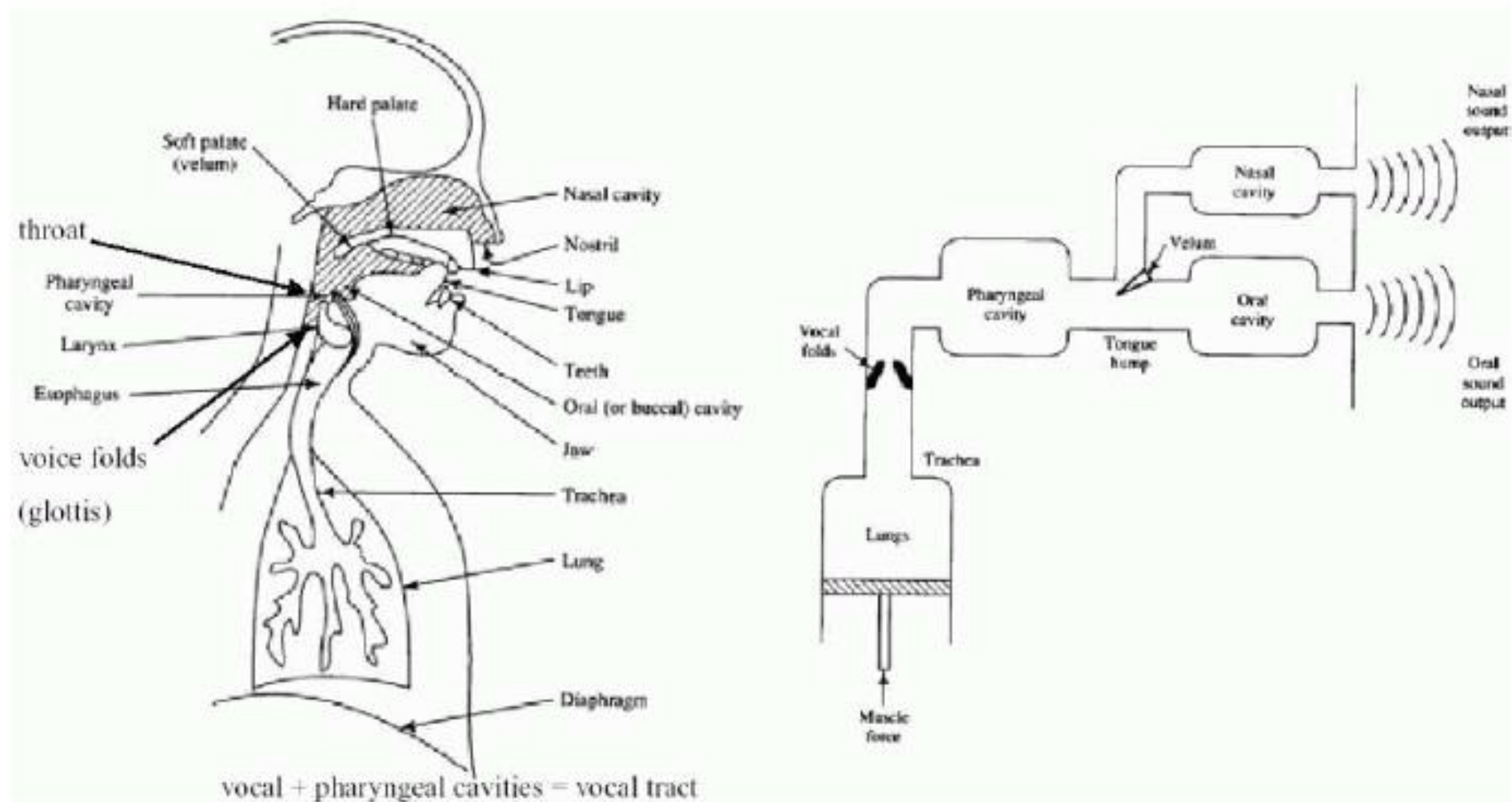# Speech Modeling, Analysis, Synthesis & Compression

