

Applied Human Language Technology

Lecture 7

Voiceprints & Acoustic Processing of Speech

Irene Murtagh

This week:

- Introduction to the acoustic processing of speech signals
(The basis of speech recognition by computers)
- Signal Analysis
- Feature Extraction
- Fourier Analysis and Linear Predictive Coding (LPC)
- Spectral Analysis and Spectra: **Human Voiceprints**
- Sound Waves
- Interpreting a waveform
- Some things to do

Acoustic Processing of Speech

This lecture presents a brief overview of the kind of acoustic processing commonly called **signal analysis** or **feature extraction**.

The term **features** refers to the **vector of numbers** which represents one time slice of a **speech signal**.

A number of kinds of features are commonly used,
e.g. **LPC** features.

These are **spectral features**, which means that they represent the waveform in terms of the distribution of different frequencies that make up the waveform.

Such a distribution of frequencies is called a **spectrum**.

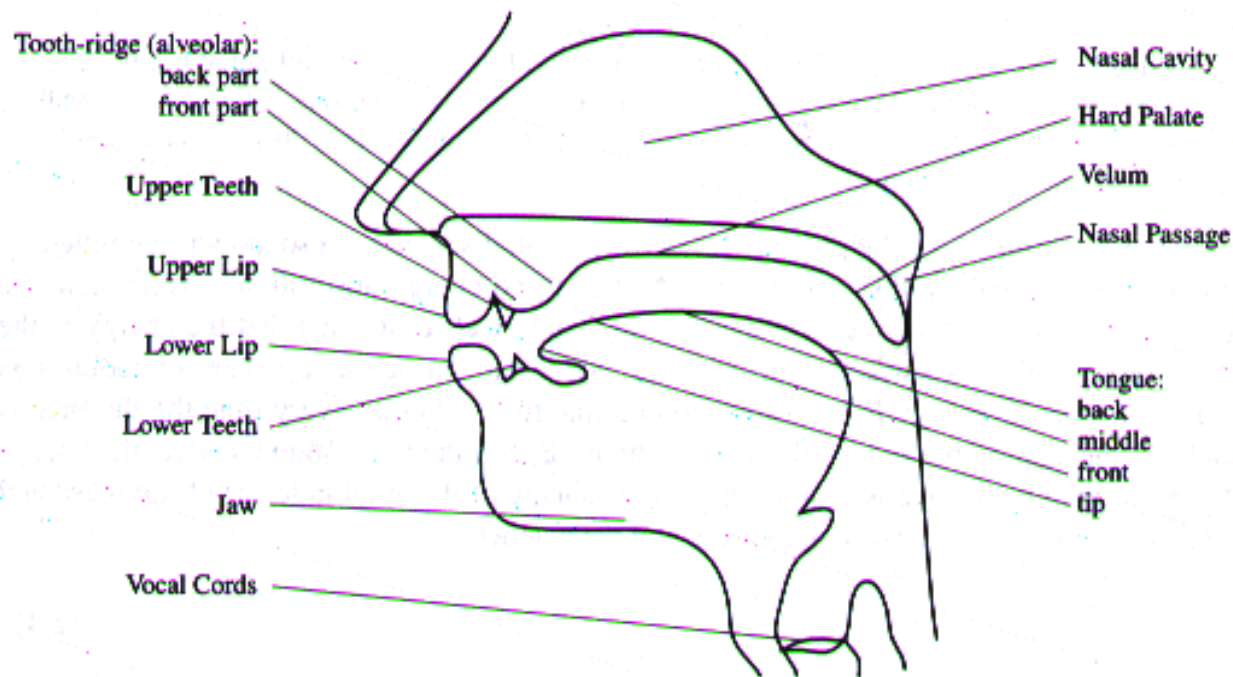
In this lecture we:

summarise the **idea of frequency analysis and spectra**, and
sketch out **different kinds of extracted features**.

Sound Waves

The input to a **speech recogniser**, like the input to the human ear, is a complex series of changes of air pressure.

These changes in air pressure originate with the speaker, and are caused by the specific way that air passes through something called the glottis and out the oral or nasal cavities.



A schematic diagram of the human speech production apparatus.

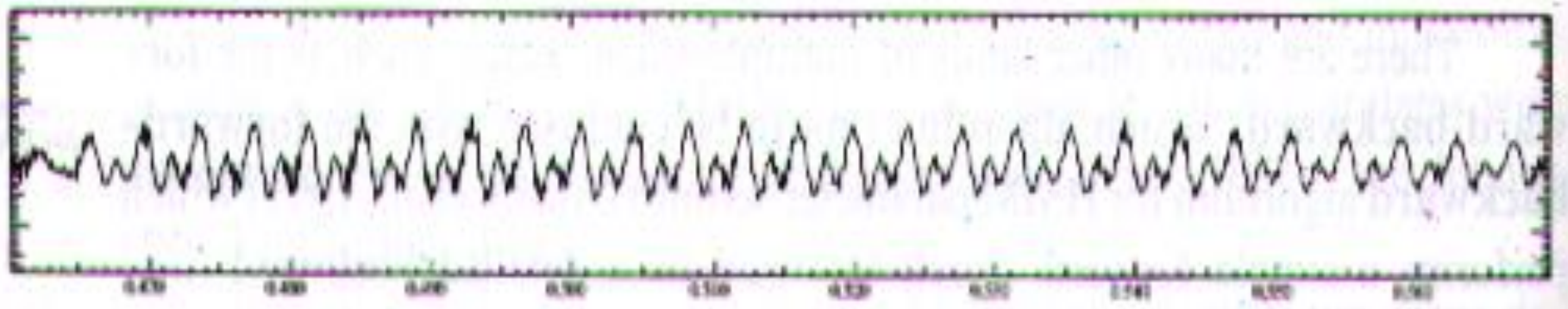
We represent sound waves by plotting the change in air pressure over time.

One way of visualising this is to imagine a graph plot of a vertical plate which is blocking the air pressure waves

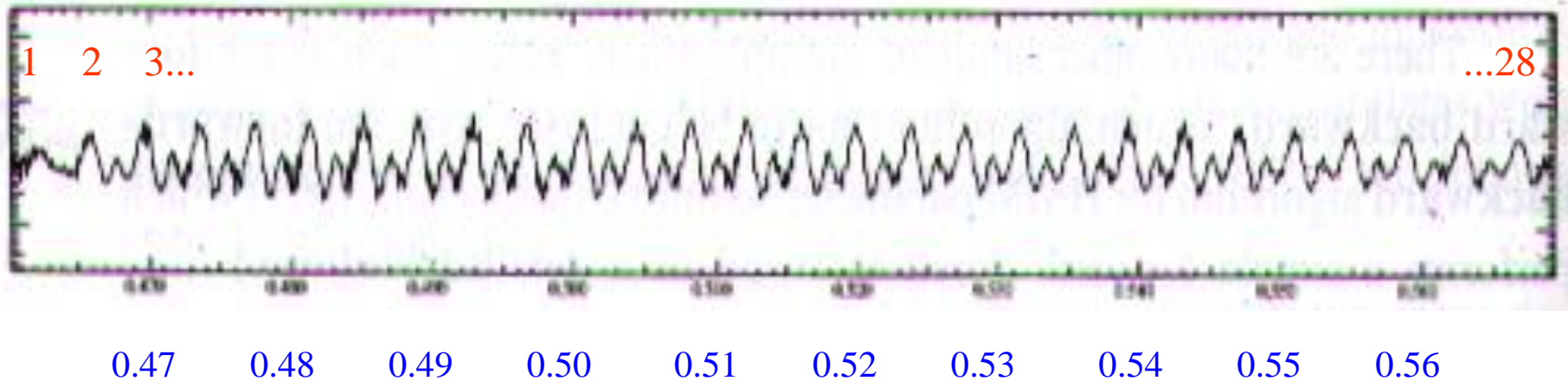
i.e. a microphone in front of the speaker, or the eardrum of the hearer.

The graph measures the amount of compression of the air molecules at this plate.

The diagram following, from the course textbook by Jurafsky and Martin, shows the waveform taken from a corpus of telephone speech of someone saying “*she just had a baby*”.



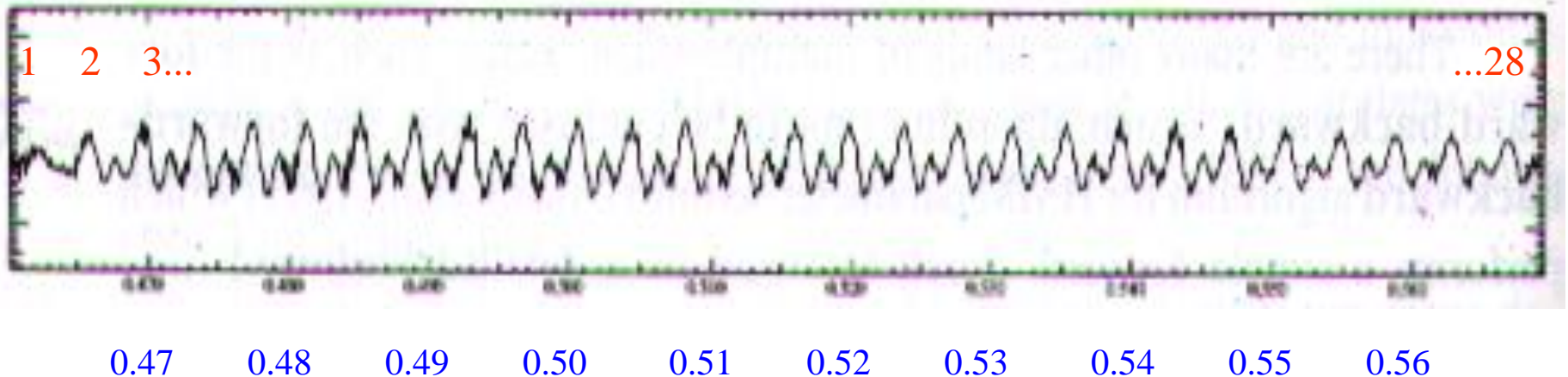
0.47 0.48 0.49 0.50 0.51 0.52 0.53 0.54 0.55 0.56



Two important characteristics of a wave are its:

1. frequency and
2. amplitude.

