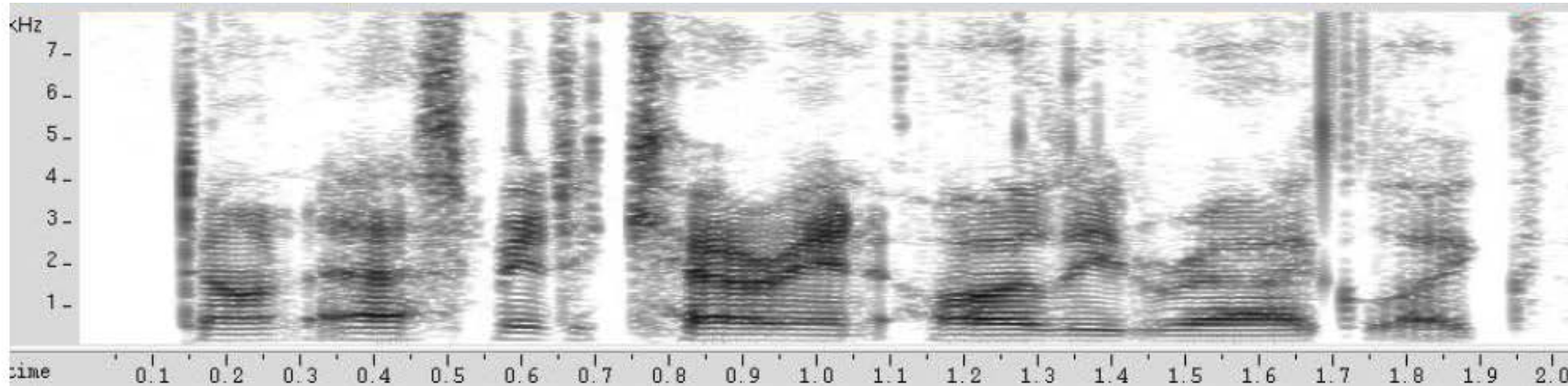


[Rabiner & Schafer, 2007]