Tutorial for SSY150:  Lab 1
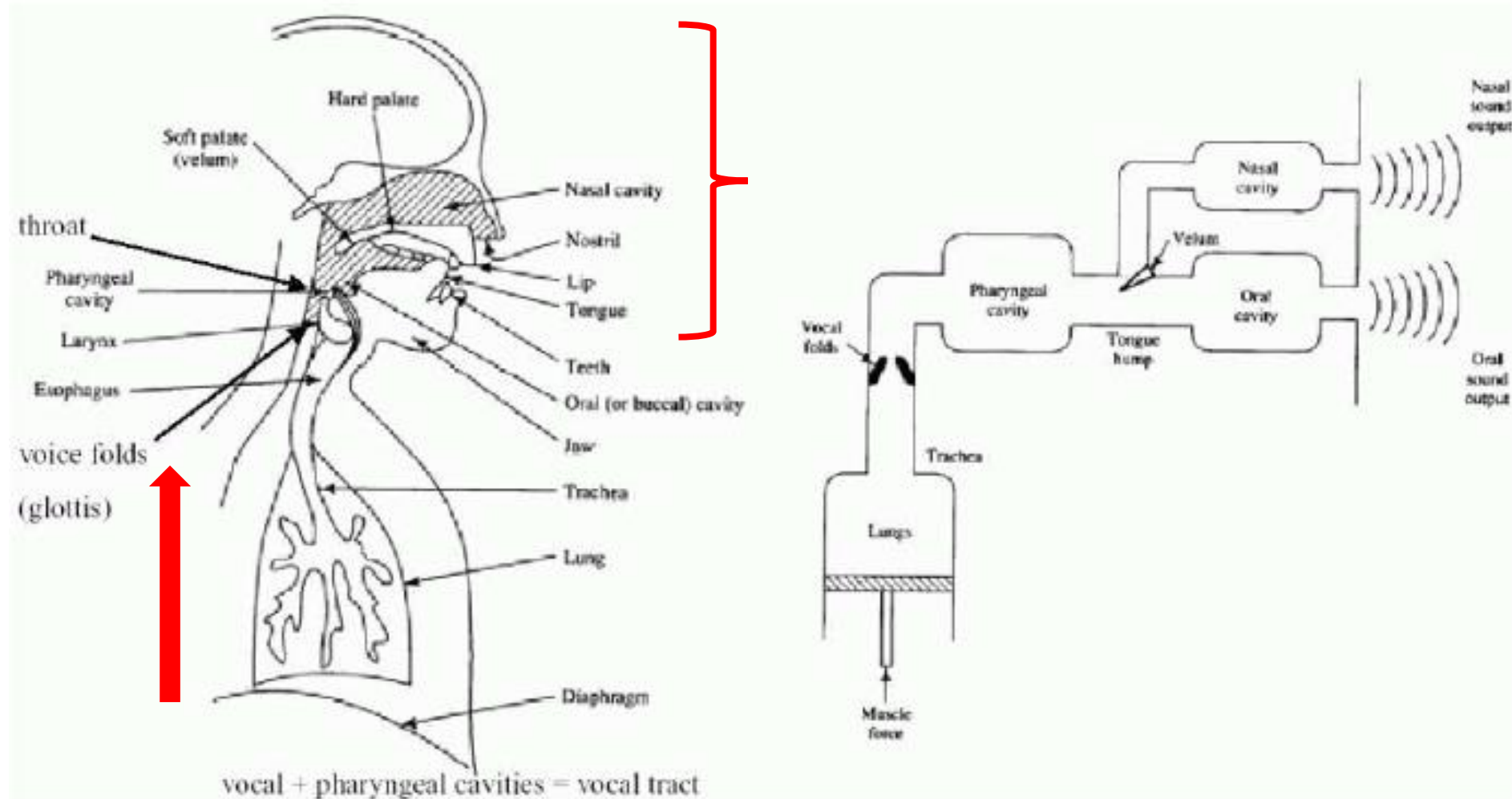
Muhaddisa Barat Ali,
Dept. Of Electrical Engineering, Chalmers University of Technology, Sweden