- The **frequency** is the number of times a second that a wave repeats itself, or more technically, cycles.

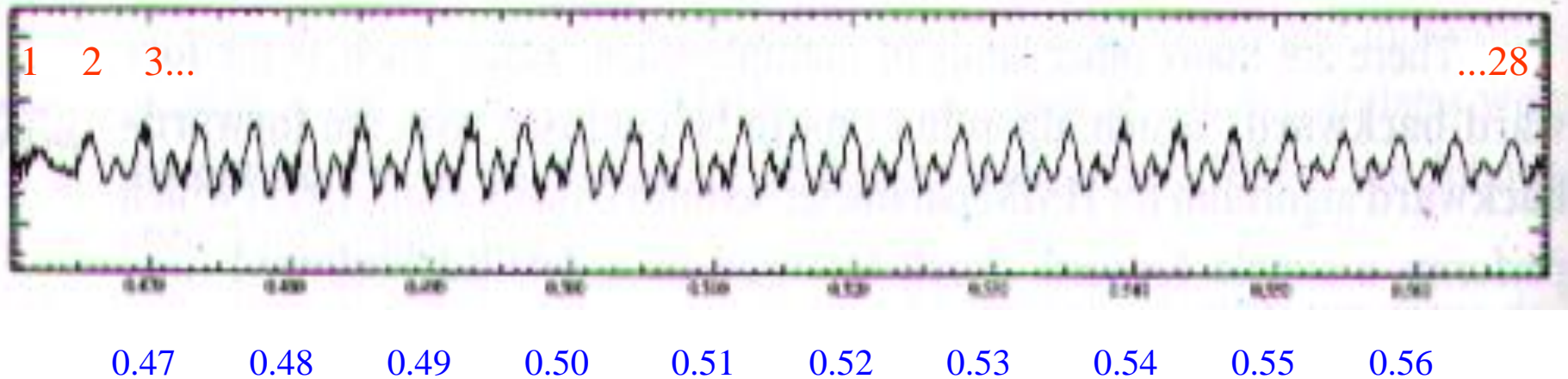


In the diagram there are **28 repetitions** of the wave in the **0.11 seconds** captured.

Therefore, the frequency of this segment of the wave is **28/0.11** or **255 cycles per second**

Frequency = repetitions/time-period

Cycles per second are called **Hertz** (Hz for short). The frequency in the diagram is **255Hz**.



The vertical axis in the diagram measures the amount of air pressure variation.

- A **high value** on the vertical axis (a high **amplitude**) indicates that there is more air pressure at that point in time,
- a **zero value** means that there is normal (atmospheric) pressure,
- a **negative value** means there is lower than normal air pressure.

Two important perceptual properties are related to frequency and amplitude.

The **pitch** of a sound is the perceptual correlate of frequency.

In general, if a **sound has a higher frequency we perceive it as having a higher pitch**, but the relationship is not linear since human hearing has different acuities for different frequencies.

Similarly, the loudness of a sound is the perceptual correlate of the power, which is related to the square of the amplitude.

Sounds with higher amplitude are perceived as louder, but again, the relationship is not linear.

How to interpret a waveform

Since humans (... and computers) can transcribe and understand speech just given the sound wave, **the waveform must contain enough information to make this task possible.**

In most cases, this information is hard to unlock just by looking at the waveform, but ... we can still learn many things by a visual inspection of the waveform.

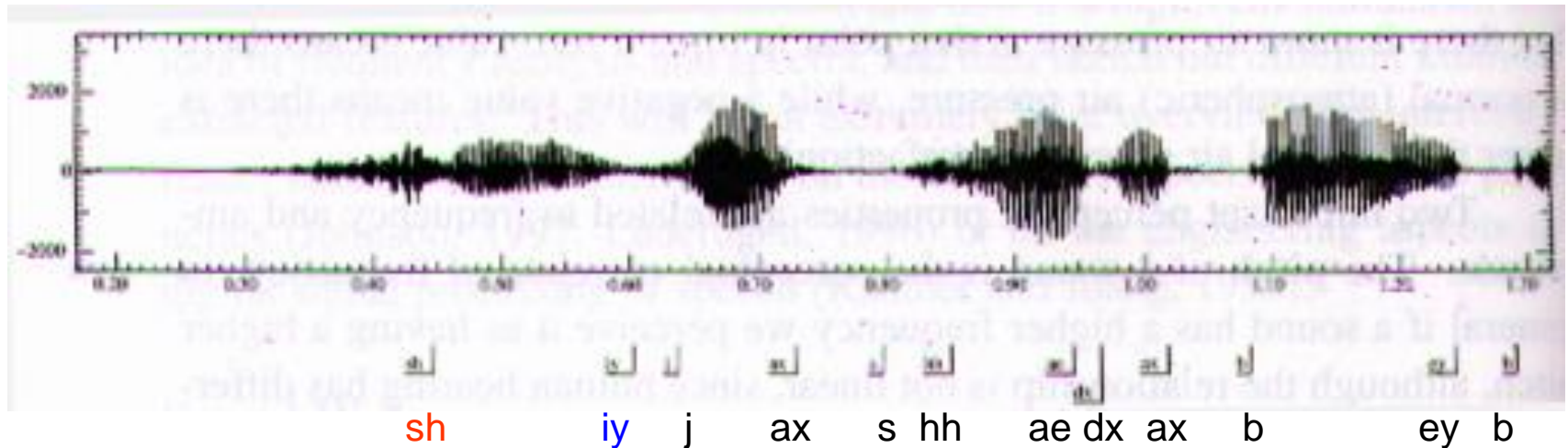
For example, the difference between vowels and consonants of spoken language is quite clear on a waveform

Vowels tend to be long and relatively loud.

Length in time manifests itself as high amplitude.

Fricatives [sh] can also be recognised in a waveform. They produce an intense irregular pattern

The diagram is a waveform by a 20 year female, speaking with an accent from the south midlands of the USA.

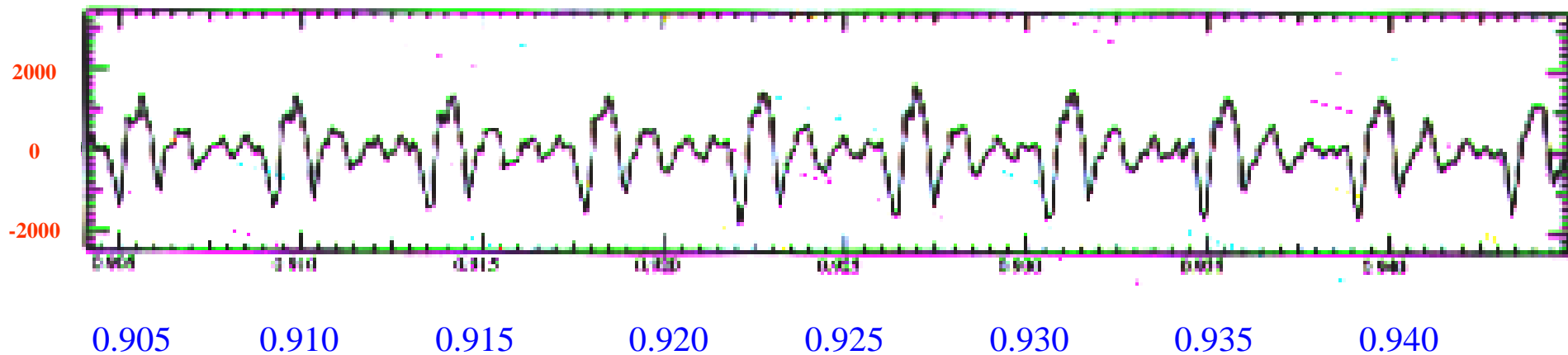


Spectra

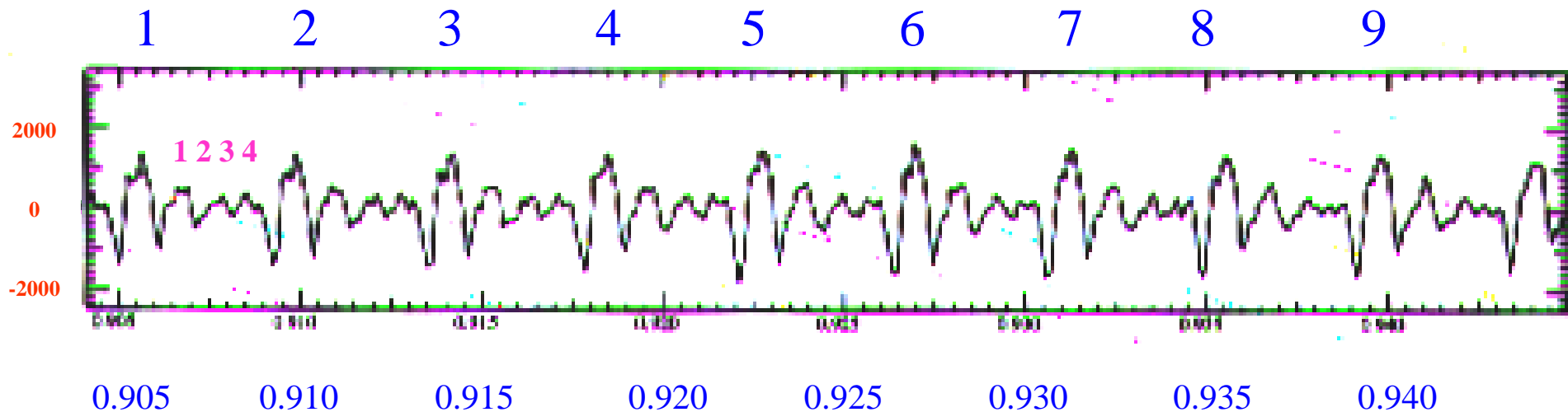
While some broad **phonetic** features can be interpreted from a waveform, more detailed classification requires a different representation of the input in terms of spectral features.

