# EE679: Speech Processing

## A preview
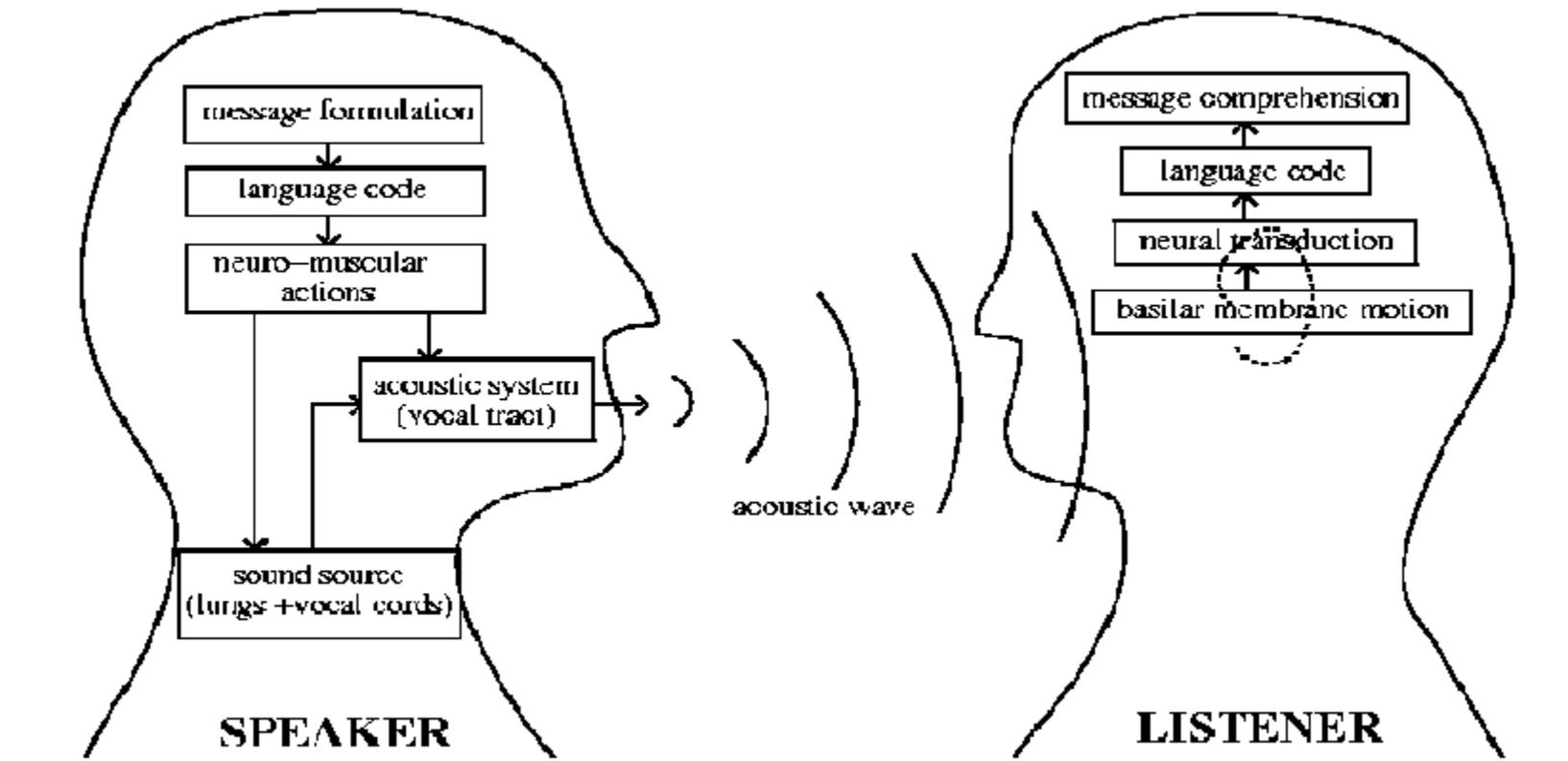
**Dept of Electrical Engineering**

**I.I.T. Bombay**

# Outline

- Speech production (physiology)

- Classification of sounds: articulatory, acoustic

- Speech analysis (signal processing methods for information extraction)

- Hearing, and speech perception

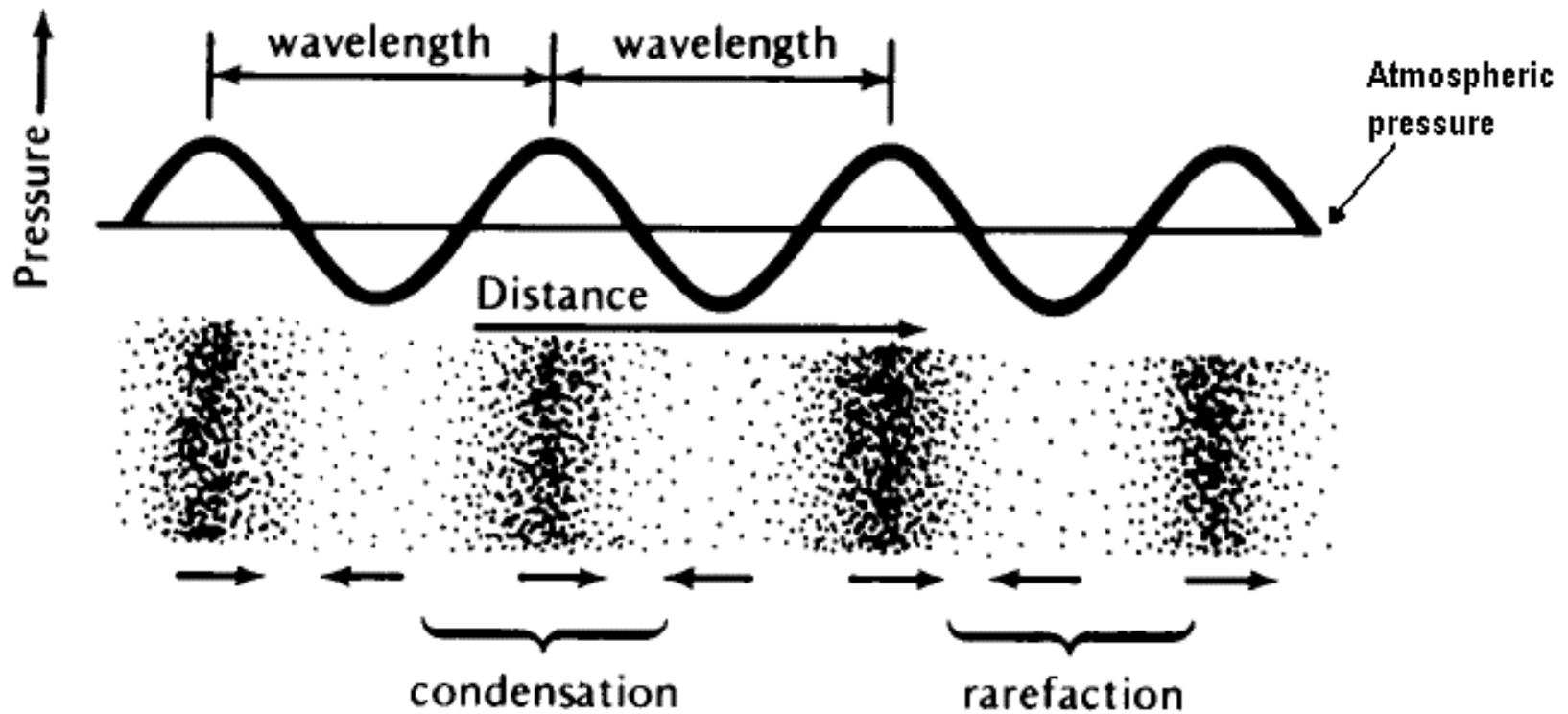- Speech technology (speech compression, ASR,TTS)