# Mechanism for Human Sound Production



vocal + pharyngeal cavities = vocal tract

# Mechanism for Human Sound Production



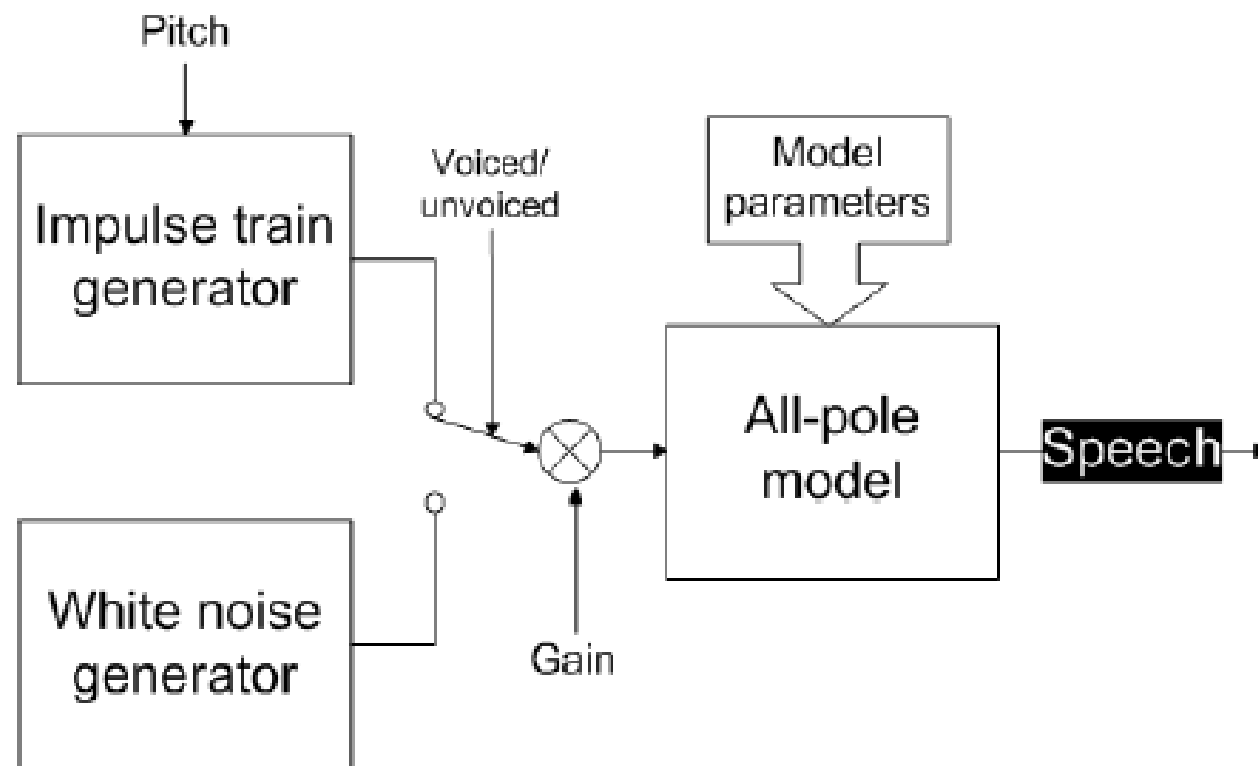Vocal tract
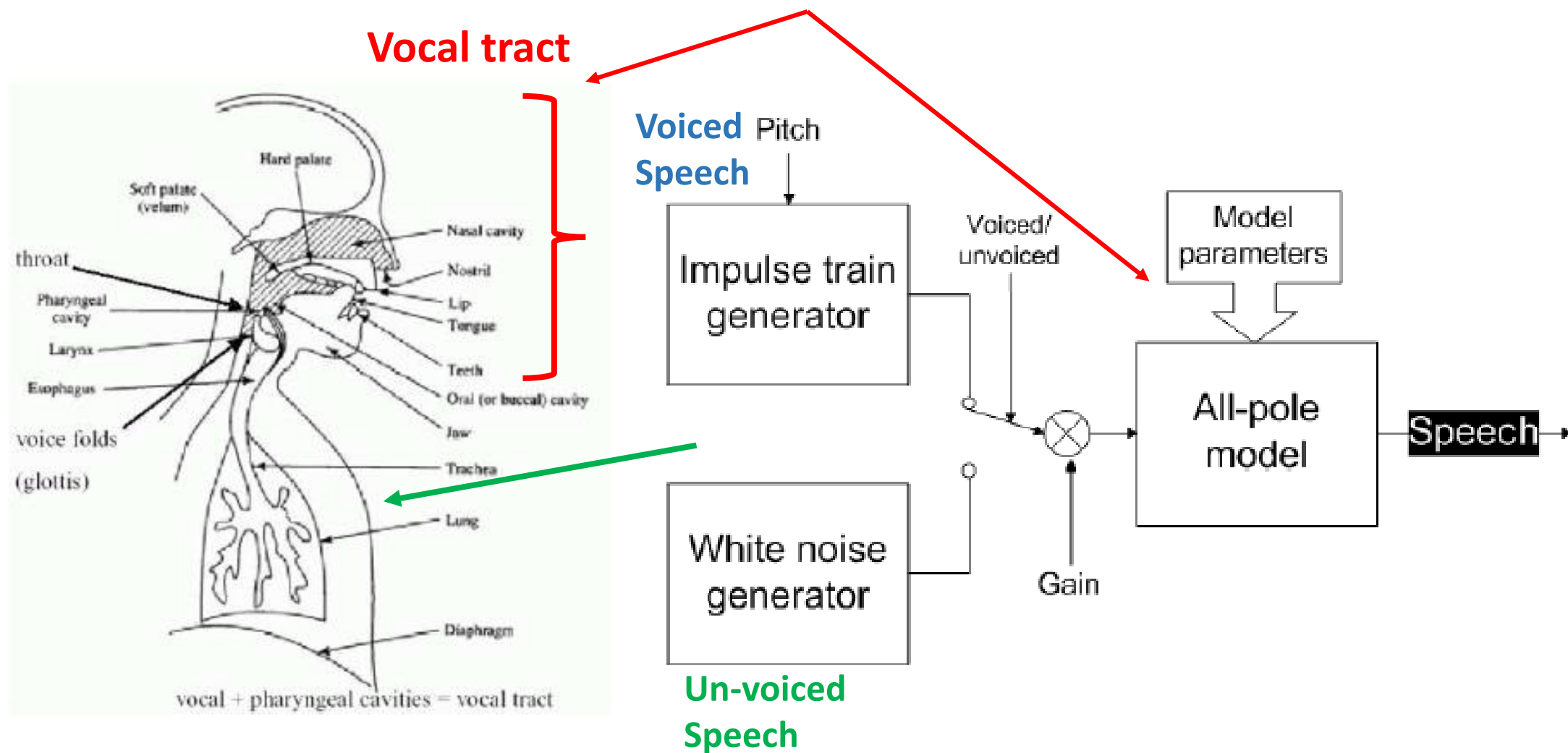
vocal + pharyngeal cavities = vocal tract