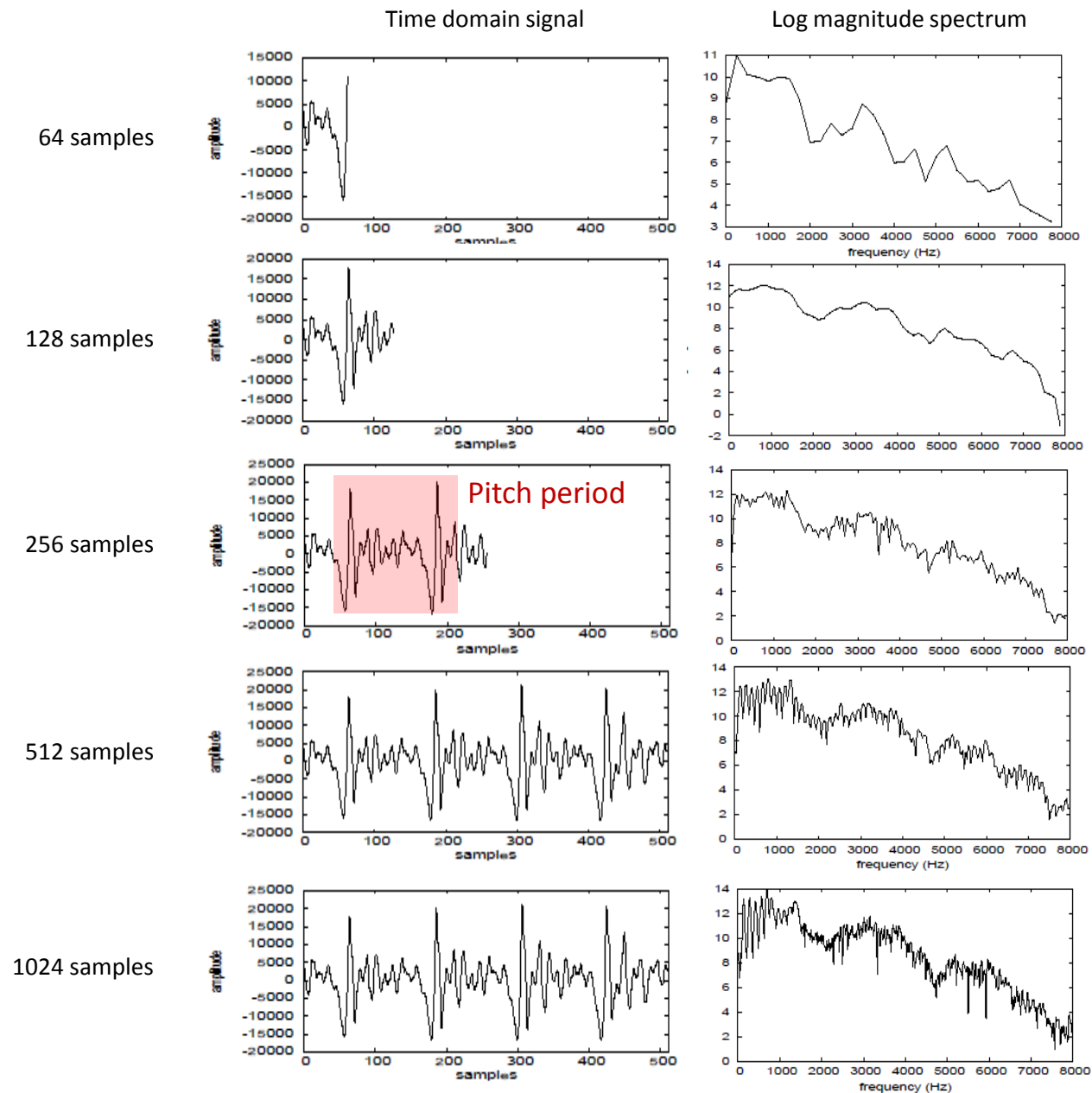
Wide-band



Narrow-band



[Taylor, 2009]



[Taylor, 2009]

Vowels

- The largest phoneme group and most interesting one
 - Carry little information in written speech, but most ASR systems rely heavily on them for performance
- Vowels are voiced (except when whispered) and have the greatest intensity and duration in the range of 50 to 400ms
- Vowels are distinguished mainly by their first three formants
 - However, there is a significant individual variability, so other cues can be employed for discrimination (upper formants, bandwidths)

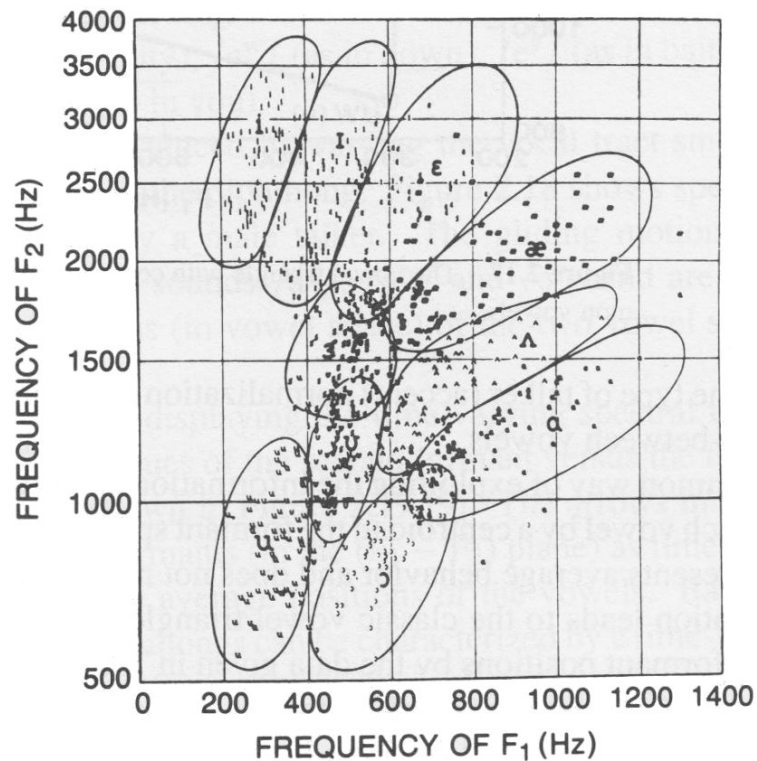
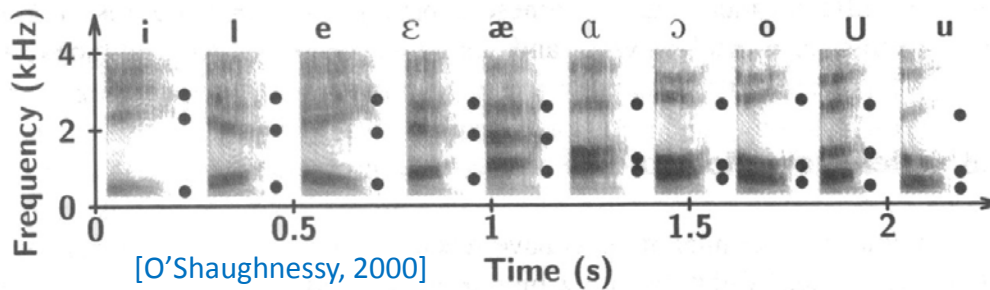


