# Pretext tasks
# 4. CPC

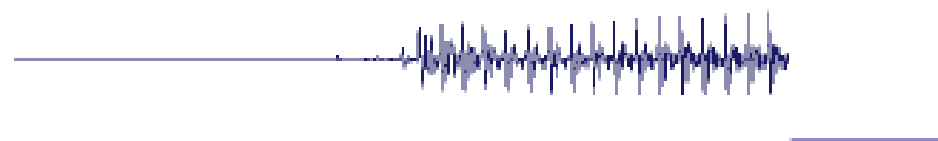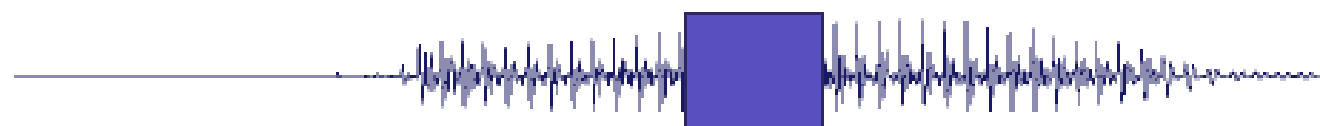| | | | |
|---|---|---|---|
| **1** **2** | SELF-PREDICTION | INNATE RELATIONSHIP (Context-based) | 1. ROTATION<br><br>2. RELATIVE POSITION | IMAGE |
| **3** | CONTRASTIVE LEARNING | INTER-SAMPLE CLASSIFICATION | 1. Instance Discrimination<br>2. SimCLR [Contrastive Loss]<br>3. Theory – Guarantees / Bounds | IMAGE |
| **4** | CONTRASTIVE LEARNING | INTER-SAMPLE CLASSIFICATION | Contrastive Predictive Coding (CPC), [NCE, InfoNCE Loss] | AUDIO/ SPEECH |
| **5** | SELF-PREDICTION | GENERATIVE (VAE) | 1. AE – Variational Bayes<br><br>2. VQ-VAE + AR | IMAGE<br><br>AUDIO/ SPEECH |
| **6** | SELF-PREDICTION | GENERATIVE (AR) | 1. AR-LM – GPT<br>2. Masked-LM – BERT | LANGUAGE |
| **7** | SELF-PREDICTION | MASKED-GEN (Masked LM for ASR) | 1. Wav2Vec / 2.0<br>2. HuBERT | AUDIO/ SPEECH |

# Learning with or without supervision – speech and audio
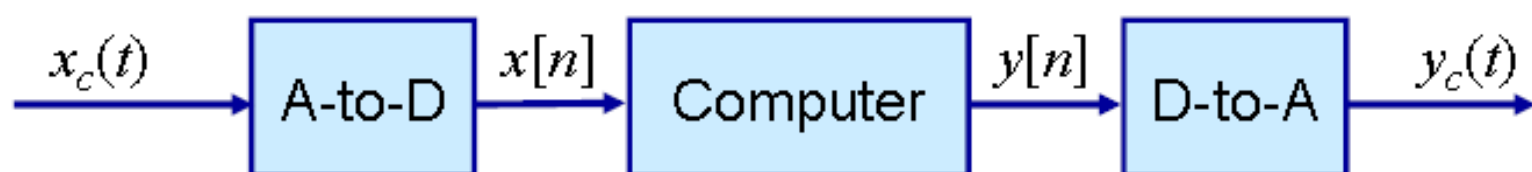
- Next frame prediction

- Masked prediction

# Speech Waveform Basics

1 Second

# Digital Processing of Analog Signals

$$x_c(t) \longrightarrow \boxed{\text{A-to-D}} \xrightarrow{x[n]} \boxed{\text{Computer}} \xrightarrow{y[n]} \boxed{\text{D-to-A}} \xrightarrow{y_c(t)}$$