- Audio/music technology

# Speech communication
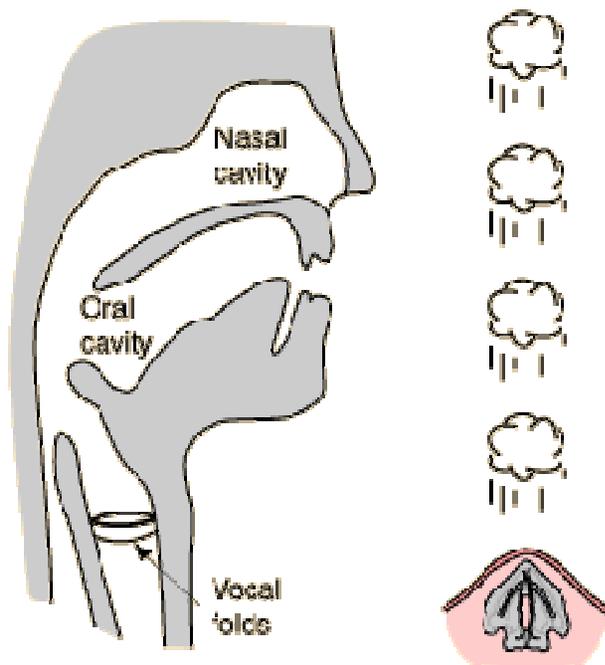
# Acoustic waves

*Speed = wavelength x frequency*

# Information in speech

- Linguistic (phone->word->sentence->message)
- Paralinguistic:

   --speaker-based (pronunciation, age, sex,etc.),

   --expressive (emotions, mood)

*The speech signal is characterised by an enormous range of perceptually contrasting sounds!*
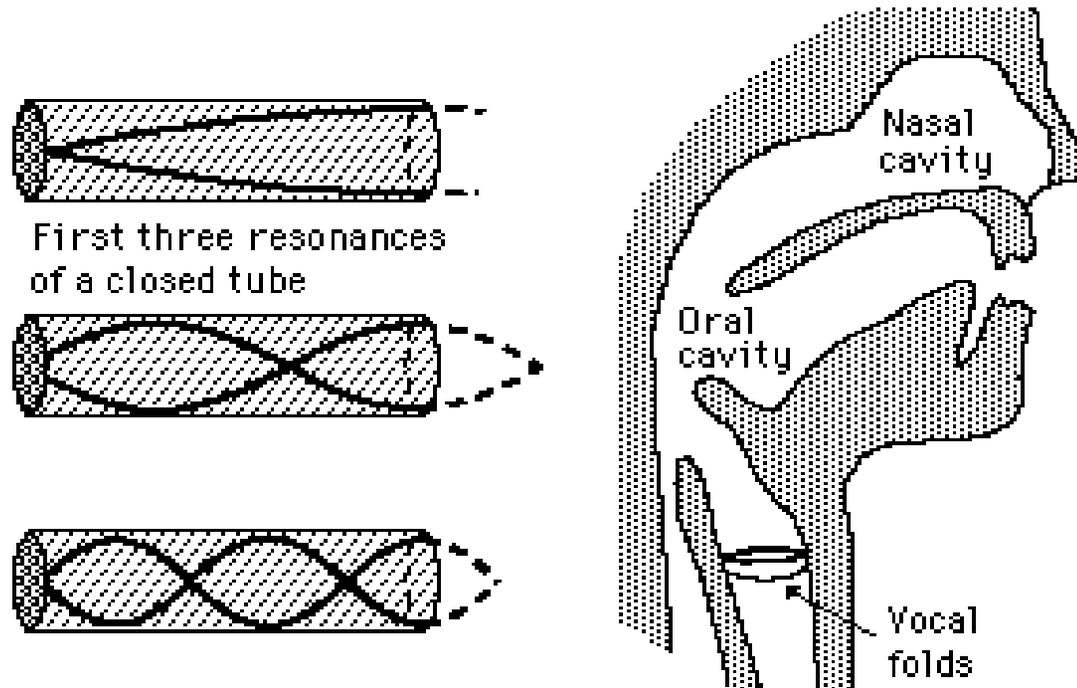
# Generating speech*



Respiration->**phonation**
->articulation

Vibrating vocal cords create puffs of air giving rise to *air pressure variations* which reach our ears.

*HyperPhysics, Sound and Hearing, Georgia State University*
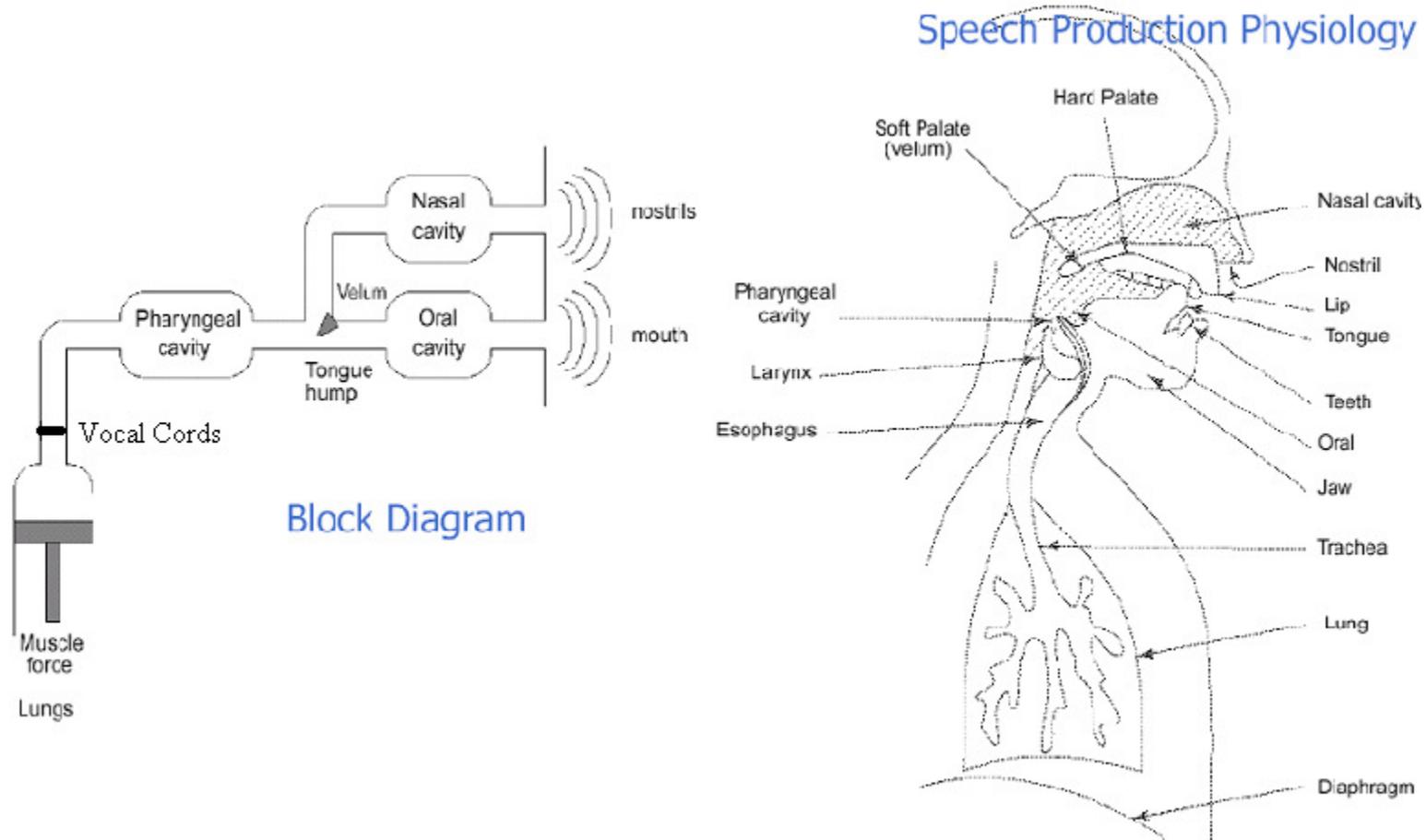
6