Spectral features are based on the insight of Fourier that every complex wave can be interpreted as a sum of many simple waves of different frequencies.

A musical analogy is a chord. Just as a chord is composed of multiple notes, any waveform is composed of the waves corresponding to its individual “notes”



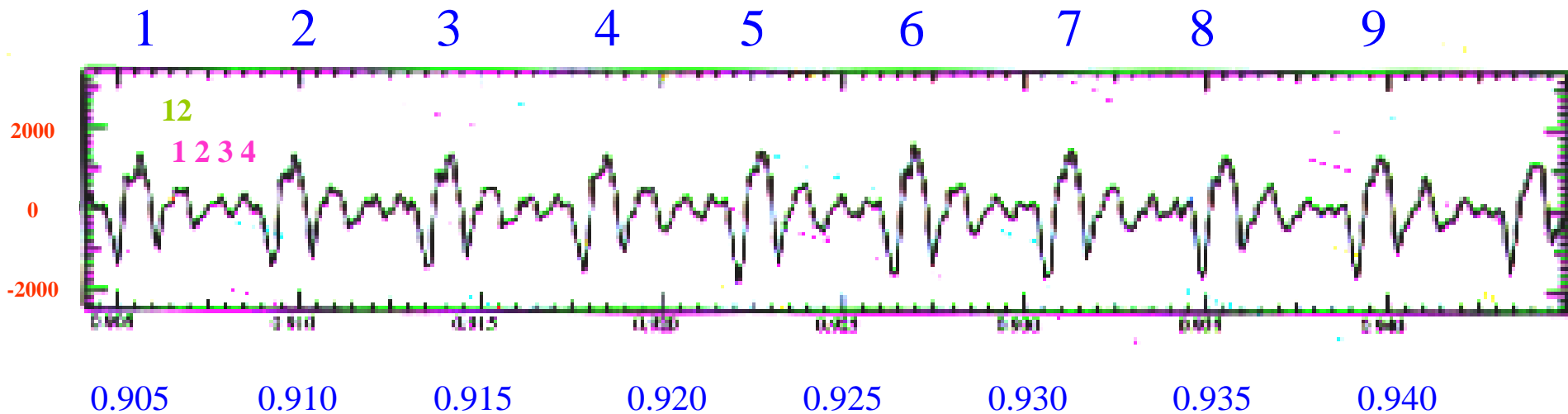
The diagram shows part of the waveform for the vowel [Q] ...in American speech... of the word **had** at **second 0.9** of the sentence "She**h**ad a baby".



Note however, that there is a complex wave which repeats about 9 times in the diagram.

Note that there is also a **smaller repeated wave which repeats 4 times** for every larger pattern (... look at the 4 small peaks inside every repeated wave).

- The complex wave has a frequency of about 250 Hz.
- We can figure this out since it repeats about 9 times in 0.036 seconds giving $9 \text{ cycles} / 0.036 = \mathbf{250 \text{ Hz}}$.

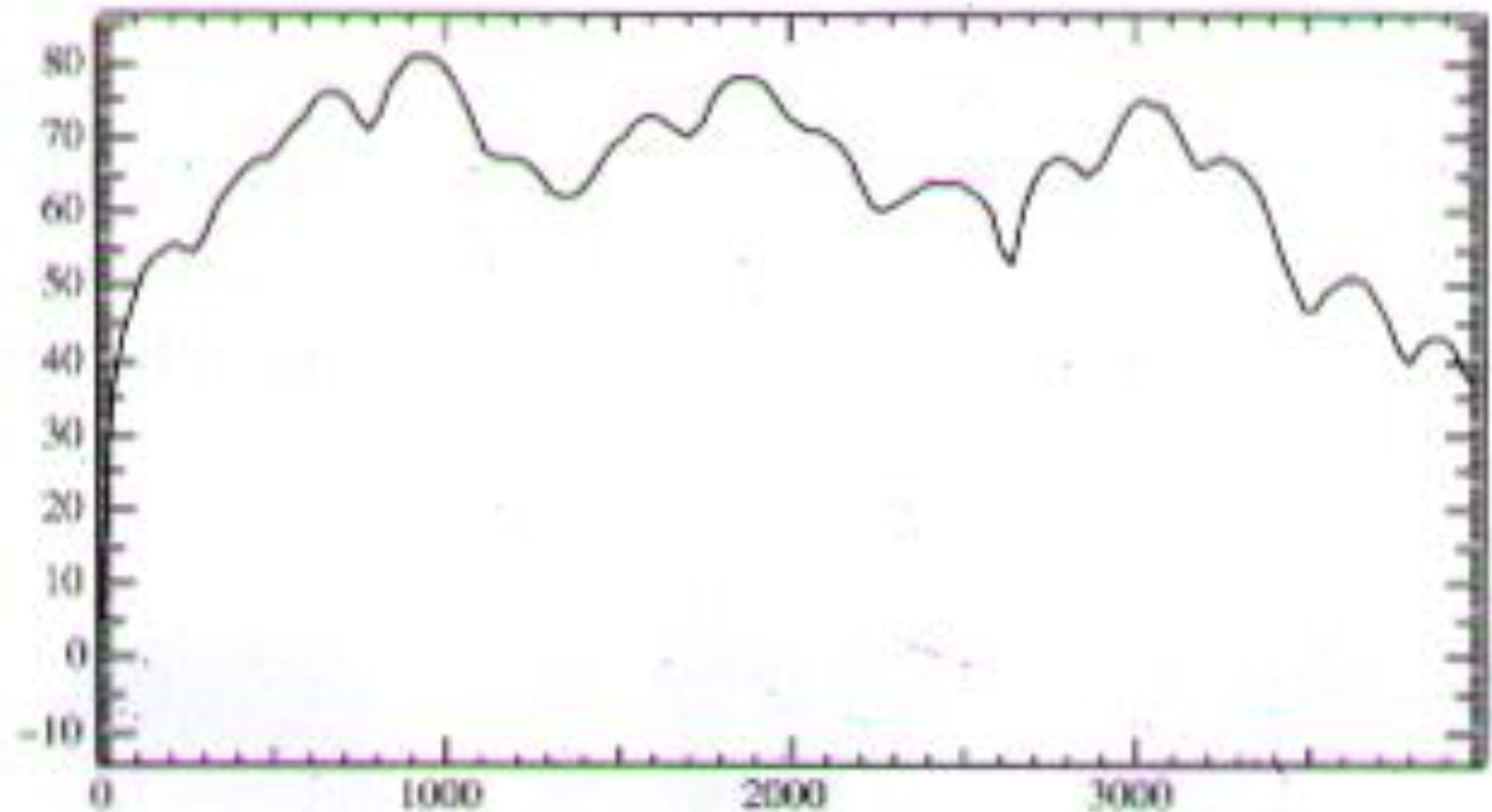


The smaller wave should then have a frequency of approx. 4 times the frequency of the larger wave, say **1000Hz**.

- Then, if you look very closely, you can see two little lines on the peak of many of the 1000Hz waves.
- The frequency of this tiny wave must be approx. 2 times that of the 1000Hz wave = about **2000Hz**.

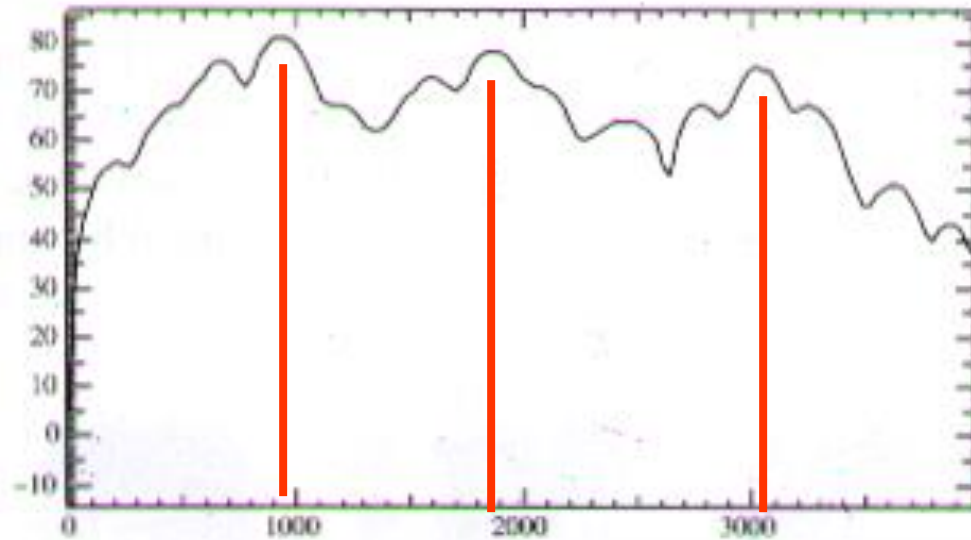
A spectrum is a representation of these different frequency components of a wave.