- **A-to-D conversion:** bandwidth control, sampling and quantization
- **Computational processing:** implemented on computers or ASICs with finite-precision arithmetic
    - **basic numerical processing:** add, subtract, multiply (scaling, amplification, attenuation), mute, …
    - **algorithmic numerical processing:** convolution or linear filtering, non-linear filtering (e.g., median filtering), difference equations, DFT, inverse filtering, MAX/MIN, …
- **D-to-A conversion:** re-quantification* and filtering (or interpolation) for reconstruction

# Discrete-Time Signals

☐ A sequence of numbers

☐ Mathematical representation:

$$x = \{x[n]\}, \quad -\infty < n < \infty$$

☐ Sampled from an analog signal, $x_a(t)$, at time $t = nT$,

$$x[n] = x_a(nT), \quad -\infty < n < \infty$$

☐ $T$ is called the **sampling period**, and its reciprocal,
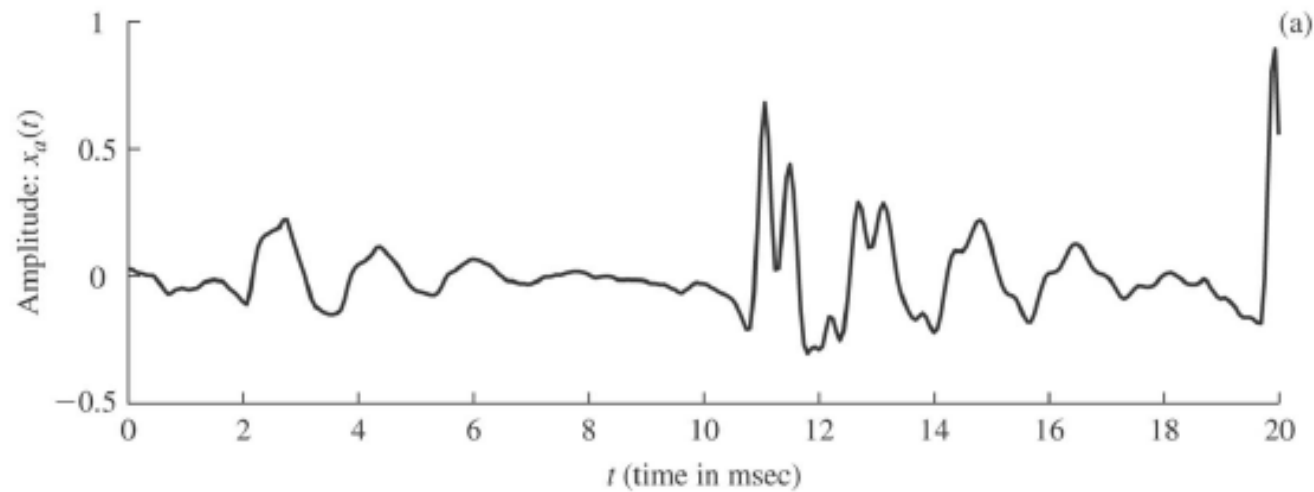
$F_S = 1/T$, is called the **sampling frequency**

$$F_S = 8000 \text{ Hz} \quad \leftrightarrow \quad T = 1/8000 = 125\,\mu\sec$$

$$F_S = 10000 \text{ Hz} \quad \leftrightarrow \quad T = 1/10000 = 100\,\mu\sec$$

$$F_S = 16000 \text{ Hz} \quad \leftrightarrow \quad T = 1/16000 = 62.5\,\mu\sec$$

$$F_S = 20000 \text{Hz} \quad \leftrightarrow \quad T = 1/20000 = 50\,\mu\sec$$

# Speech Waveform Display



plot( );

stem( );

# Discrete Signals



sample

Sampled Sinusoid
$5\sin(2\pi nT)$

quantize

Analog sinusoid,
$5\sin(2\pi x)$

Discrete sinusoid
$\text{round}[5\sin(2\pi nT)]$

Quantized sinusoid
$\text{round}[5\sin(2\pi x)]$

quantize

sample

# Discrete-Time (DT) Signals are Sequences



- $x[n]$ denotes the "sequence value at 'time' $n$"
- Sources of sequences:
  - Sampling a continuous-time signal
    $$x[n] = x_c(nT) = x_c(t)|_{t=nT}$$
  - Mathematical formulas – generative system
    e.g.,  $x[n] = 0.3 \cdot x[n\text{-}1] \text{-}1; \quad x[0] = 40$