# Vocal tract: Acoustic resonances*

First three resonances of a closed tube

Nasal cavity

Oral cavity

Vocal folds

$$f_1 = \frac{c}{4L} \quad ; \quad f_2 = \frac{3c}{4L} ; \quad f_3 = \frac{5c}{4L} ; \ldots\ldots$$

**Department of Electrical Engineering , IIT Bombay**

# Speech production (Childers, Speech Overview, 1993)



Block Diagram

Speech Production Physiology

# **Articulation**: producing the *<u>various sounds</u>* of speech*

Velum

Teeth

Lips

Tongue   Jaw

Vocal cords

Nasal sound output

Nasal cavity

Velum

Pharyngeal cavity

**Vocal cavity**

Oral Cavity

Articulators

Oral sound output

Trachea connection to lungs

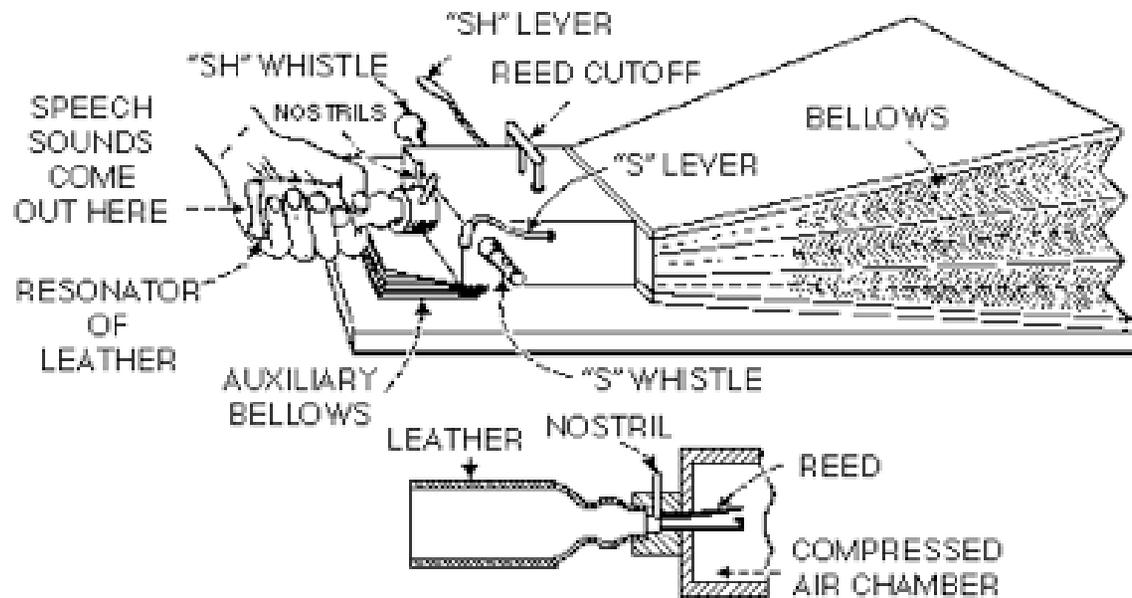Moving muscles which alter the resonant cavities

Dynamic cavity

Static cavity

*Securivox tutorial

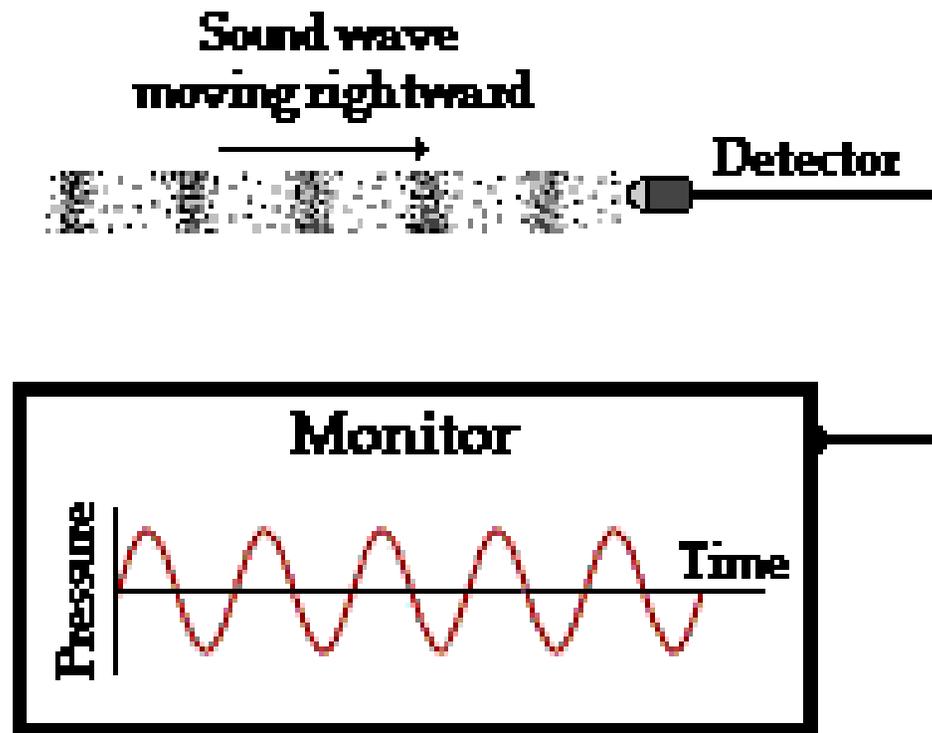# Von Kempelen's talking machine
# 1791

# 1875

- Alexander Bell invents the method of, and apparatus for, "transmitting vocal or other sounds telegraphically ... by causing electrical undulations, similar in form to the vibrations of the air accompanying the said vocal or other sound".