# Phonemes

- Phonemes
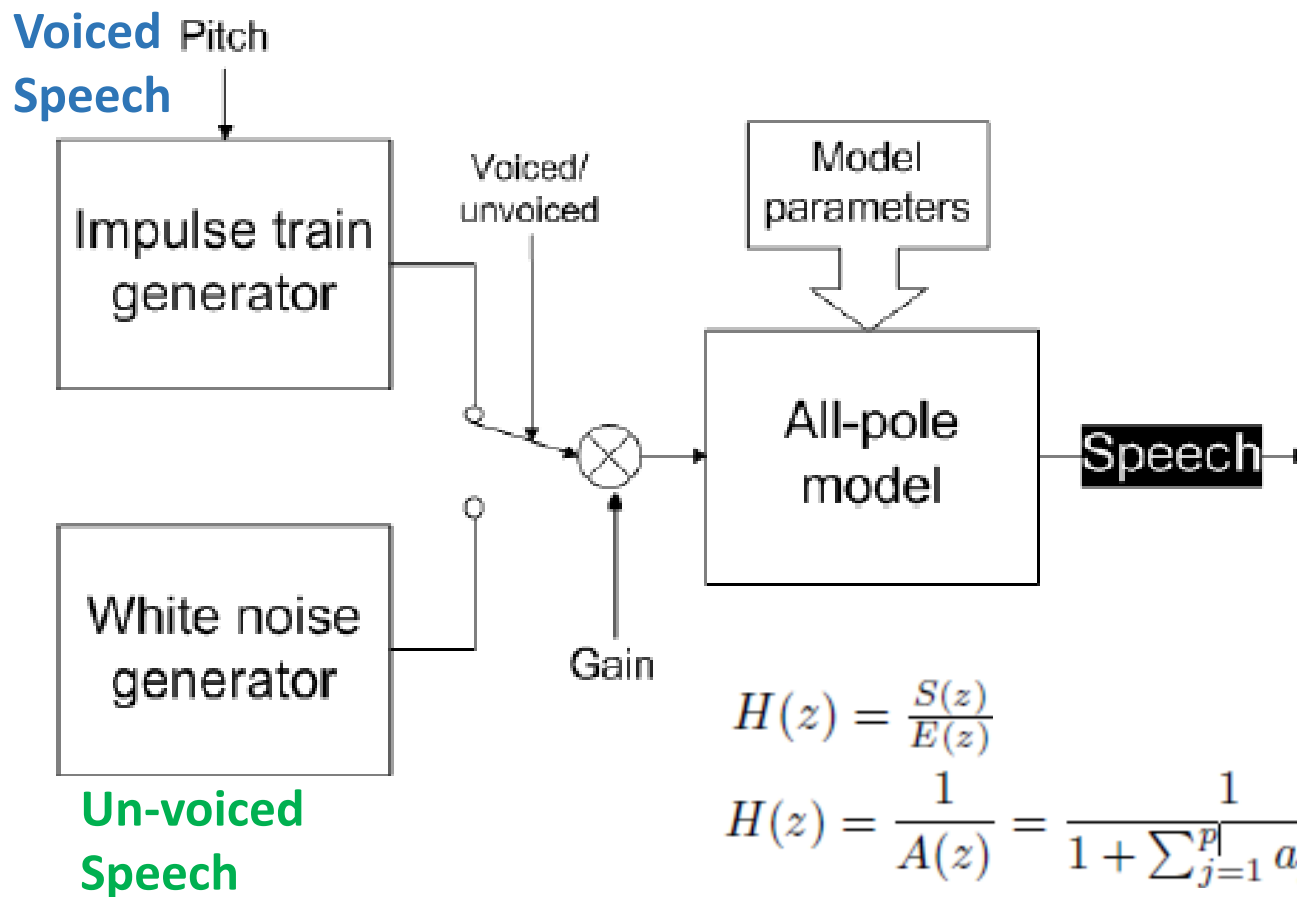  - Voiced sound (vowels: a, e, i, o, u)
  - Unvoiced sound (consonants: t, h, p …)