- It can be computed by a Fourier transform, a mathematical procedure which separates out each of the frequency components of a wave.
- Many speech applications use an LPC (Linear Predictive Coding) spectrum because this makes it easier to see where the peaks are.



An LPC spectrum for the vowel [Q] waveform of ÒSe just had a babyÓ at the point in time shown in the previous diagram.

LPCs makes it easier to see **formants**



The x-axis of a spectrum shows **frequency** while the y-axis shows some measure of the **magnitude** of each frequency component (in decibels **dB**, a logarithmic measure of amplitude).

The diagram shows that there are important frequency components at 930Hz, 1860Hz, and 3020Hz, along with many other lower-magnitude frequency components.

These components at approx. 1000Hz and 2000Hz are what we predicted earlier.

Why is a Spectrum Useful?

It turns out that these spectral peaks that are easily visible in a spectrum are very characteristic of different sounds.

Phones have characteristic spectral “signatures”.

For example;

Chemical elements give off different wavelengths of light when they burn, allowing scientists to detect elements in stars that are light-years away by looking at the spectrum of the light.

Similarly, by looking at the spectrum of a waveform, we can detect the characteristic signature of the different phones that are present.

This use of spectral information is essential to both **human** and **machine speech recognition**.

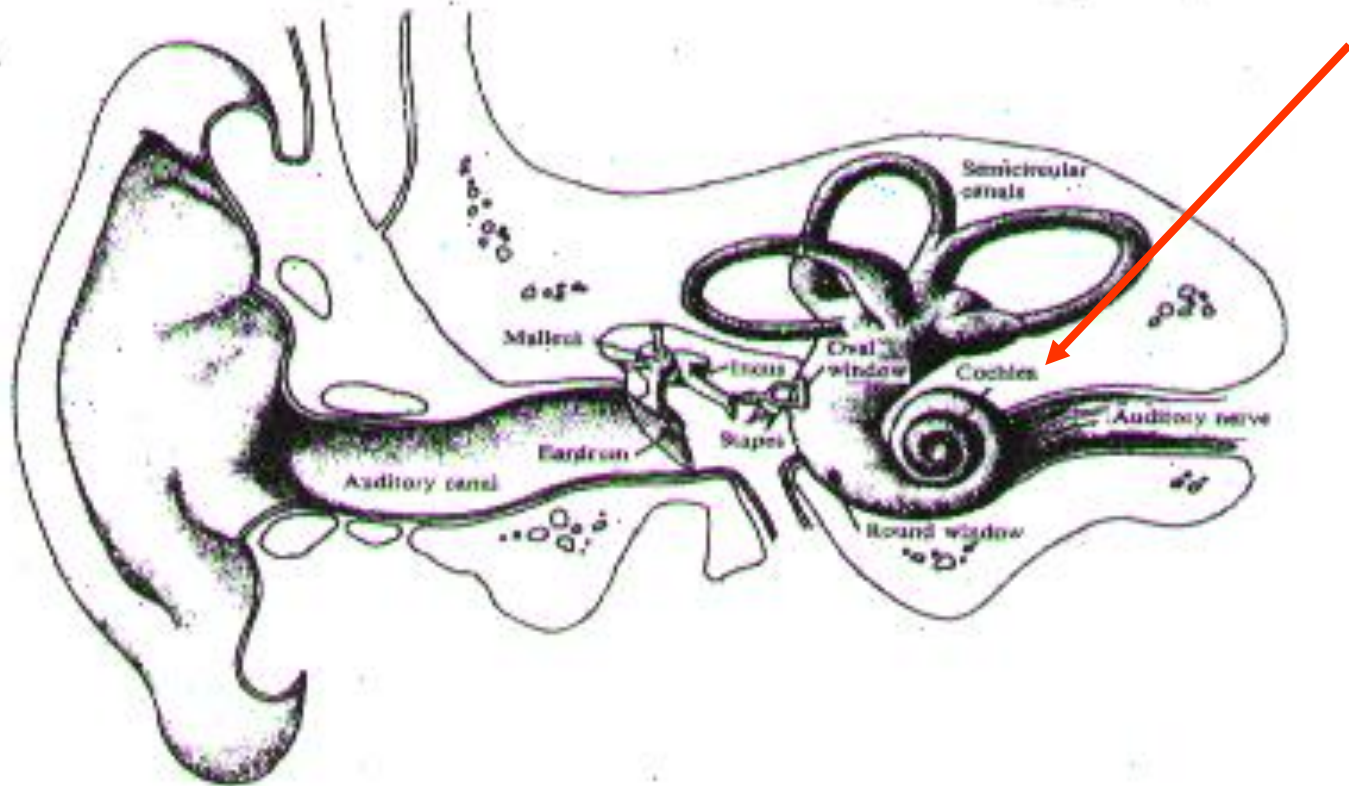
Spectrogram

While a **spectrum** shows the frequency components of a wave at one point in time, ...

a **spectrogram** is a way of envisioning how the different frequencies which make up a waveform change over time.

In human audition, the function of the **cochlea** or **inner ear** is to compute a spectrum of the incoming waveform.

Similarly, the features that are input to **HMMs in speech recognition** are all representations of spectra.

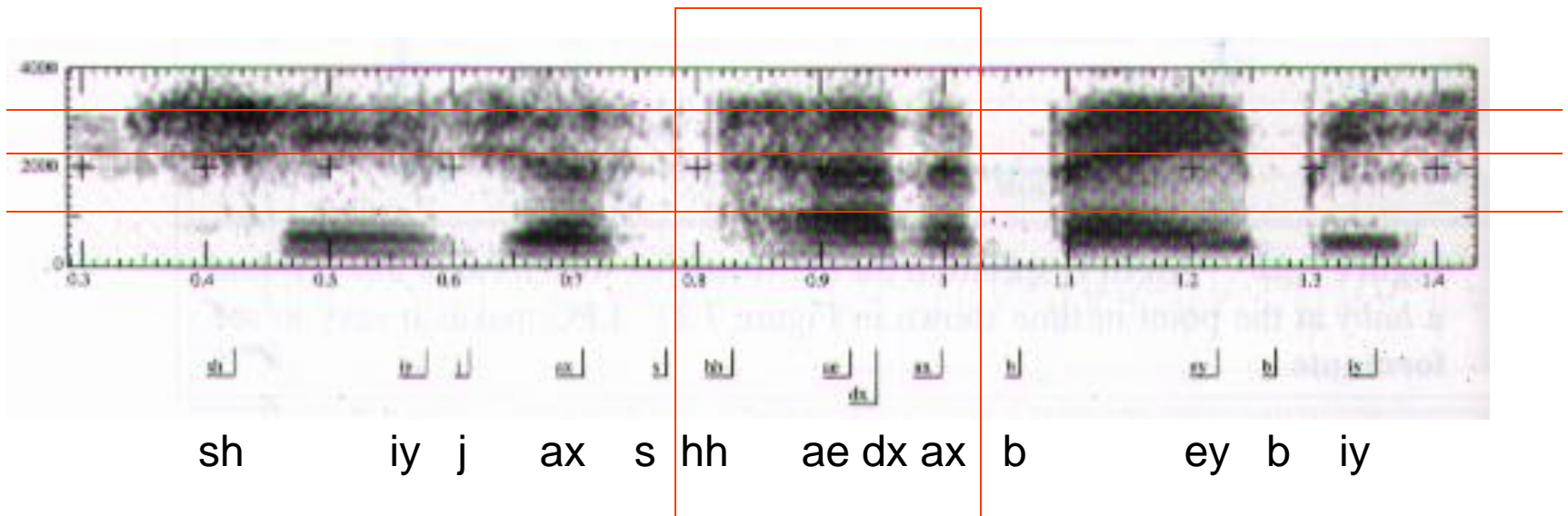


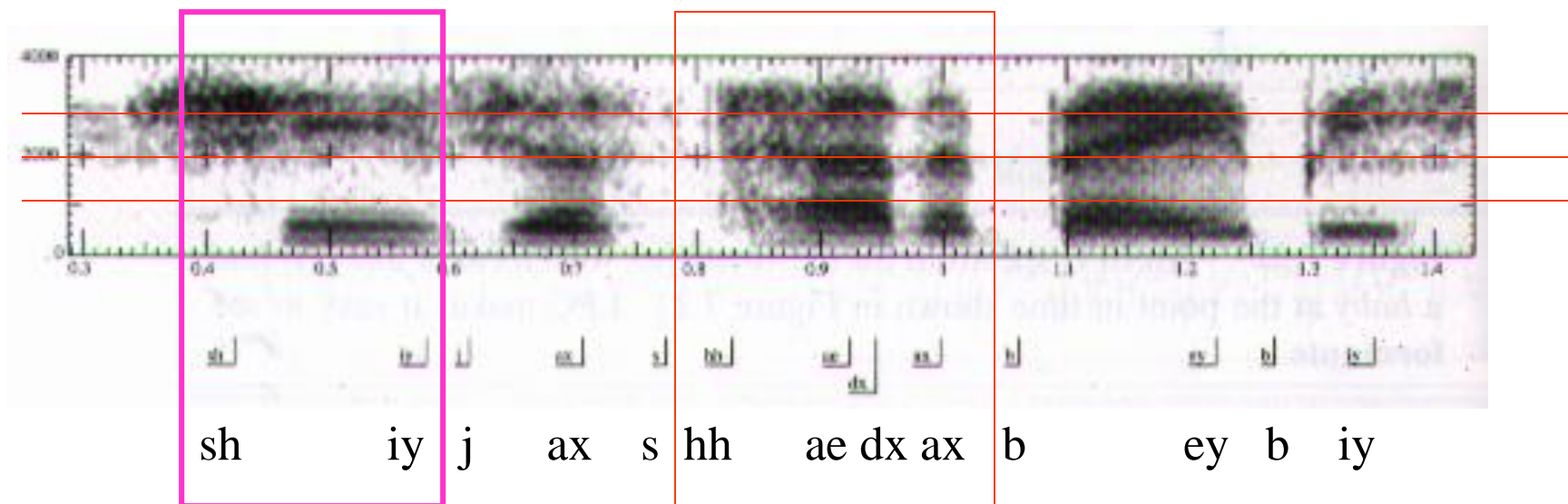
The structure of the peripheral auditory system with the outer, middle, and inner ear.