=> Major impetus to modern speech processing.

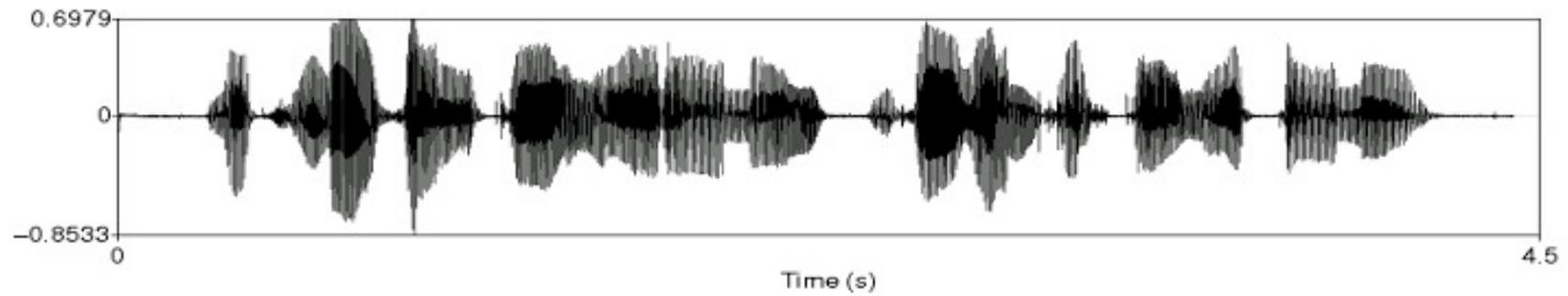- 1930s: Electrical synthesis of speech by Dudley's vocoder

# Sound -> electrical form*

# Speech "waveform"

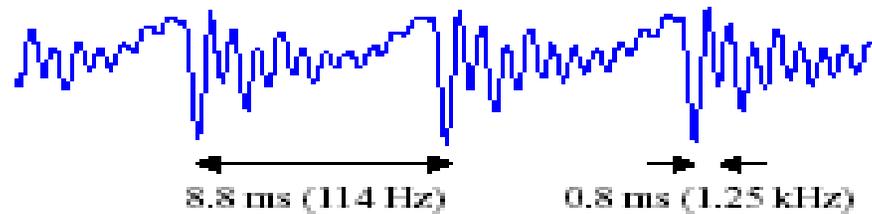# Speech Waveforms from "my speech"

(a) start of "y" vowel

8.8 ms (114 Hz)    0.8 ms (1.25 kHz)

(b) "ee" vowel

8.8 ms (114 Hz)    4.3 ms (233 Hz)

(c) "s" consonant

**Department of Electrical Engineering , IIT Bombay**

14

# Components of sound

A sound is usually comprised of *several* frequency components.

Depending on the relationships of the frequency components, the sound can elicit a sensation of pitch.

# Speech production

- **<u>Vocal cords</u>** (larynx) modulate the airflow from the lungs by rapid opening-closing; the *rate of vibration* is determined by their mass and tension.

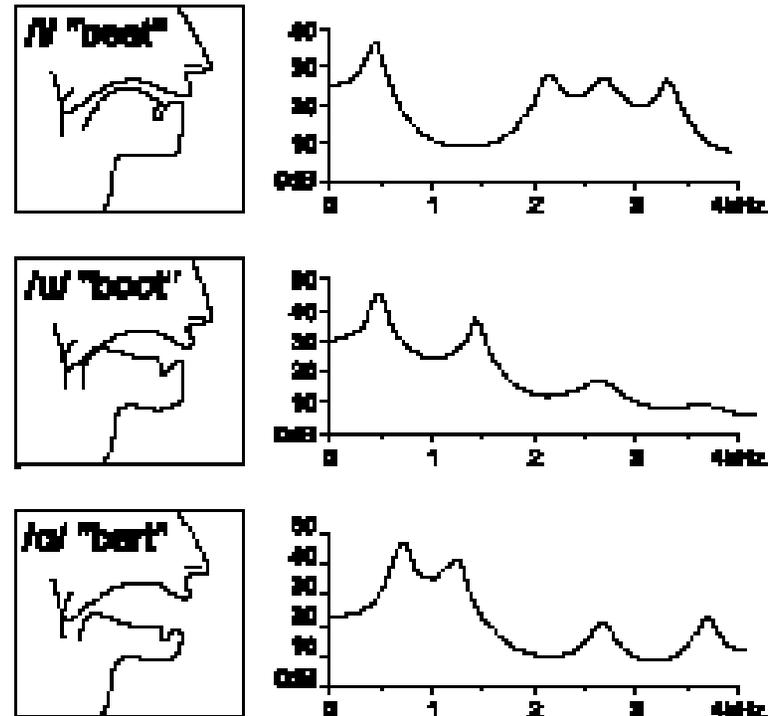  Pitch frequency ranges:

  male: 80-160 Hz; female:160-320 Hz;

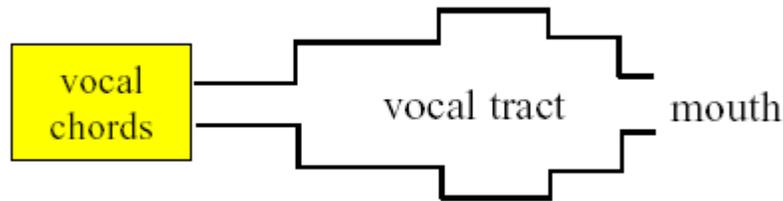  singers: over 2 octaves.


- **<u>Vocal tract</u>** shapes the vocal cord vibrations into the intricate sounds of speech via *changes in shape* to produce various *acoustic resonances*.

# Vocal tract "filter"*

• The sound spectrum is modified by the shape of the vocal tract.
• The resonant frequencies of the vocal tract cause peaks in the spectrum called *formants*.



*Childers, Speech Overview

# Most important aspects of speech…

- The intelligence in speech is encoded in the *power spectrum* of the acoustic pressure wave.

- Different articulatory configurations result in signals with different spectra, esp. different resonance frequencies called formants, which are perceived as different sounds.

- The different spectra make up the finite alphabet of symbols (linguistic code) governed by a hierarchy of linguistic rules.