# A Mathematical Model for Speech Production

# A Mathematical Model for Speech Production

# A Mathematical Model for Speech Production



**Voiced** Pitch
**Speech**

Impulse train generator

Voiced/ unvoiced

Model parameters

All-pole model

Speech

White noise generator

Gain

**Un-voiced Speech**

$$H(z) = \frac{S(z)}{E(z)}$$

$$H(z) = \frac{1}{A(z)} = \frac{1}{1 + \sum_{j=1}^{p} a_j z^{-j}}$$

$$s(n) = -\sum_{j=1}^{p} a_j s(n-j) + e(n)$$

$a_j = model\ parameter$

p = model order
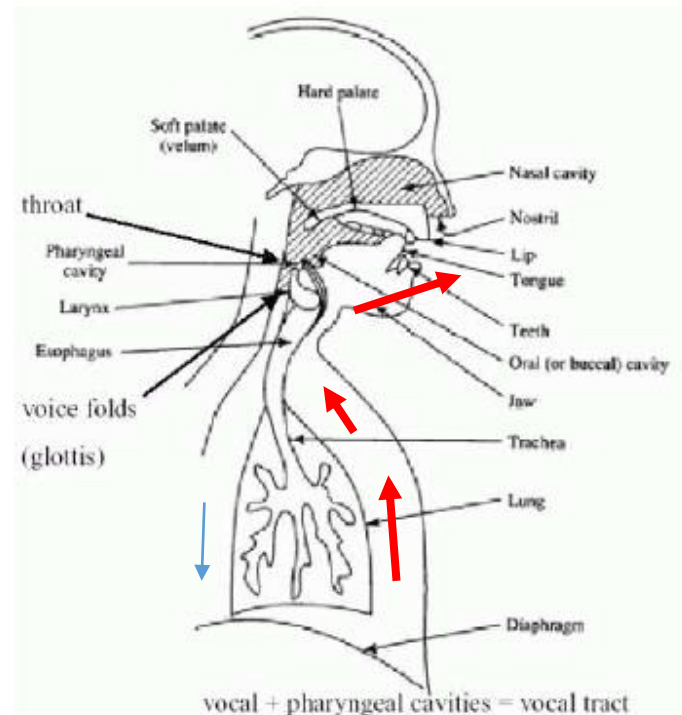
$$S(z) = G \cdot H(z)$$

# Selecting p

- **Selecting model order p:**

Speech spectrum consists of 3-4 formants ( resonant peaks)