The x-axis shows time, as it did for the waveform, but the y-axis now shows frequencies in Hertz (Hz).

The darkness of a point on a spectrogram corresponds to the amplitude of the frequency of the component.

For example, in the diagram at the point in time of second 0.9, notice the dark bar at around 1000Hz.





For example, in the diagram at the point in time of second 0.9, notice the dark bar at around 1000Hz.

This means that the vowel [i_y] of the word “she” has an important component around 1000Hz.

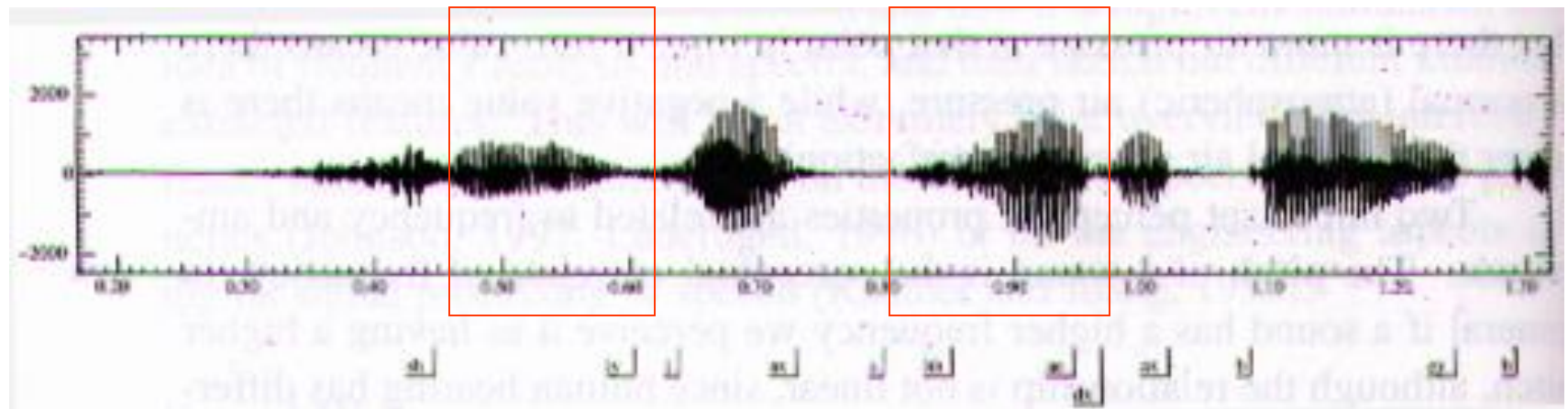
The dark horizontal bars on a spectrogram, representing spectral peaks, usually of vowels, are called **formants**.

What specific clues can spectral representations give for phone identification?

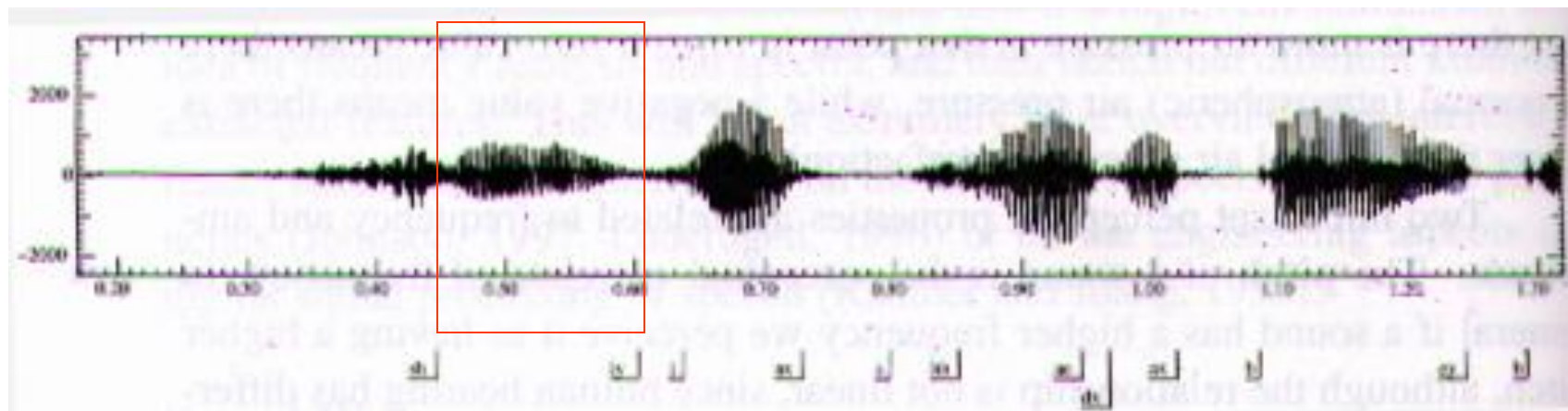
First, different vowels have their formants at characteristic places.

We have seen that the vowel [Q] in the simple waveform had formants at 930Hz, 1860Hz and 3020Hz.

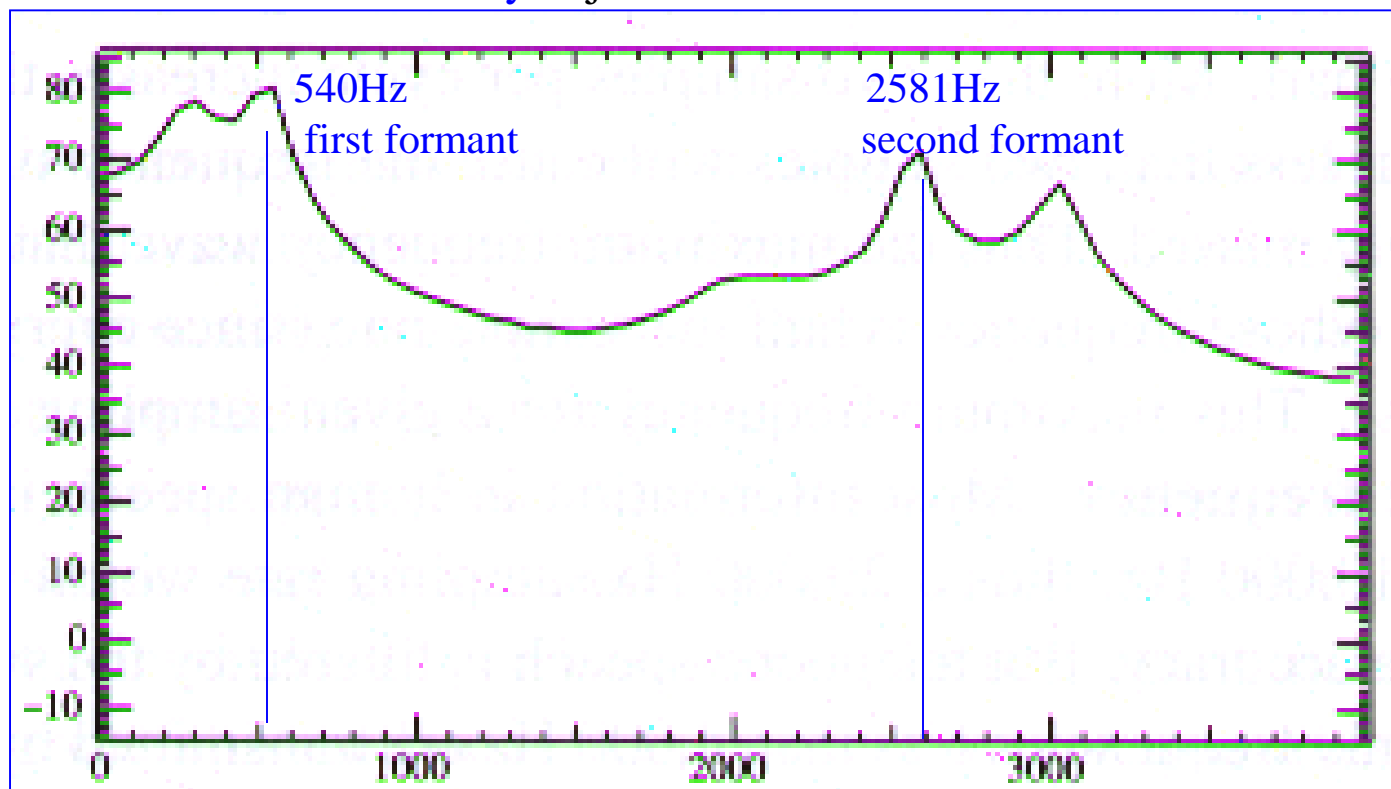
Consider the vowel [iy], at the beginning of the utterance in the first diagram following of Òshèhad a babyÓ.