Figure 2.16 Measured frequencies of first and second formants for a wide range of talkers for several vowels (after Peterson & Barney [7]). [Rabiner & Juang, 1993]

Complete the following sections of text

Th_y n_t_d s_gn_f_c_nt _mpr_v_m_nts _n th_
c_mp_ny's _m_g_, s_p_rv_s_n, th_r wrk_ng
c_nd_t__ns, b_n_f_ts _nd _pp_rt_n_t__s

A__i__u__e__ _o__a__ _ay __a__e__ e__e__ia__y __e
__a__e, __i__ __e__ _o__e__ o__ o__u__a__io__a__ e__oyee__
__i__y__ _e__ea__i__

They noted significant improvements in the company's image, supervision, their working conditions, benefits and opportunities

Attitudes toward pay stayed essentially the same, with the scores of occupational employees slightly decreasing

Diphthongs

- Diphthongs consist of a dynamic vowel sound in which the tongue and lips move between two vowel positions
- Three diphthongs are universally accepted
 - These can be described as a vowel followed by a glide, though the articulators rarely reach the full constriction of a glide
 - Examples: [aj] bite, [oy] boy, [aw] bout

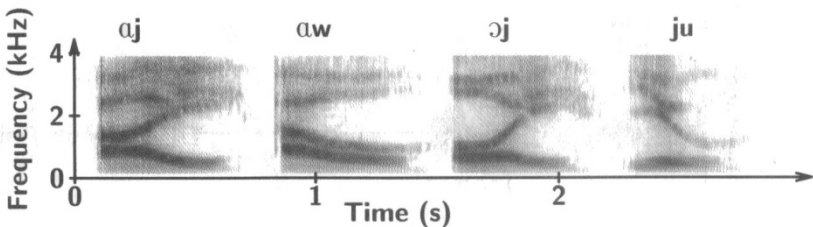
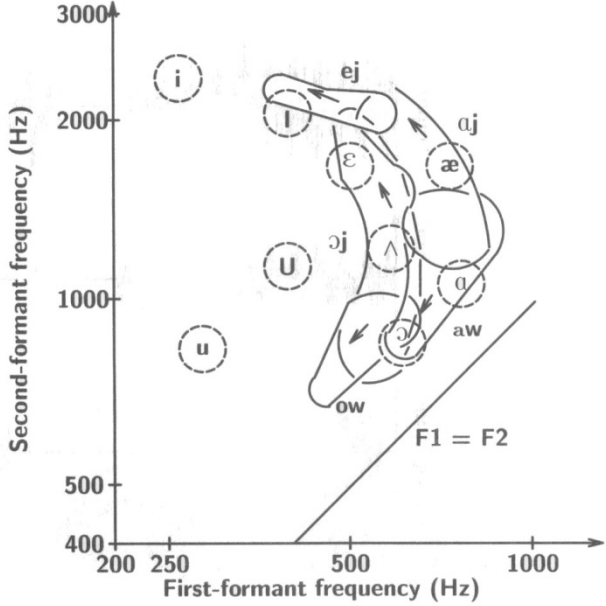


Figure 3.16 Spectrograms of four English diphthongs.

[O'Shaughnessy, 2000]

Figure 3.15 F1 and F2 time movements for diphthongs. The dashed circles indicate typical F1-F2 positions for vowels, and the solid contours enclose formant trajectories during diphthongs. (After Holbrook and Fairbanks [78].)