# Basic sounds of speech: Phones

- The speech signal can be divided into sound segments
  with *fixed articulation and acoustics over short intervals.*

  *i.e. articulatory configuration <=> acoustic properties*

Smallest meaningful sound unit: "***phone***"

(i.e. set of distinctive sounds of a language)

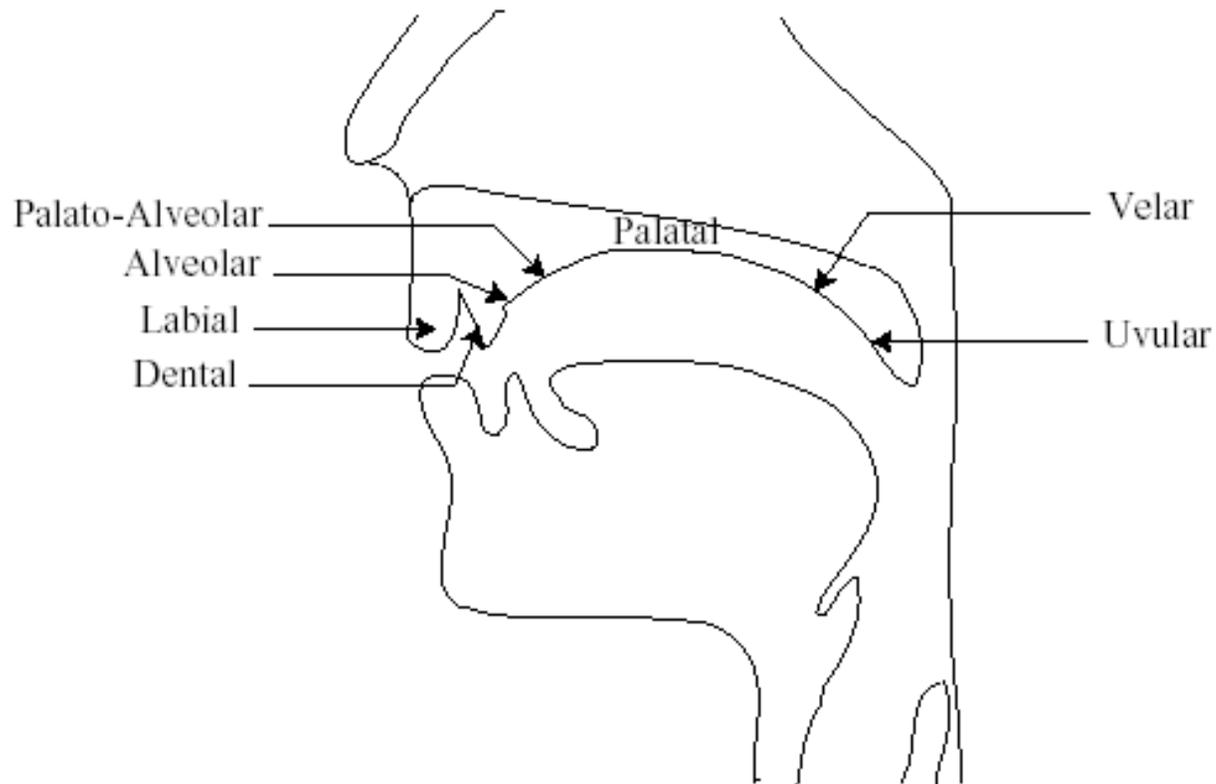In Indian written scripts, one symbol represents one phone.

# Classification of speech sounds

## Vowels and Consonants

- <u>Vowels</u>: steady sounds specified by position of the articulators (typically, tongue)

- <u>Consonants</u>: are (dynamic) sounds classified by **place** and **manner** of articulation

# Place of articulation
## (*constriction of vocal tract*)

| MoA \ PoA | | Velar | Palatal/ Palato Alvelor | Retroflex | Dentals | Bilabials |
|---|---|---|---|---|---|---|
| Stop and Affricates | UvUa | क | च | ट | त | प |
| | UvAs | ख | छ | ठ | थ | फ |
| | VoUa | ग | ज | ड | द | ब |
| | VoAs | घ | झ | ढ | ध | भ |
| Nasal | | ङ | ञ | ण | न | म |

UvUa: Unvoiced Unaspirated     MoA: Manner of Articulation
UvAs: Unvoiced Aspirated     PoA:  Place of Articulation
VoUa: Voiced Unaspirated
VoAs: Voiced Aspirated

*Table 1: The table shows some consonants used in Hindi language and their classification depending on their phonetic properties.*
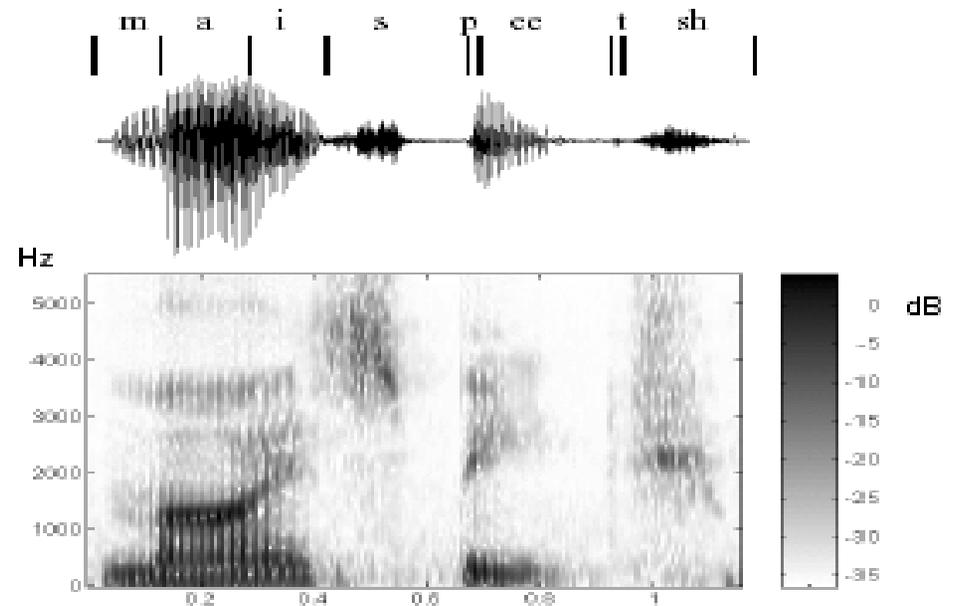
# "*Decoding*" the speech signal: visible speech