sh iy j ax s hh ae dx ax b ey b



sh iy j ax s hh ae dx ax b ey b



F1 (formant 1) and F2 (formant 2) play a large role in determining vowel identity,

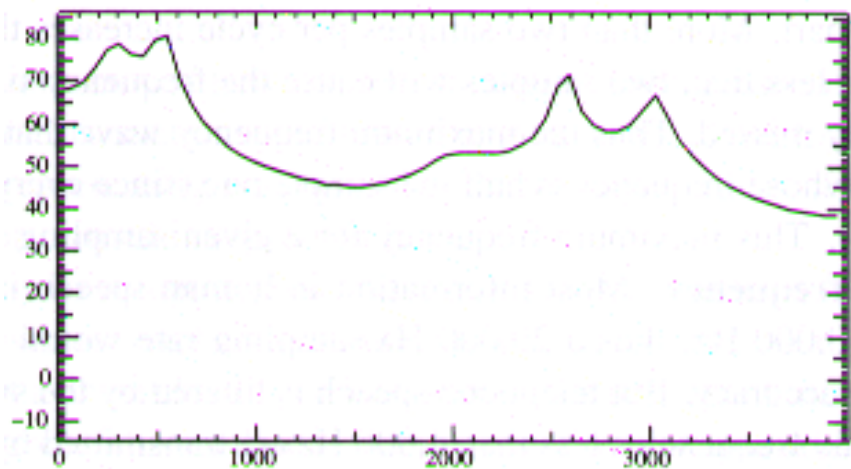
although

the formants still differ from speaker to speaker.

Formants can also be used to identify:

the nasal phones **[n] etc..**

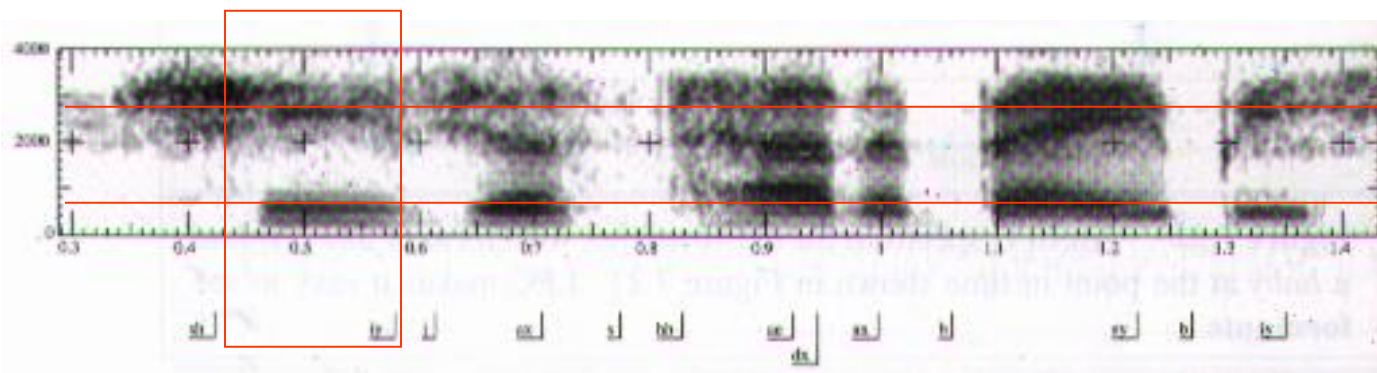
the lateral phone **[l] etc..**



The spectrum for this particular vowel is shown next.

The first formant of [iy] is 540 Hz, much lower than the first formant for [ae], while the second formant at 2581 Hz is much higher than the second formant for [ae].

We can see these as dark bars on the spectrogram diagram

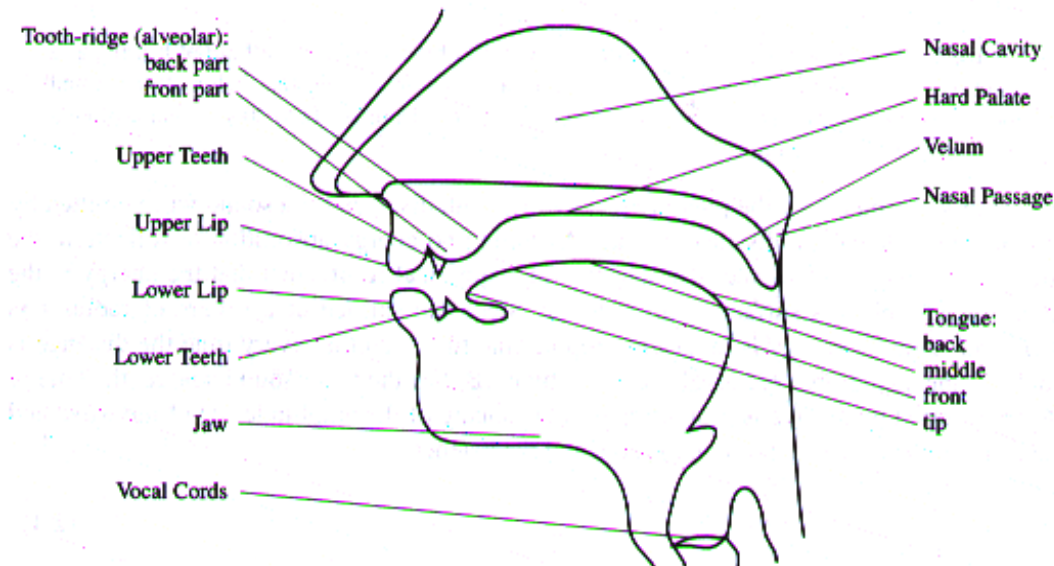


Why do different vowels have different spectra?

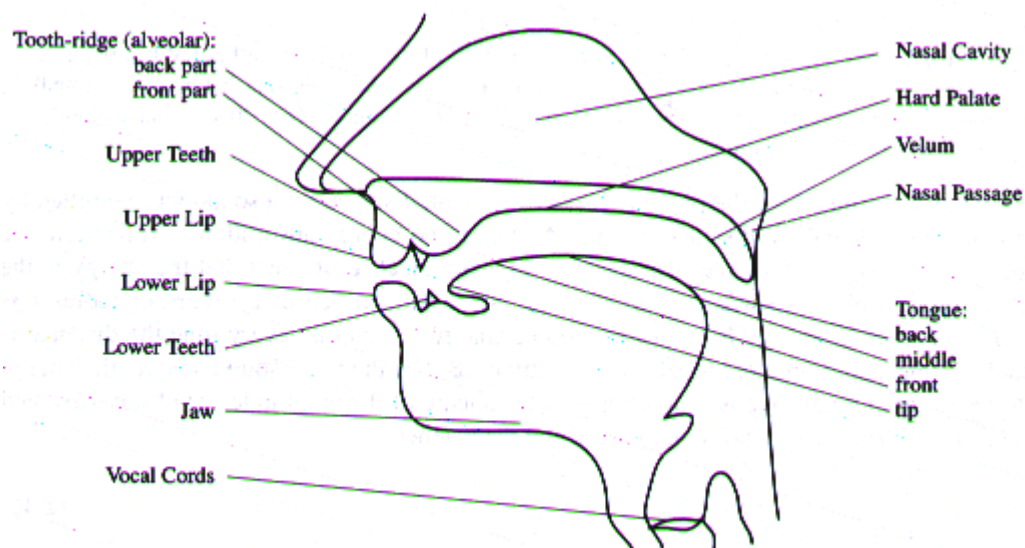
The formants are caused by the resonant cavities of the mouth.

The oral cavity can be thought of as a filter which selectively passes through some of the harmonics of the vocal cord vibrations.

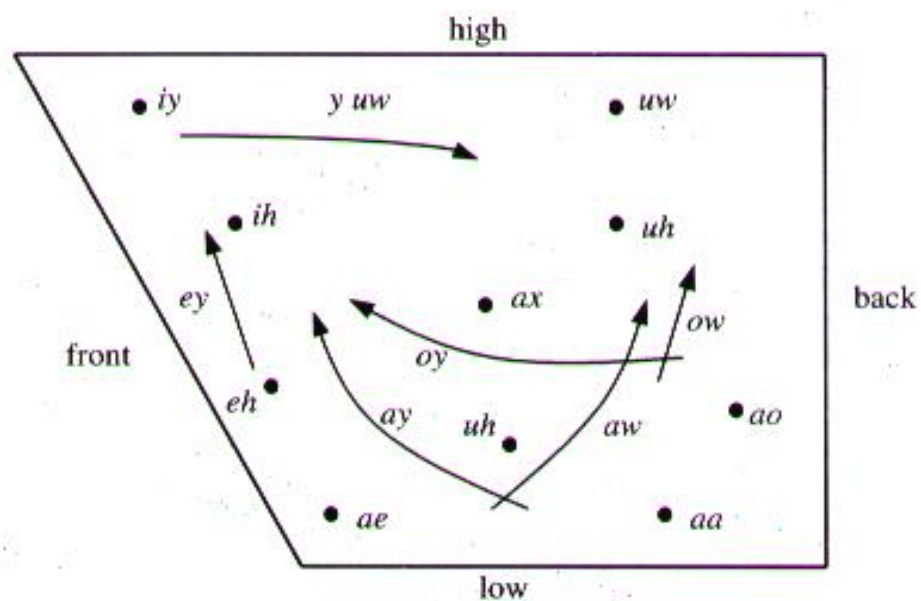
Moving the tongue creates spaces of different size the mouth which selectively amplify waves of the appropriate wavelength, hence amplifying different frequency bands.