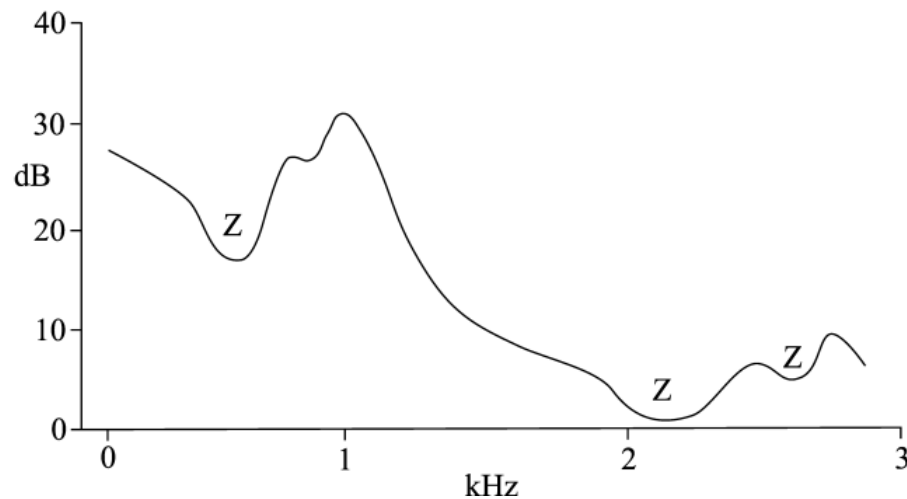
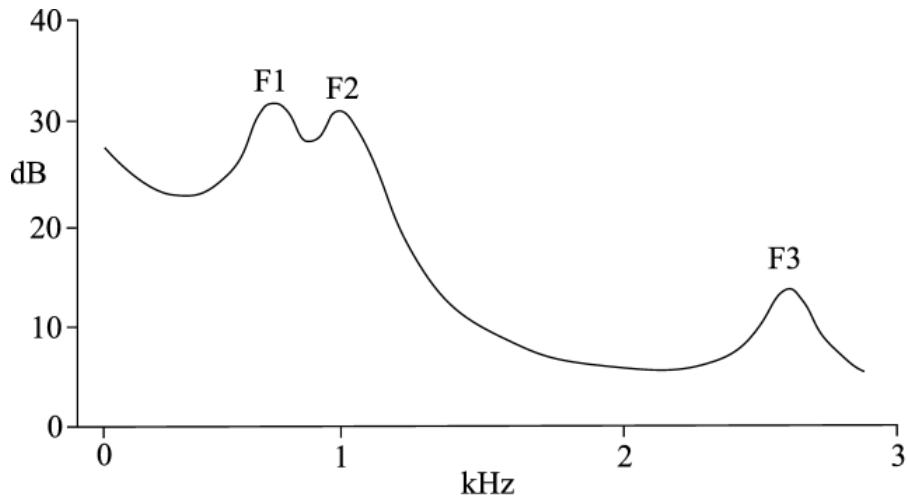


Glides and liquids

- Very similar to vowels: periodic, intense waveforms: most energy in low formant; hence, they are also known as semivowels
 - However, they are weaker than vowels since they require a greater constriction of the vocal tract
 - They are also transitional sounds, so their characteristics are strongly influenced by context
- Liquids ([l], [r]) tend to color preceding vowels with their sound
- Glides ([y], [w]) often lengthen vowels to create diphthongs

Nasals

- Produced with glottal excitation (i.e., voicing)
- Vocal tract totally constricted at some point in the oral cavity
 - Thus, the mouth serves as a resonant cavity that traps energy at certain acoustic frequencies, which appear as anti-resonances (zeros of the transfer function)
- Velum is lowered, so air flows through the nasal cavity
 - Since nasal cavity has large area, formant bandwidths are generally broader than for other sonorants
- Nasal consonants are distinguished by their place of constriction
 - [m] at the lips
 - [n] at the back of the teeth
 - [ŋ] forward of the velum