Typical choice p = 10
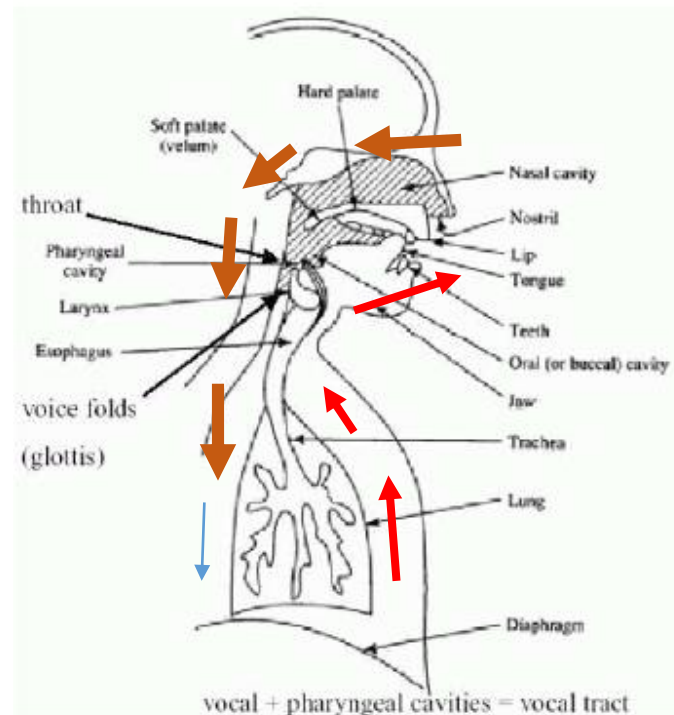
- **Computing the residuals e(n) and the variance $\sigma_e^2$ (Gain):**

$$s(n) = -\sum_{j=1}^{p} a_j s(n-j) + e(n)$$

- **Computing the residuals e(n) and the variance $\sigma_e^2$:**

  Gain = $|\sigma_e|$

$$\hat{e}(n) = s(n) + \sum_{j=1}^{p} \hat{a}_j s(n-j)$$



$$s(n) = -\sum_{j=1}^{p} a_j s(n-j) + e(n)$$

Hard palate

Soft palate (velum)

throat

Nasal cavity

Nostril

Pharyngeal cavity

Lip

Tongue

Larynx

Teeth

Esophagus

Oral (or buccal) cavity

Jaw

voice folds

(glottis)

Trachea

Lung

Diaphragm

vocal + pharyngeal cavities = vocal tract

# Vocal cord excitation for voiced and unvoiced



amplitude of speech signal of "six"

frame of /s/

frame of /I/

# Speech Compression:

- Non-perceivable loss by human ears.

**For a 10- 20ms block of speech**

- Set of parameters for LPC speech model

$\{a_j \; for \; j = 1 \ldots p = 10, \sigma_e^2, \text{pitch period}\}$

# Speech Compression:

- Non-perceivable loss by human ears.

**For a 10-20ms block of speech**

- Set of parameters for LPC speech model
$\{a_j \ for \ \ j = 1 \ldots p = 10, \sigma_e^2, \text{pitch period}\}$

If each parameter is encoded using 10-bits, total bits to be transmitted will be 12*10=120 bits

- For direct encoding, first sample at 8kHz that gives 160 samples. Convert to bits 160*8= 1280 bits.

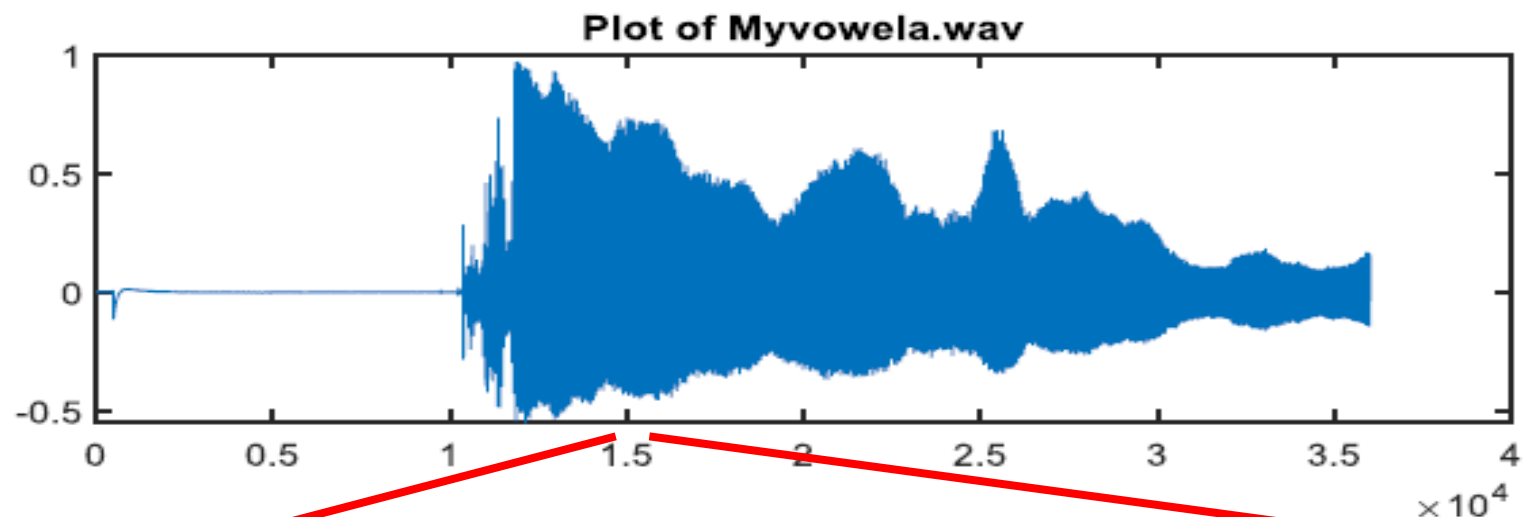# Speech Compression:

- Non-perceivable loss by human ears.

**For a 10-20ms block of speech**
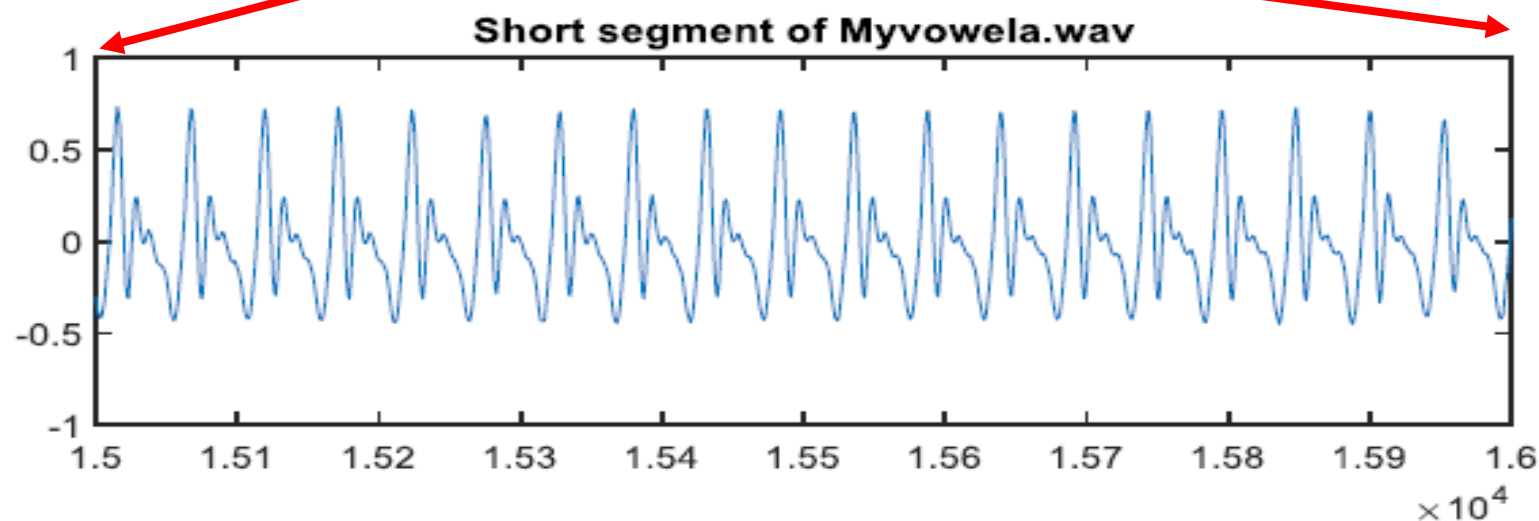- Set of parameters for LPC speech model $\{a_j \ for \ j = 1 \ldots. p = 10, \sigma_e^2, \text{pitch period}\}$

If each parameter is encoded using 10-bits, total bits to be transmitted will be 12*10=120 bits

- For direct encoding, first sample at 8kHz that gives 160 samples. Convert to bits 160*8= 1280 bits.
- **LPC gives 10 times compression.**