A schematic diagram of the human speech production apparatus.



A schematic diagram of the human speech production apparatus.



Relative tongue positions of English vowels

Phonological (abstract) feature decomposition of basic English vowels.

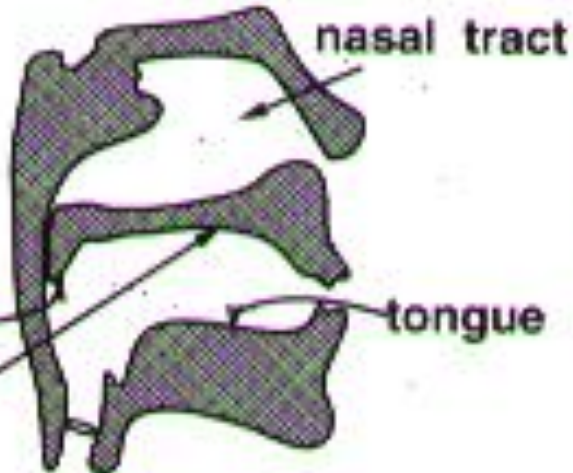
owel	high	low	front	back	round	tense
<i>iy</i>	+	-	+	-	-	+
<i>ih</i>	+	-	+	-	-	-
<i>ae</i>	-	+	+	-	-	+
<i>aa</i>	-	+	-	-	-	+
<i>ah</i>	-	-	-	-	-	+
<i>ao</i>	-	+	-	+	+	+
<i>ax</i>	-	-	-	-	-	-
<i>eh</i>	-	-	+	-	-	-
<i>ow</i>	-	-	-	+	+	+
<i>uh</i>	+	-	-	+	-	-
<i>uw</i>	+	-	-	+	-	+

back



heed [iy]

closed
velum
palate



had [ae]

nasal tract

tongue



who'd [uw]

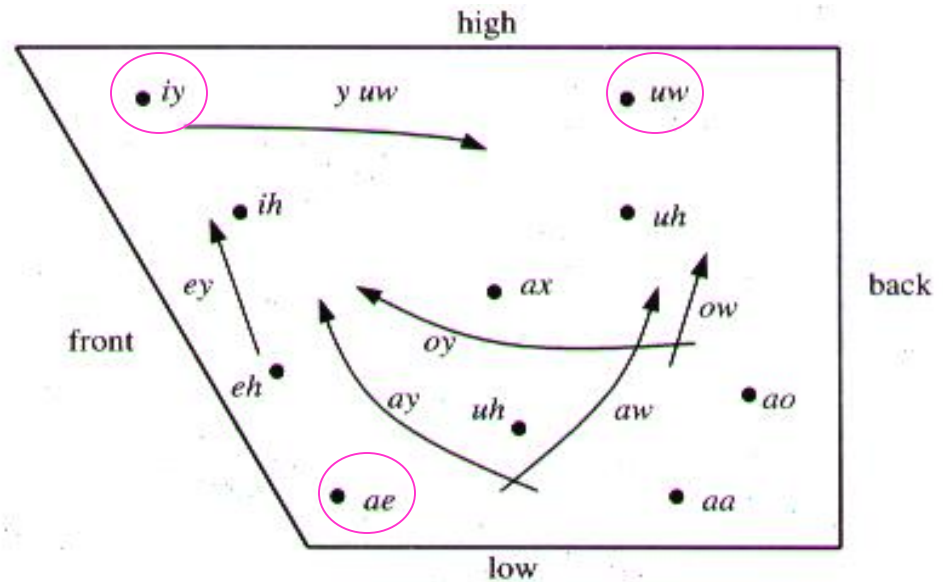
front

Positions of the tongue for three English vowels...

high front [iy]

low front [ae]

high back [uw]



Relative tongue positions of English vowels

IPA Symbol	ARPAbet Symbol	Word	IPA Transcription	ARPAbet Transcription
[i]	[iy]	lily	['lɪli]	[l ih l iy]
[ɪ]	[ih]	lily	['lɪli]	[l ih l iy]
[eɪ]	[ey]	da <u>is</u> y	['deɪzi]	[d ey z i]
[ɛ]	[eh]	po <u>in</u> settia	[pɔɪn'setɪə]	[p oy n s eh dx iy ax]
[æ]	[ae]	a <u>s</u> ter	['æstə]	[ae s t axr]
[ɑ]	[aa]	po <u>p</u> py	['pɑpi]	[p aa p i]
[ɔ]	[ao]	o <u>r</u> chid	['ɔrkɪd]	[ao r k ix d]
[ʊ]	[uh]	w <u>oo</u> druff	['wʊdrʌf]	[w uh d r ah f]
[oʊ]	[ow]	lot <u>u</u> s	['ləʊrəs]	[l ow dx ax s]
[u]	[uw]	tul <u>i</u> p	['tulɪp]	[t uw l ix p]
[ʌ]	[uh]	bu <u>tt</u> er <u>c</u> u <u>p</u>	['bʌtə'kʌp]	[b uh dx axr k uh p]
[ɜ]	[er]	bi <u>r</u> d	['bɜd]	[b er d]
[aɪ]	[ay]	ir <u>i</u> s	['aɪrɪs]	[ay r ix s]
[aʊ]	[aw]	sun <u>f</u> low <u>e</u> r	['sʌnflaʊə]	[s ah n f l aw axr]
[ɔɪ]	[oy]	po <u>in</u> settia	[pɔɪn'setɪə]	[p oy n s eh dx iy ax]
[ju]	[y uw]	fe <u>v</u> er <u>f</u> ew	['fɪvə'fju]	[f iy v axr f y u]
[ə]	[ax]	wo <u>o</u> dr <u>u</u> ff	['wʊdrəf]	[w uh d r ax f]
[ɪ]	[ix]	tul <u>i</u> p	['tulɪp]	[t uw l ix p]
[ə]	[axr]	he <u>a</u> th <u>e</u> r	['hɛðə]	[h eh dh axr]
[ʊ]	[ux]	du <u>d</u> e ²	[dʊd]	[d ux d]

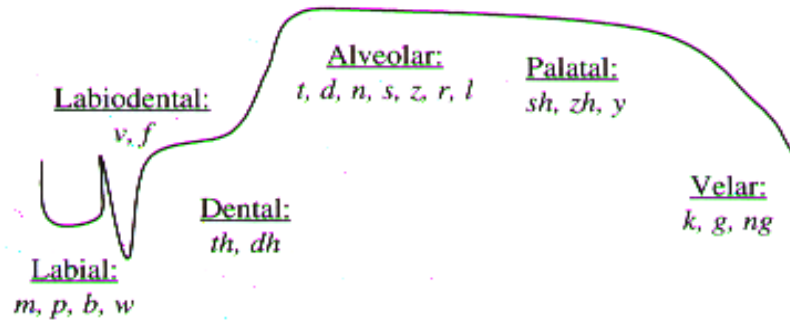
IPA and ARPAbet
symbols for
transcription of
English vowels

IPA Symbol	ARPAbet Symbol	Word	IPA Transcription	ARPAbet Transcription
[p]	[p]	<u>p</u> arsley	['parsli]	[p aa r s l iy]
[t]	[t]	<u>t</u> arragon	['tærəgən]	[t ae r ax g aa n]
[k]	[k]	<u>c</u> atnip	['kætnip]	[k ae t n ix p]
[b]	[b]	<u>b</u> ay	[beɪ]	[b ey]
[d]	[d]	<u>d</u> ill	[dɪl]	[d ih l]
[g]	[g]	<u>g</u> arlic	['gɑrlɪk]	[g aa r l ix k]
[m]	[m]	<u>m</u> int	[mɪnt]	[m ih n t]
[n]	[n]	<u>n</u> utmeg	['nʌtmeg]	[n ah t m eh g]
[ŋ]	[ng]	<u>g</u> inseng	['dʒɪnsɪŋ]	[jh ih n s ix ng]
[f]	[f]	<u>f</u> ennel	['fenl]	[f eh n el]
[v]	[v]	<u>c</u> love	[klov]	[k l ow v]
[θ]	[th]	<u>t</u> histle	['θɪsl]	[th ih s el]
[ð]	[dh]	<u>h</u> eather	['hæðə]	[h eh dh axr]
[s]	[s]	<u>s</u> age	[seɪdʒ]	[s ey jh]
[z]	[z]	<u>h</u> azelnut	['heɪz nʌt]	[h ey z el n ah t]
[ʃ]	[sh]	<u>s</u> quash	[skwɒʃ]	[s k w a sh]
[ʒ]	[zh]	<u>a</u> mbrosia	[æm'brʊʒə]	[ae m b r ow zh ax]
[tʃ]	[ch]	<u>c</u> hicory	['tʃɪkəɪ]	[ch ih k axr iy]
[dʒ]	[jh]	<u>s</u> age	[seɪdʒ]	[s ey jh]
[l]	[l]	<u>l</u> icorice	['lɪkəɪf]	[l ih k axr ix sh]
[w]	[w]	<u>k</u> iwi	['kiwi]	[k iy w iy]
[r]	[r]	<u>p</u> arsley	['parsli]	[p aa r s l iy]
[j]	[y]	<u>y</u> ew	[ju]	[y uw]
[h]	[h]	<u>h</u> orseradish	['hɔrsrædɪʃ]	[h ao r s r ae d ih sh]
[ʔ]	[q]	uh-oh	[ʔʌʔou]	[q ah q ow]
[ɹ]	[dx]	<u>b</u> utter	['bʌtə]	[b ah dx axr]
[ɹ]	[nx]	<u>w</u> intergreen	[wɪntəgrɪn]	[w ih nx axr g r i n]
[l]	[el]	<u>t</u> histle	['θɪsl]	[th ih s el]