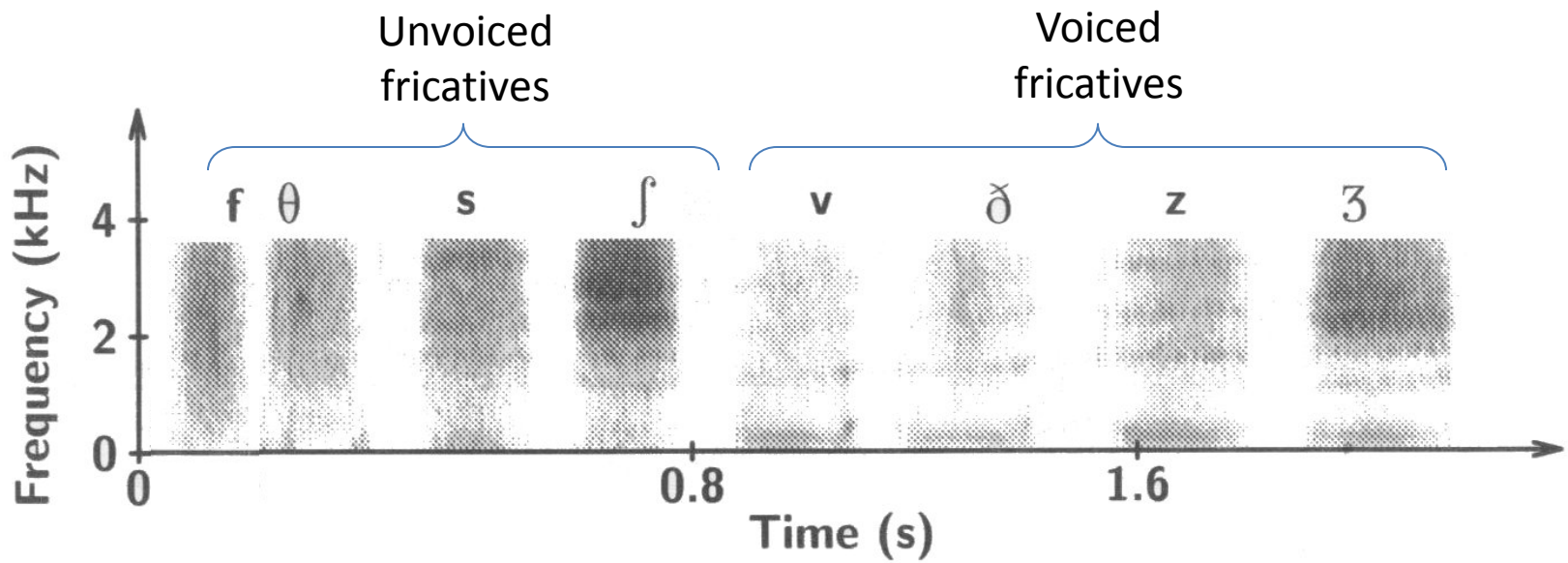
Effect of vowel nasalisation

- Spectral envelope of a normal [a] and a nasalised [a]
- The first three formants are labeled in the normal vowel
- Reductions in spectrum level in the nasalised vowel are due to the addition of zeroes (labeled as "z")

http://clas.mq.edu.au/physiology/nasal/nasality_review.html

Fricatives

- Characterized by a major place of constriction
- Come in two flavors: voiced and unvoiced
 - Place of constriction is identical for corresponding pairs of phones
- Unvoiced
 - Vocal tract is excited by a steady air flow, which becomes turbulent at a place of constriction
 - Sorted by place of constriction: [f] *fish*, [θ] *thin*, [s] *sound*, and [ʃ] *shout*
 - Spectrum looks “noisy”
- Voiced
 - Glottis is vibrating, but the airflow becomes turbulent near the constriction, so these sounds display two distinct components
 - Sorted by place of constriction : [v] *vote*, [ð] *then*, [z] *zoo*, [ʒ] *azure*
 - Spectrum contains both “noisy” and harmonic components



[O'Shaughnessy, 2000]

Stops (plosives)

- Unlike other sounds, which can be described by steady-state spectra, stops are transient phones
 - They are produced by building up pressure behind a constriction and then suddenly releasing it
- As with fricatives, stops come in two flavors: voiced and unvoiced
 - Place of constriction is identical for corresponding pairs of phones: bilabial, alveolar, velar
- Voiced
 - Vocal folds vibrate while pressure builds up, which can be heard as a low-frequency energy radiating through the walls of the throat (voice bar)
 - Sorted: [b] *bee*, [d] *day*, [g] *guy*
- Unvoiced
 - Vocal folds do not vibrate as pressure builds up
 - Sorted: [p] *pea*, [t] *tea*, [k] *key*