# Speech waveform example



THE NEW BRICKS FELL OVER

# Speech spectrogram example



**Figure 2.10** Wide-band spectrogram of the speech waveform shown in Figure 2.9. The dynamic range of the grey scale in the display is 50 dB, so very weak sounds are clearly visible.

# 🔊 She had your dark suit in…

# "Wideband" Spectrogram

She had your dark suit in.



SH  IY HH   AE  D  AXR D  AA R   K   S   UW   IH N

Y                              T

frequency

time

# Spectrogram

- ❑ Speech is a continuous evolution of the vocal tract
- ❑ Spectrogram shows time-frequency evolution
- ❑ Represented as a time-series of short-time spectra

# Spectrum Basics

# Fourier Series (Calculus required)

Continuous functions are often approximated by linear combinations of sine and cosine functions. For instance, a continuous function might represent a sound wave, an electric signal of some type, or the movement of a vibrating mechanical system.

For simplicity, we consider functions on $0 \le t \le 2\pi$. It turns out that any function in $C[0, 2\pi]$ can be approximated as closely as desired by a function of the form

$$\frac{a_0}{2} + a_1 \cos t + \cdots + a_n \cos nt + b_1 \sin t + \cdots + b_n \sin nt \tag{4}$$

for a sufficiently large value of $n$. The function (4) is called a **trigonometric polynomial**. If $a_n$ and $b_n$ are not both zero, the polynomial is said to be of **order $n$**. The connection between trigonometric polynomials and other functions in $C[0, 2\pi]$ depends on the fact that for any $n \ge 1$, the set

$$\{1, \cos t, \cos 2t, \ldots, \cos nt, \sin t, \sin 2t, \ldots, \sin nt\} \tag{5}$$