Dark areas of *spectrogram* show high intensity

– Voiced segments are much louder than unvoiced

– Horizontal dark bands are the formant peaks

– "s" has high frequency content
– Vertical bands are individual larynx closures
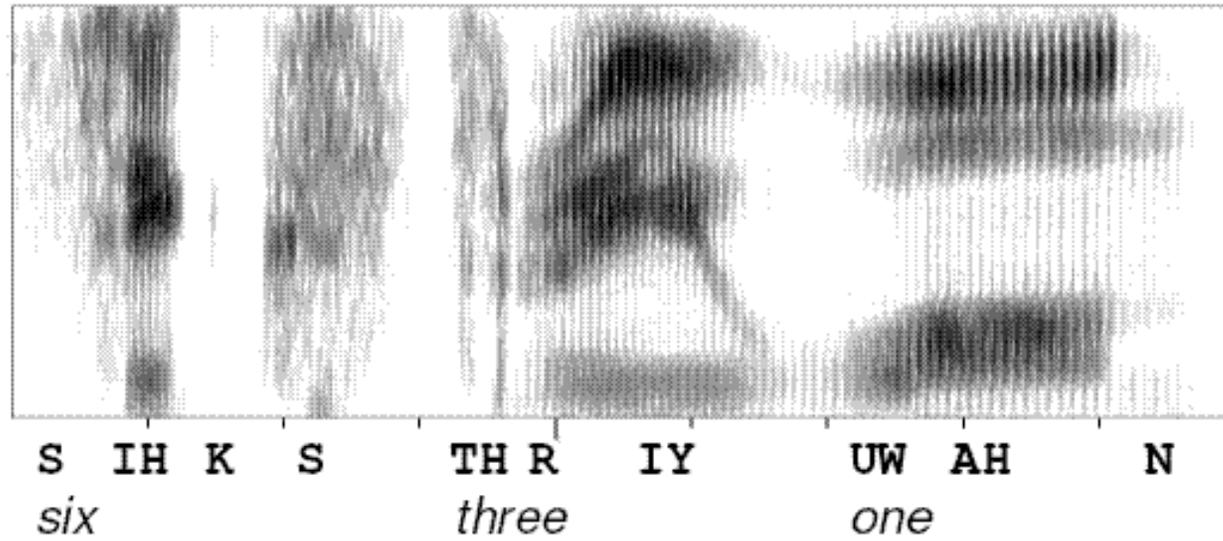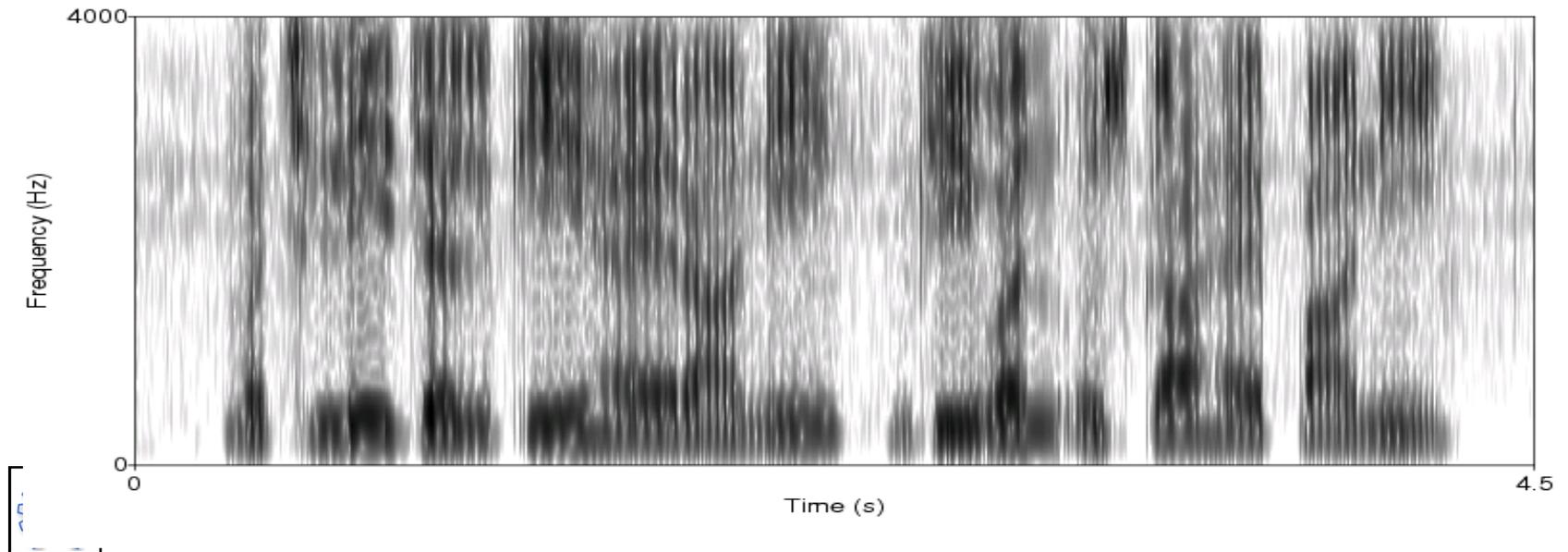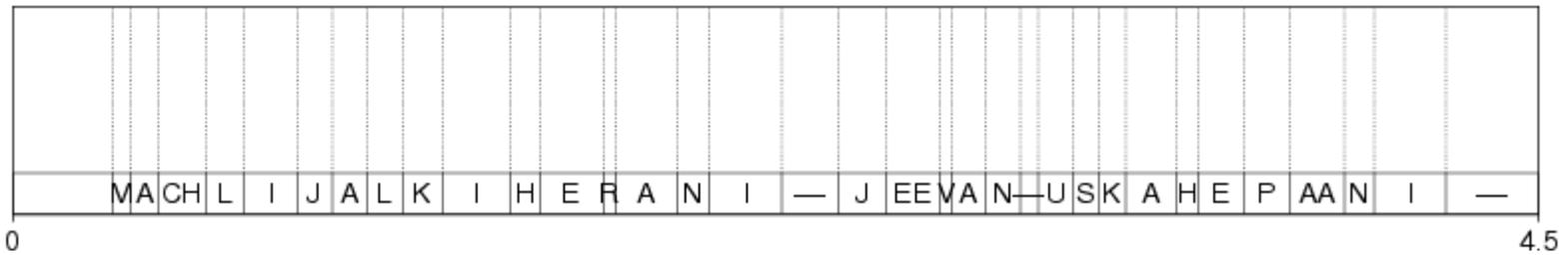– The "y" of "my" is a diphthong: two successive vowels

"my speech"