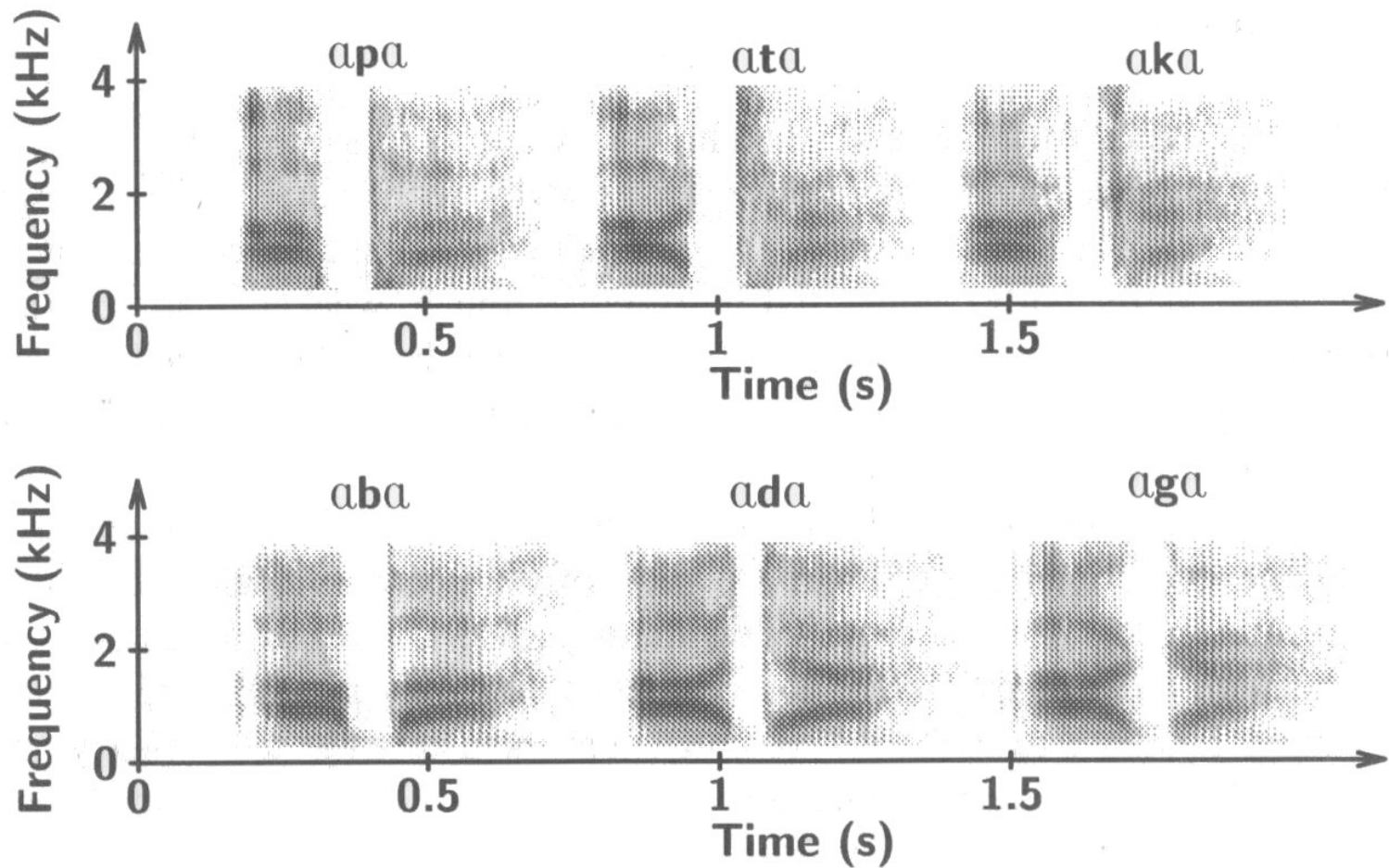


Figure 3.21 Spectrograms of the six English stops in vocalic context:
/apa,ata,aka,aba,ada,aga/.