is orthogonal with respect to the inner product

$$\langle f, g \rangle = \int_0^{2\pi} f(t) g(t) \, dt \tag{6}$$

$y_1 = a_1 \sin(2\pi{*}ft)$

$y_2 = a_2 \sin(2\pi{*}2ft)$

$y_5 = a_5 \sin(2\pi{*}5ft)$

$y_6 = a_6 \sin(2\pi{*}6ft)$

$y_3 = a_3 \sin(2\pi{*}3ft)$

$y_4 = a_4 \sin(2\pi{*}4ft)$

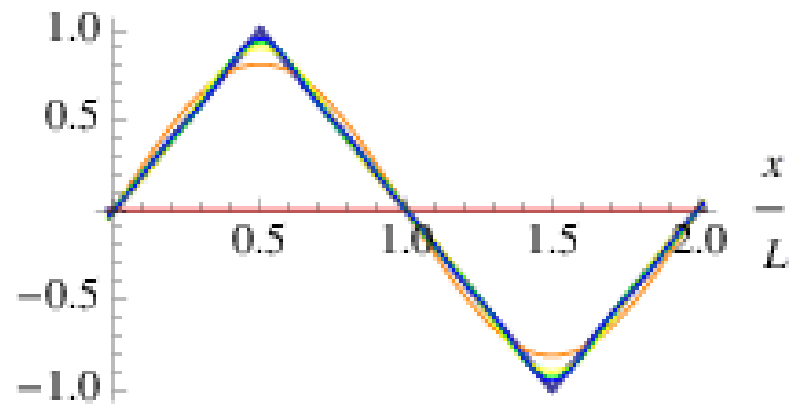$y = a_1 \sin(2\pi{*}ft) + a\ \sin(2\pi{*}2ft) + a\ \sin(2\pi{*}3ft)$
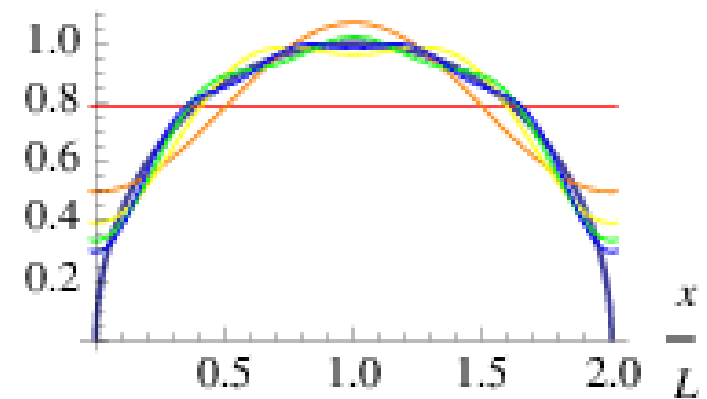