# Block based processing of Speech Signal

## Plot of Myvowela.wav

## Short segment of Myvowela.wav

Model a vocal tract by an all pole-filter for each block.
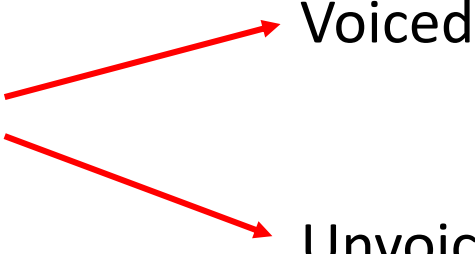
$$\hat{e}(n) = s(n) + \sum_{j=1}^{p} \hat{a}_j^{(i)} s(n-j)$$

# Re-synthesize the Speech

- Use LPC residuals as the input of the filter:

$$\hat{s}(n) = -\sum_{j=1}^{p} \hat{a}_j^{(i)} \hat{s}(n-j) + \hat{e}(n)$$

- Signal
  - Voiced
  - Unvoiced

# Re-synthesize the Speech

- Use LPC residuals as the input of the filter:

$$\hat{s}(n) = -\sum_{j=1}^{p} \hat{a}_j^{(i)} \hat{s}(n-j) + \hat{e}(n)$$

- Signal

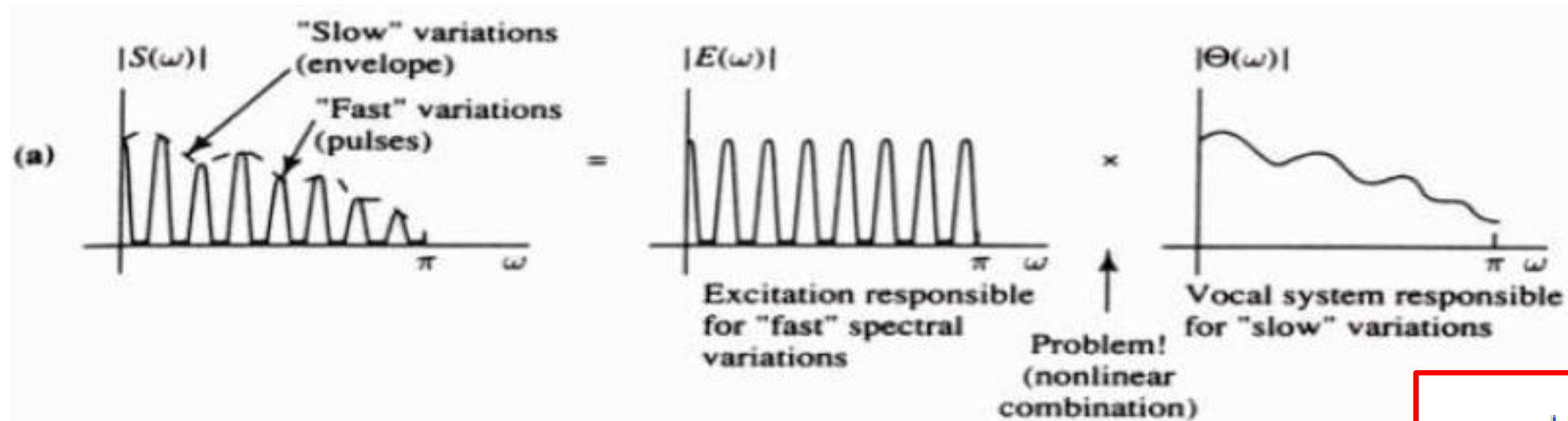Voiced → Pitch period

Unvoiced → Variance $\sigma_e^2$

# Methods for automatically estimating pitch periods

- Pitch estimation using ACF

$$R_s(\tau) = \sum_{j=(i-1)L+1}^{iL} \hat{e}(j)\hat{e}(j+\tau),$$

- Pitch estimation using Cepstrum domain analysis.

$$\mathbf{c}_i = \left| FFT^{-1}\{log(|FFT(\mathbf{x}_i)|\} \right|$$

(a) $|S(\omega)|$ — "Slow" variations (envelope), "Fast" variations (pulses)

$|E(\omega)|$ — Excitation responsible for "fast" spectral variations

$|\Theta(\omega)|$ — Vocal system responsible for "slow" variations

Problem! (nonlinear combination)

$$c_i = \left| FFT^{-1}\{log(|FFT(\mathbf{x}_i)|\} \right|$$

(b) $log\,|S(\omega)| = C_s(\omega)$

$log\,|E(\omega)|$

$log\,|\Theta(\omega)|$

Problem solved! (linear combination)

(c) $c_s(n)$ — Low quefrency energy, High quefrency

$c_e(n)$ [approx]

$c_\Theta(n)$