[O'Shaughnessy, 2000]

Speech perception

Vowels

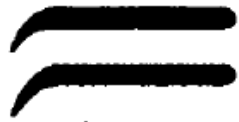








- Vowel perception is relatively simple: formant frequencies are the main factors in vowel identification
- However, formant frequencies scale with vocal tract length
 - There is evidence that listeners “normalize” formant location by making formant spacing essential features of vowel identification
- Vowel nasalization is cued primarily by
 - Increase in the bandwidth of F1, and
 - The introduction of zeros

Consonants

- More complex, and depends on a number of factors
 - Formant transition into the following vowel
 - Formant location
 - Voice onset time
 - Voicing

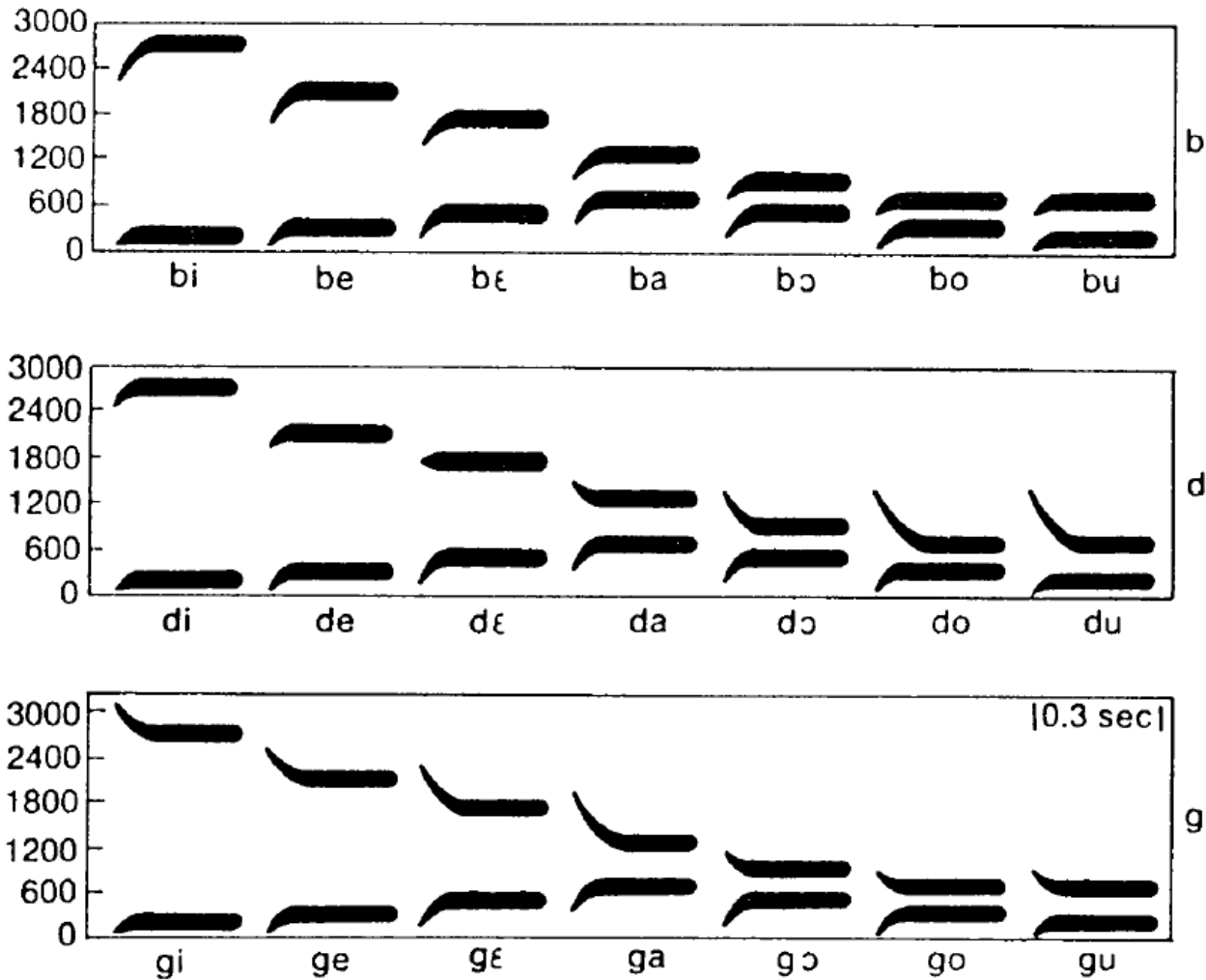
- Direction of formant transition
 - See examples in the next few slides
- Rate of formant transition
 - It is possible to transform the perception of a plosive into a semi-vowel by decreasing the formant transition rate
 - Example: contrast between [b] and [w]
 - On the word ‘be’, the transition between [b] and [i] is 10 ms
 - As the transition increases beyond 30 ms, ‘be’ is transformed into ‘we’
- Formant locus
 - Vocal tract configuration in front of the closure (for stops)
- Voice onset time (VOT)
 - Length of time between release of a closure and the start of voicing
 - Critical for the perception of stop consonants
 - Example: contrast between [t] and [d]
 - On the word ‘do’, segment the sound [d] and increase the delay wrt [o]
 - When VOT exceeds about 25 ms, the word is perceived as ‘to’