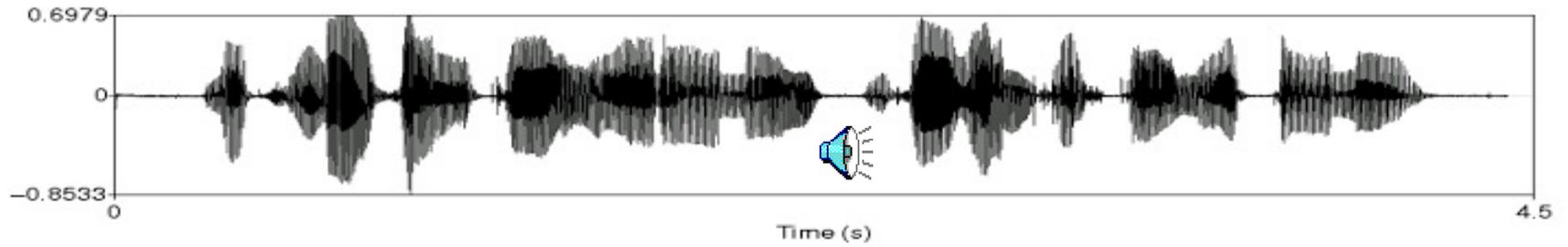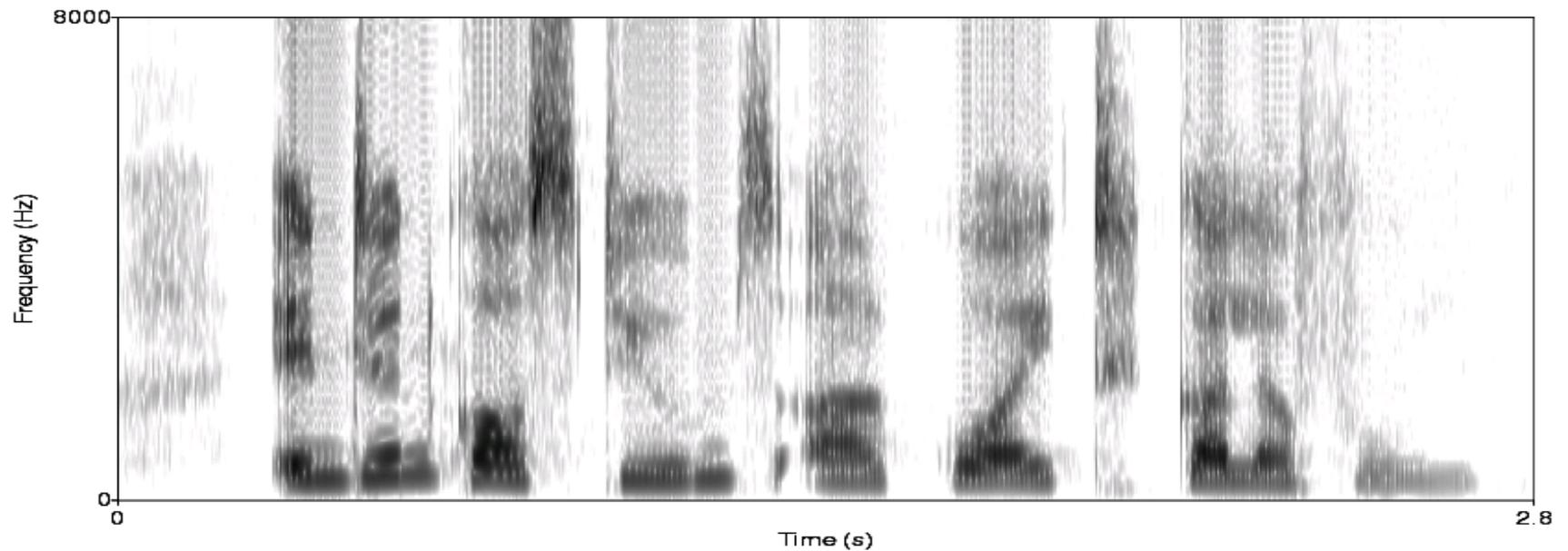**Department of Electrical Engineering , IIT Bombay**

Figure 1: Speech spectral characteristics

# Machli jal ki hai raani   jeevan uska he paani

# Indian costumes are quite colourful



| IY | N | D | IA | N | K | O | S | T | U | M | S | – | Ar | K | A | I | T | ¨ | K | A | L | Ar | F | U | L |

# Speech perception

Distinct stages of physiological processing
in the auditory system:

Peripheral auditory system   (Ears)        ← analysis

Auditory nervous system    (Brain)        ←synthesis

# Audible sound

# Sound and Sensation

A sound of given frequency components and sound pressure levels leads to perceived sensations that can be distinguished in terms of:

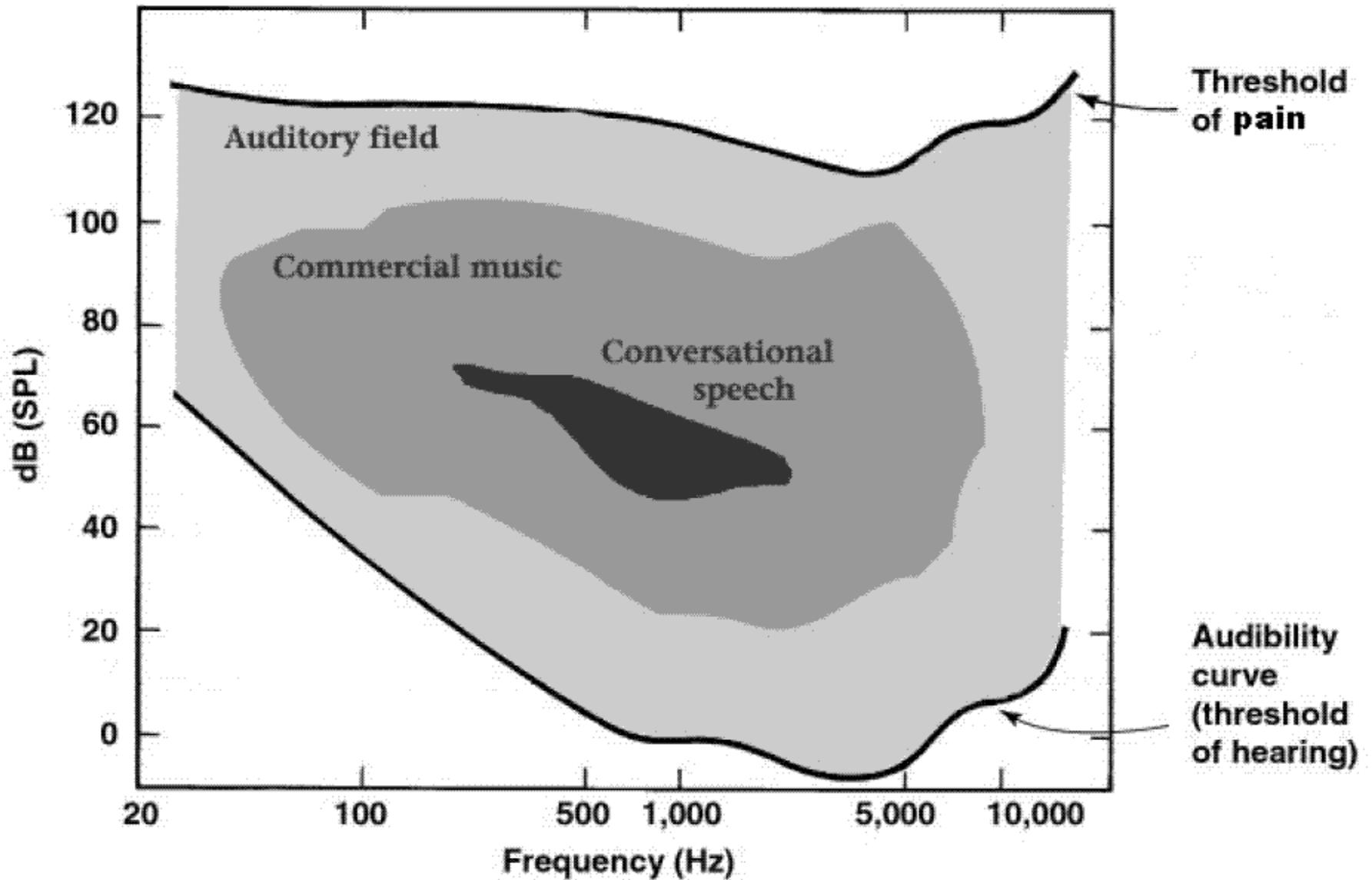- *loudness*     *<-- intensity*
- *pitch*     *<-- fundamental frequency*
- *timbre ("quality" or "colour")*
    - *<--ther spectro-temporal properties*

# Our auditory apparatus



Cochlea:
Ear's microphone

*HyperPhysics, Sound and Hearing, Georgia State University*
*(http://hyperphysics.phy-*
*astr.gsu.edu/hbase/sound/soucon.html#soucon)*