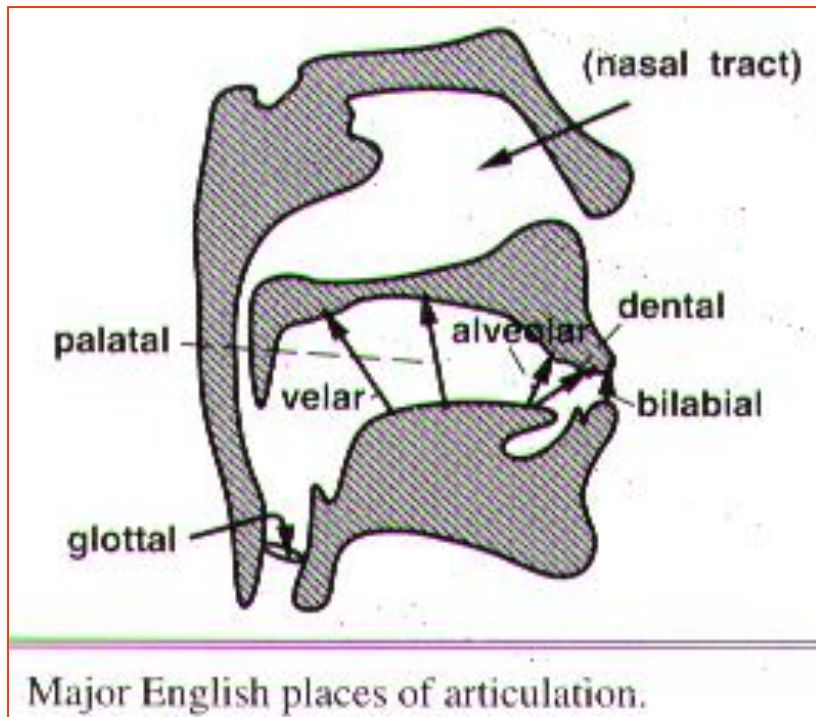
IPA and **ARPAbet** symbols for transcription of English consonants

Manner of articulation of English consonants.

Consonant Labels	Consonant Examples	Voiced?	Manner
<i>b</i>	big, able, tab	+	plosive
<i>p</i>	put, open, tap	-	plosive
<i>d</i>	dig, idea, wad	+	plosive
<i>t</i>	talk, sat	-	plosive
<i>g</i>	gut, angle, tag	+	plosive
<i>k</i>	cut, oaken, take	-	plosive
<i>v</i>	vat, over, have	+	fricative
<i>f</i>	fork, after, if	-	fricative
<i>z</i>	zap, lazy, haze	+	fricative
<i>s</i>	sit, cast, toss	-	fricative
<i>dh</i>	then, father, scythe	+	fricative
<i>th</i>	thin, nothing, truth	-	fricative
<i>zh</i>	genre, azure, beige	+	fricative
<i>sh</i>	she, cushion, wash	-	fricative
<i>jh</i>	joy, agile, edge	+	affricate
<i>ch</i>	chin, archer, march	-	affricate
<i>l</i>	lid, elbow, sail	+	lateral
<i>r</i>	red, part, far	+	retroflex
<i>y</i>	yacht, onion, yard	+	glide
<i>w</i>	with, away	+	glide
<i>hh</i>	help, ahead, hotel	+	fricative
<i>m</i>	mat, amid, aim	+	nasal
<i>n</i>	no, end, pan	+	nasal
<i>ng</i>	sing, anger, drink	+	nasal



The major places of consonant articulation with respect to the human mouth.



Feature Extraction

We can now summarise the process of extraction of spectral features, beginning with the sound wave itself and ending with a feature vector.

An input sound-wave is first digitised.

This process of analogue-to-digital conversion has two steps:

- Sampling and
- Quantisation.

A signal is sampled by measuring its amplitude at a particular time.

The sampling rate is the number of samples taken per second.

Common sampling rates are 8,000Hz and 16,000Hz.

In order to accurately measure a wave, it is necessary to have at least two samples in each cycle. One measuring the positive part of the wave and the other measuring the negative part.

More than two samples per cycle increases the amplitude accuracy, but less than two samples will cause the frequency of the wave to be completely missed.

Therefore, the maximum frequency wave that can be measured is one whose frequency is half that of the sample rate (since every cycle needs two samples).

The maximum frequency for a given sampling rate is called the **Nyquist** frequency.

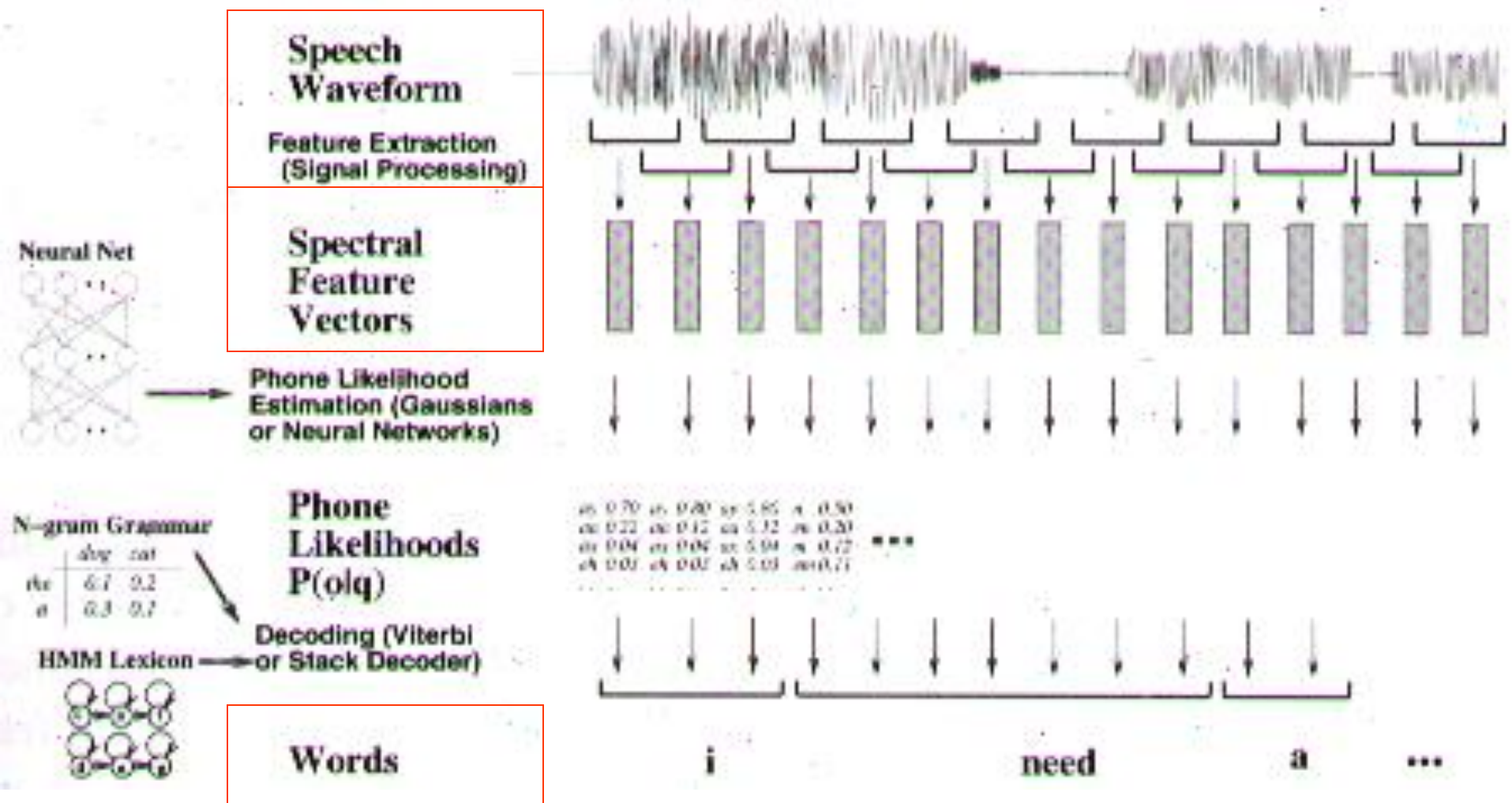
Most information in human speech is in frequencies below 10,000 Hz, therefore a 20,000 Hz sampling rate would be necessary for complete accuracy.

A sampling rate of, for example, 8,000 Hz, will require 8000 amplitude measurements for each second of speech.

It is very important to store amplitude measurements efficiently.

They are usually stored as integers, either 8-bit with values from -128 to 127 or 16-bit with values from -32768 to 32767 .

- This process of representing a real-valued number as an integer is **quantisation** because there is a minimum granularity (the quantum size) and all values which are closer together than this quantum size are represented identically.
- Once a waveform has been digitised, it is converted to some set of spectral features.



TO DO THIS WEEK

1. Read chapter 7 of the Jurafsky and Martin textbook