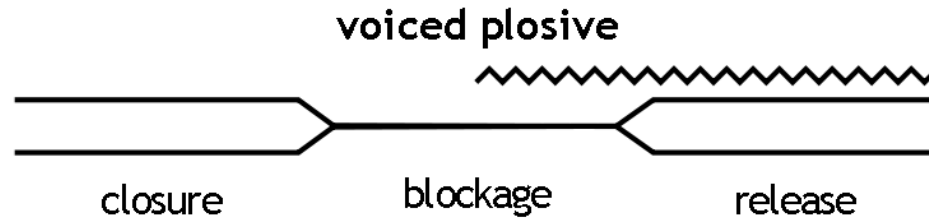
Schematic formant transition patterns varying in place and manner

	FRONTAL	MIDDLE	BACK
VOICED STOPS	 <p>ba</p>	 <p>da</p>	 <p>ga</p>
UNVOICED STOPS	 <p>pa</p>	 <p>ta</p>	 <p>ka</p>
NASALS	 <p>ma</p>	 <p>na</p>	 <p>ŋa</p>

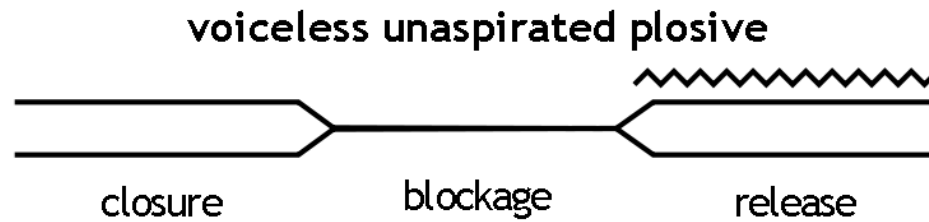
<http://www.phon.ucl.ac.uk/courses/spsci/b214/week2-5.pdf>

Schematic formant transition patterns for voiced stop-vowel syllables

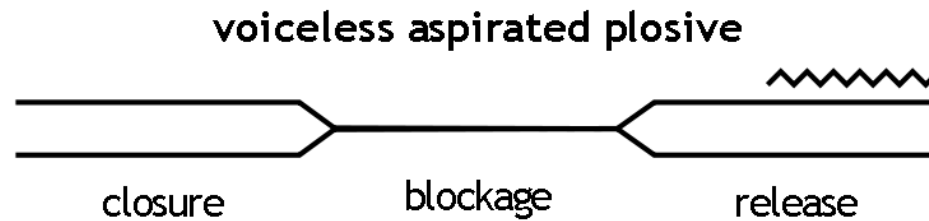




Negative VOT
[b,d,g]



Zero VOT
[p,t,k]



Positive VOT
[p^h,t^h,k^h]

http://en.wikipedia.org/wiki/Voice_onset_time

Prosody

Long-term variations (over more than one phoneme) in

- Pitch (intonation)
- Amplitude (loudness), and
- Timing (articulation rate or rhythm)

Roles of prosody

- Helps highlight the spoken message
 - Alternation of stressed and unstressed syllables identifies the words that the speaker considers more important
- Helps segment the spoken message
 - Provides cues to syntactic boundaries (e.g., main vs. subordinate clauses) and syntactic structure (e.g., declaratives statements vs. questions)
 - Serves as a “continuity guide” to track speakers in noisy environments
- Provides cues to the state of the speaker
 - F0 and amplitude patterns vary with emotions
- **Interestingly, however, prosody is typically ignored in ASR**