# Basilar Membrane
## Location-dependent frequency "resonance"

•Thickness and tension vary along its length

•Traveling wave has maximum vibration amplitude at a location depending on its frequency



WIDTH AT APEX 0.5 mm
WIDTH AT BASE 0.04 mm
AVERAGE WIDTH 0.21 mm BASAL TURN
0.34 mm MIDDLE TURN
0.36 mm APICAL TURN
LENGTH 32 mm

# Basilar Membrane

*Frequency-to-place transformation*
*(Fourier analysis)*



High Frequencies

Low Frequencies

nerve signals

# Applications

- Automatic speech recognition/ understanding
- Text-to-speech synthesis
- Speaker verification (biometric)
- Digital storage/transmission of speech
- Aids to the handicapped
- Enhancement of quality

**Transmission/storage**

**Waveform coding:**

**distortion** vs **bit rate**



*What distortion is "acceptable" depends on the application and on human perception.*

# Digital audio bit rates: Waveform coding

| Format | Sample Rate (kHz) | Bits/sample |
|---|---|---|
| **Telephony** | **8** | **12** (=> 96 kbps) |
| **Wideband audio** | **16** | **16** |
| **Hi-fidelity audio** | **44.1** | **16** |

# Source-filter model parameters



Pitch and vocal tract shape vary slowly in time

# Frame-based coding of speech



Feature vector (short-time spectrum) extraction from speech

# Automatic speech recognition

- To extract the linguistic code (a *structured sequence of discrete symbols*) from an *analysis* **of the acoustic speech signal**.

- That is, <u>continuous</u>, <u>noisy</u> measurements of a non-stationary function of time only are available.

# Automatic speech recognition

- Feature calculation (to a more distinctive domain)

- Pattern classification with respect to previously trained models of phones/words

- Improved transcription based on language model

# ASR: block diagram*



Signal

Feature Extraction

Training — Testing

Acoustic Model

Matching (acoustic domain)

Symbol sequence

Language Model

Matching (symbolic domain)

Sentence Hypothesis

*K.Samudravijaya, A Tutorial on Speech and Speaker Recognition*

# ASR: Challenges

- Inter- and intra-speaker variations
- Effects of coarticulation in continuous speech
- Background noise and variable channels

# Categories of speech recognition tasks

Human to machine:

- Database query/ information retrieval
- Dictation

Human to human:

- Broadcast news
- Lectures
- Voice mail
- Meeting
- Telephone conversation

# Speaker recognition
## (*voice-based biometric*)

- The voice signal is considered relatively easy to acquire/collect.

- Speech enables an (indirect) measurement of physiological features (i.e. characteristics of the speaker's *voice production system*).

- Applications:

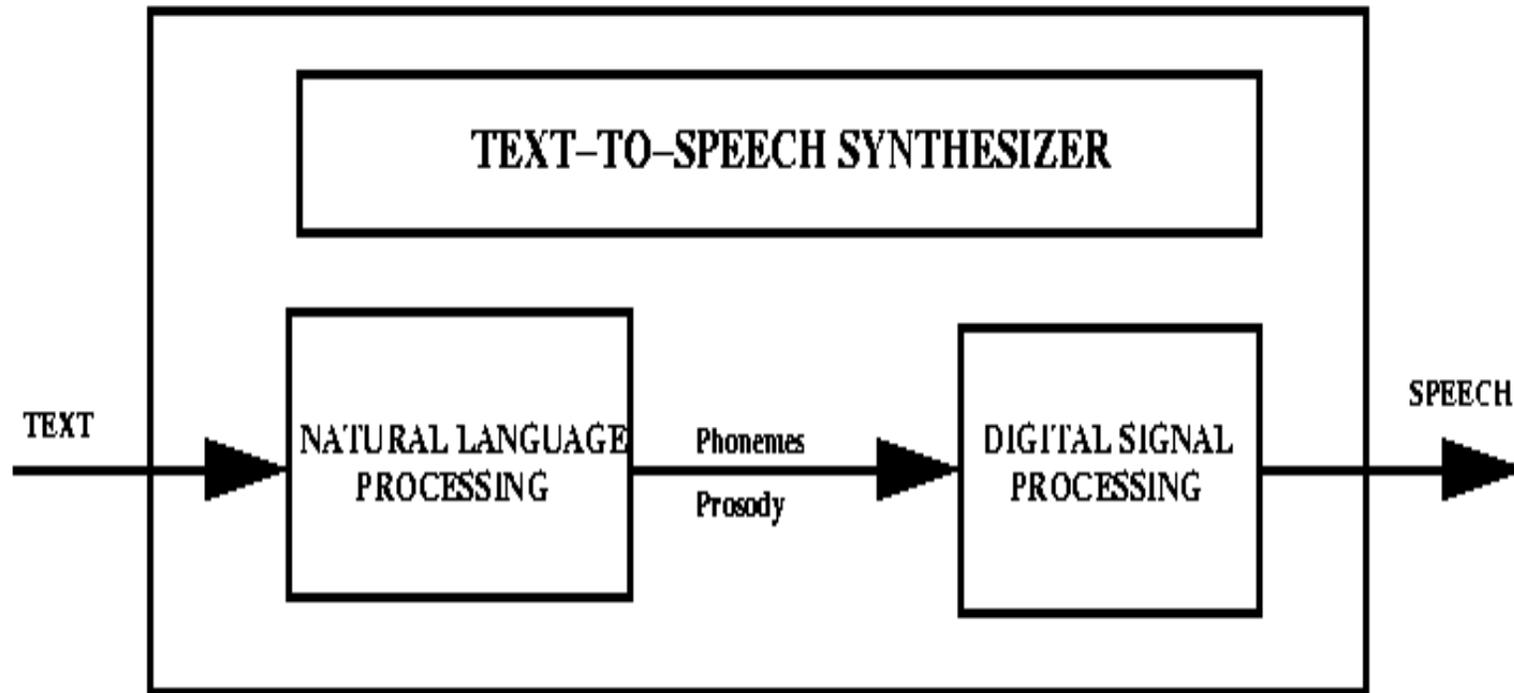  Commercial (access control, segmentation)

  Military, Forensic

# Speech Synthesis

**What**: To convert a text string into a speech waveform

**Why**: For technology to communicate when a display would
   be inconvenient.

# Basic TTS System



Prosody => A phone is *long*/*short*, *loud*/*soft*, *high*/*low*-pitched

# Outline

- Speech production (physiology)

- Classification of sounds: articulatory, acoustic

- Speech analysis (signal processing methods for information extraction)

- Hearing, and speech perception

- Speech technology (speech compression, ASR,TTS)

- Audio/music technology

# Text / References

- Douglas O'Shaughnessy, Speech Communications: Human and Machine, Universities Press (India) Ltd., 2001

- Rabiner and Schafer, Digital Processing of Speech Signals

- IITB Moodle for all course-related hand-outs

# Recognition: "Vowel triangle"



The vowel triangle. plot of F1 versus F2 for vowels in American English

# Speaker variability: due to differences in vocal physiology