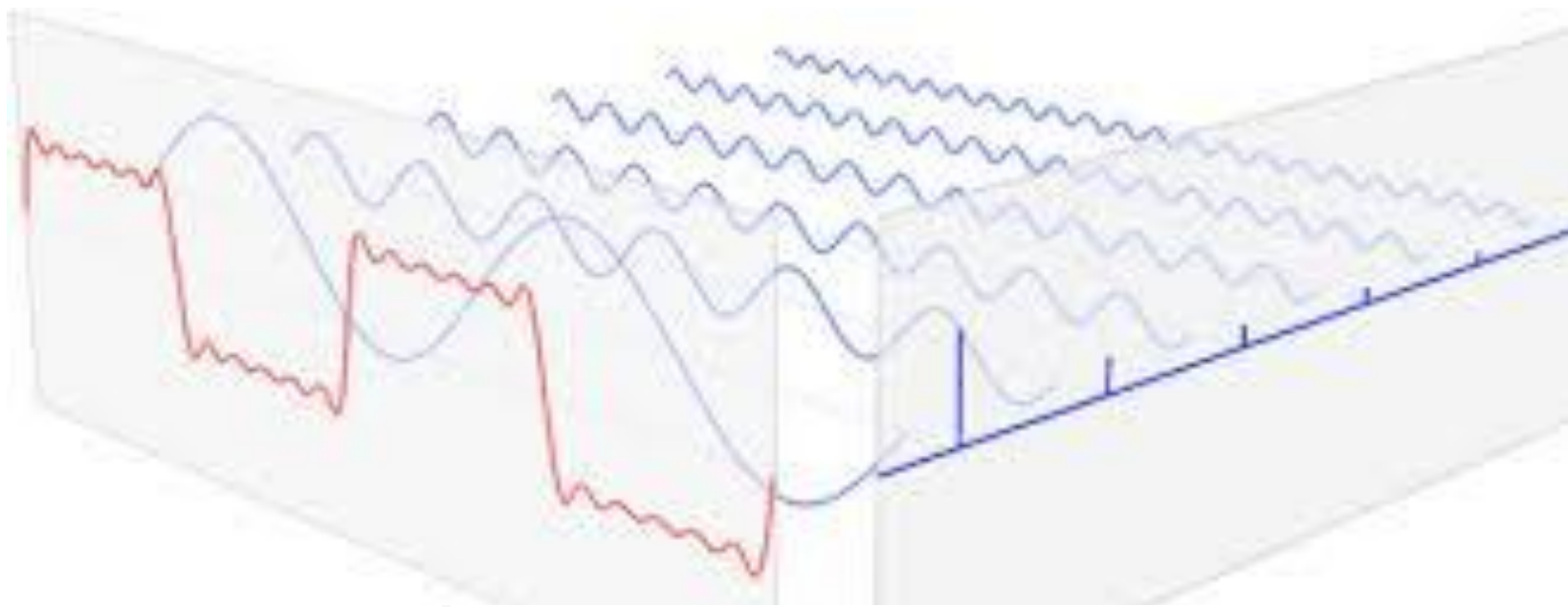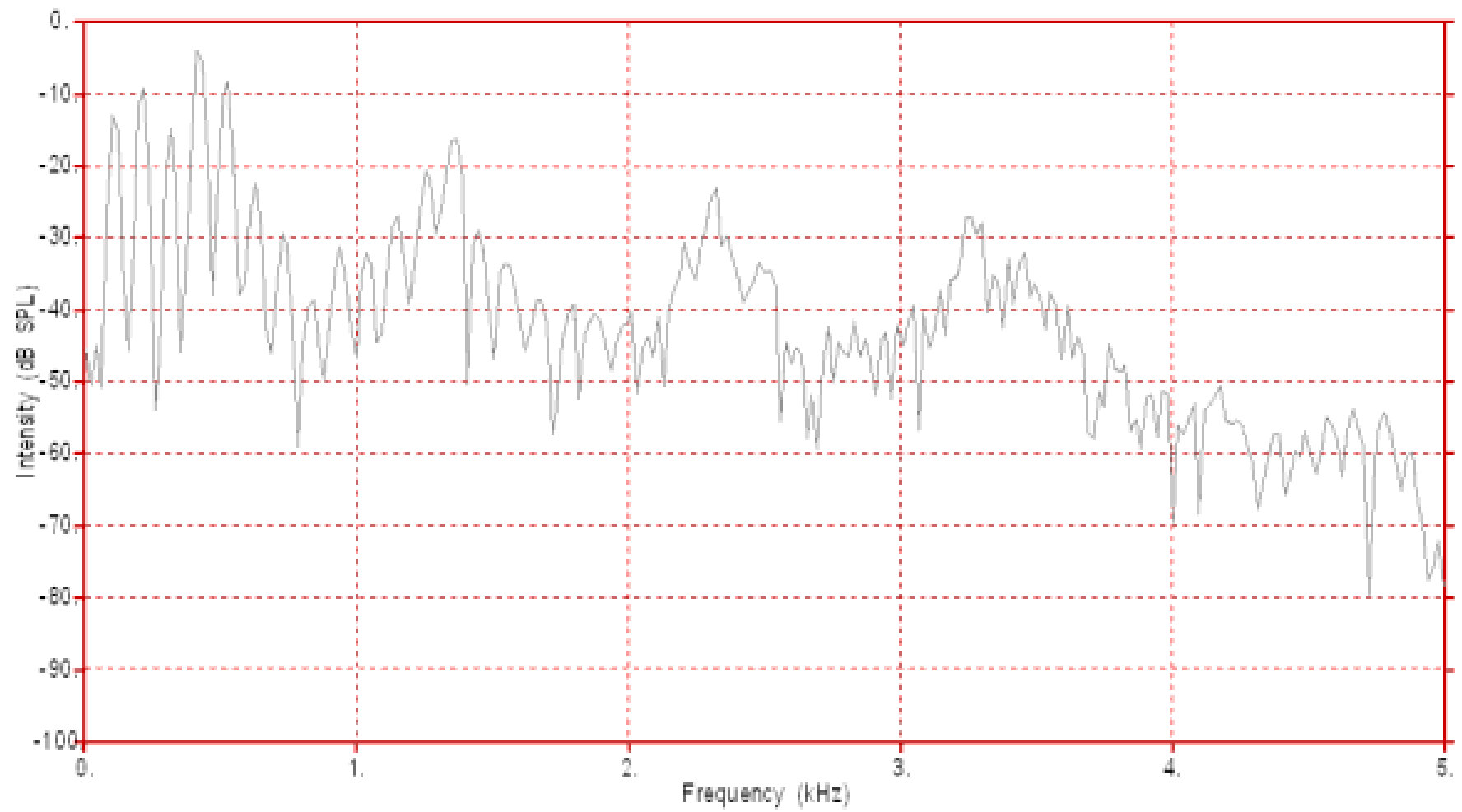$+ a_4 \sin(2\pi{*}4ft) + a_5 \sin(2\pi{*}5ft) + a_6 \sin(2\pi{*}6ft)$
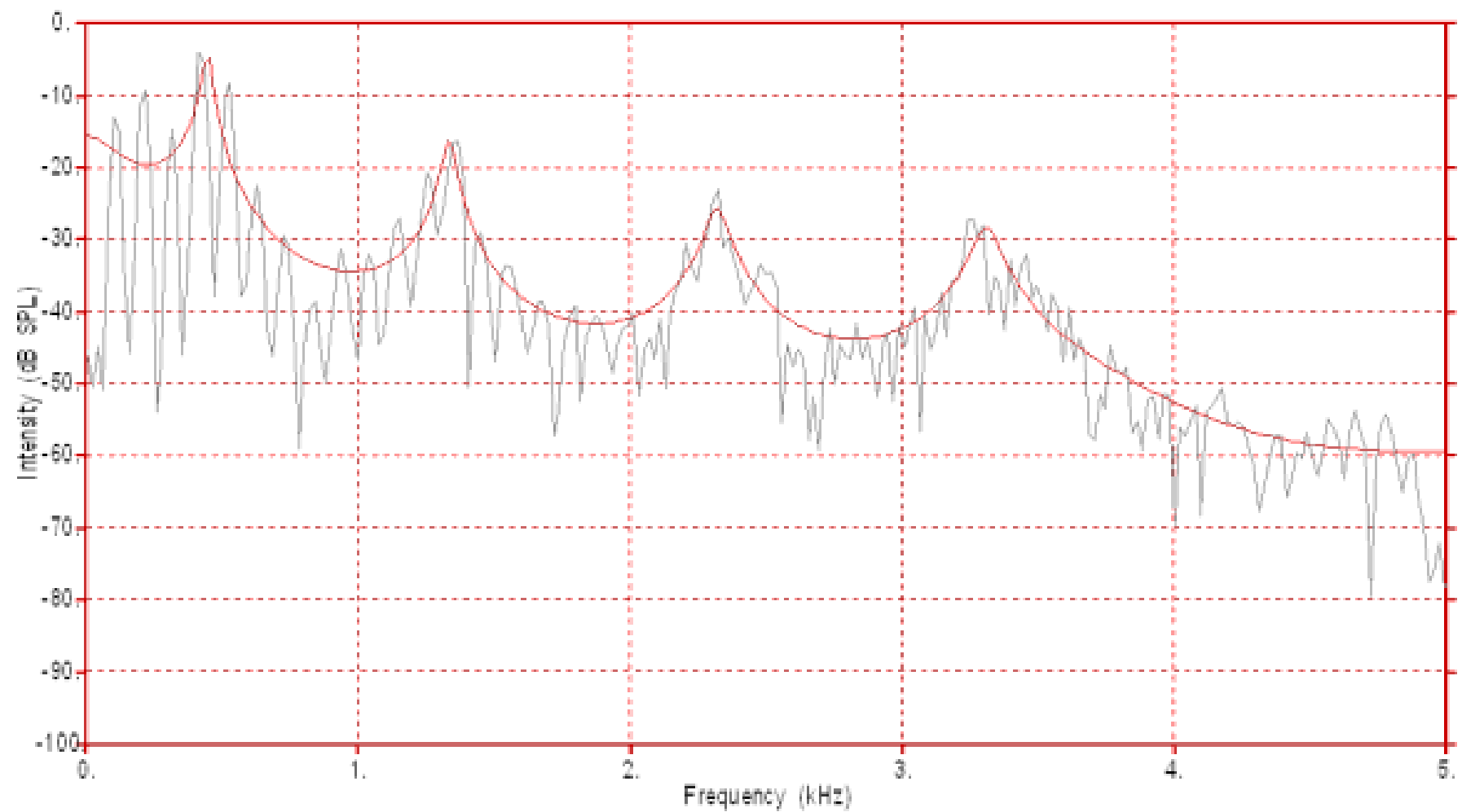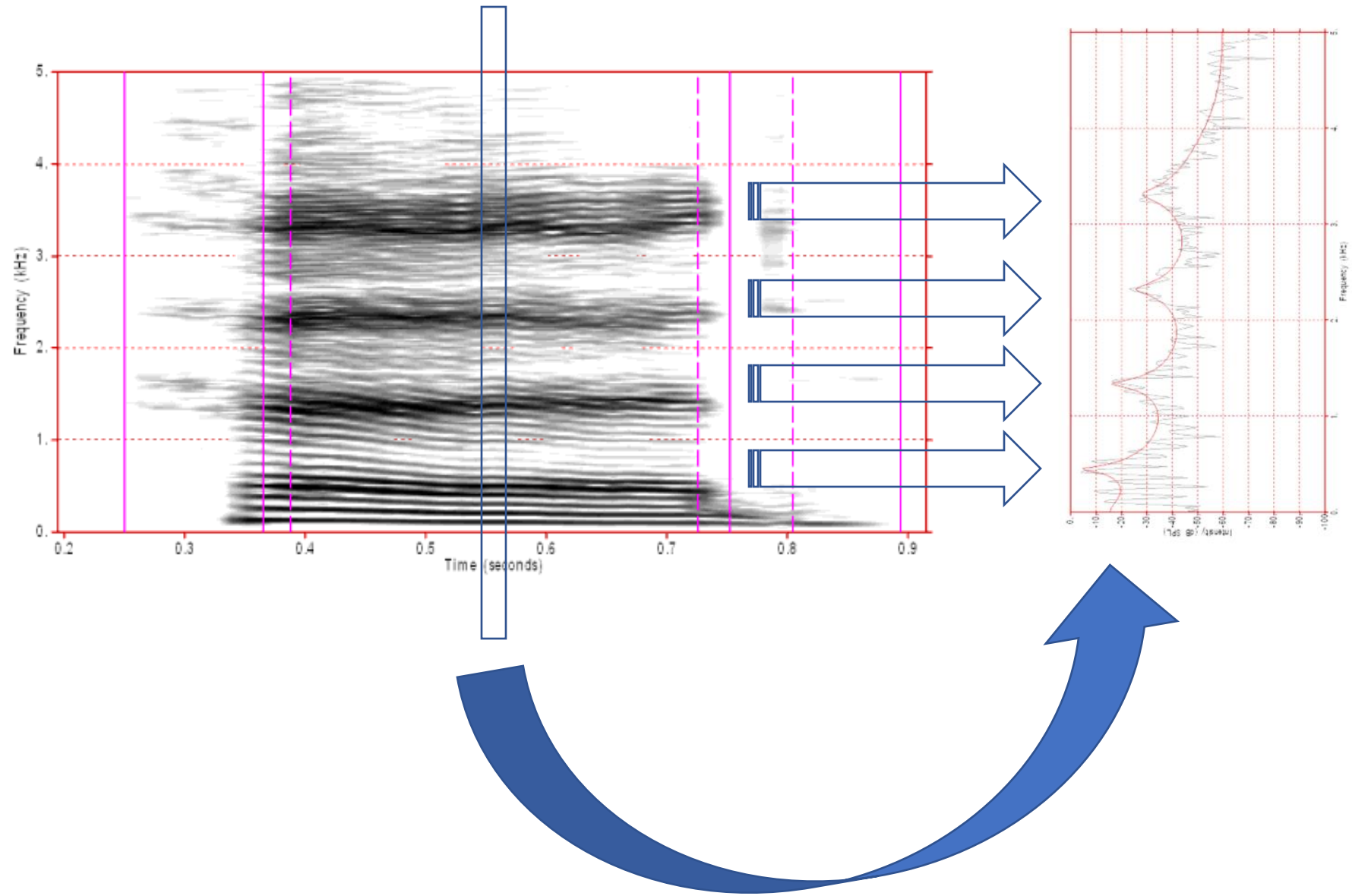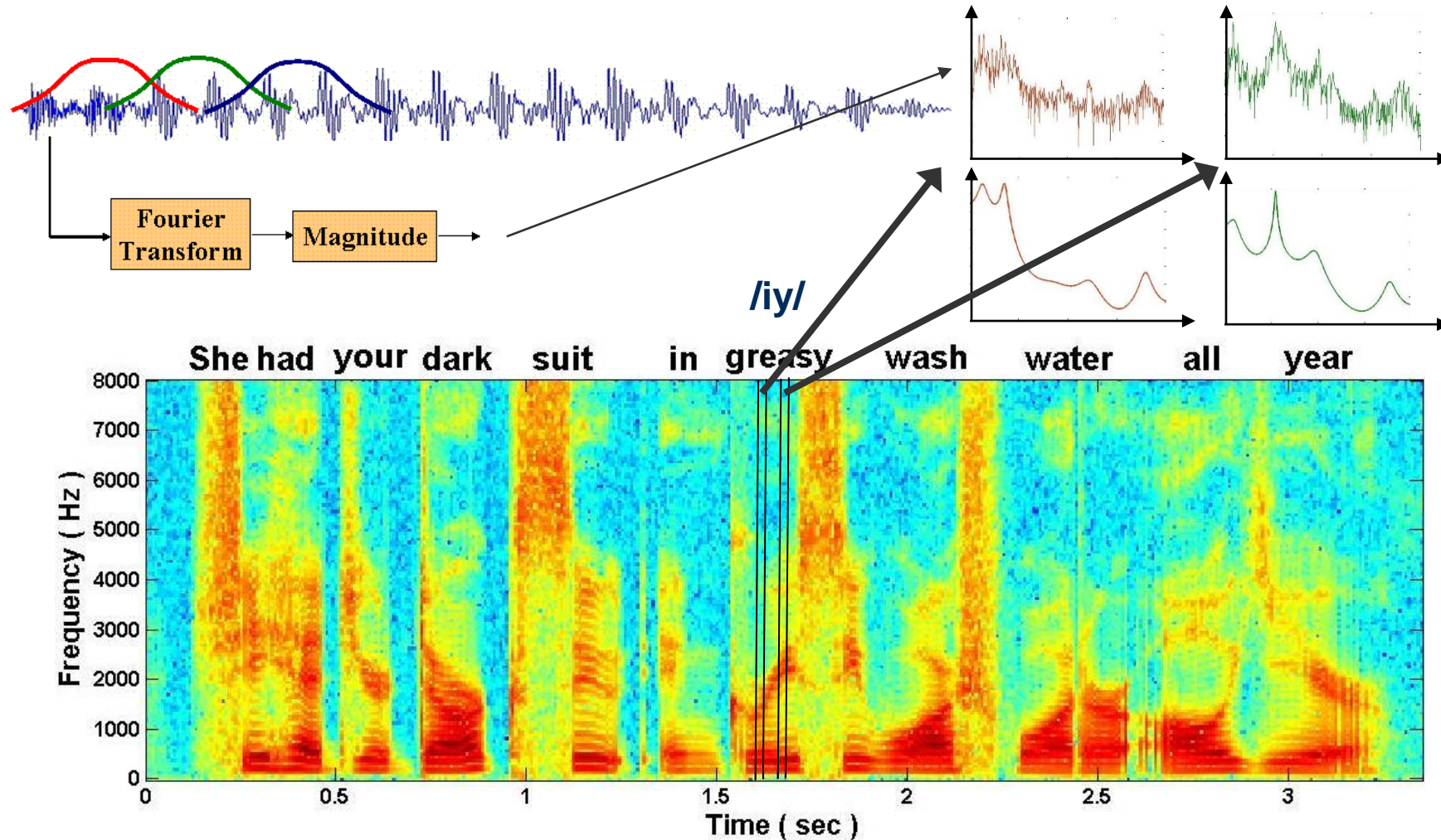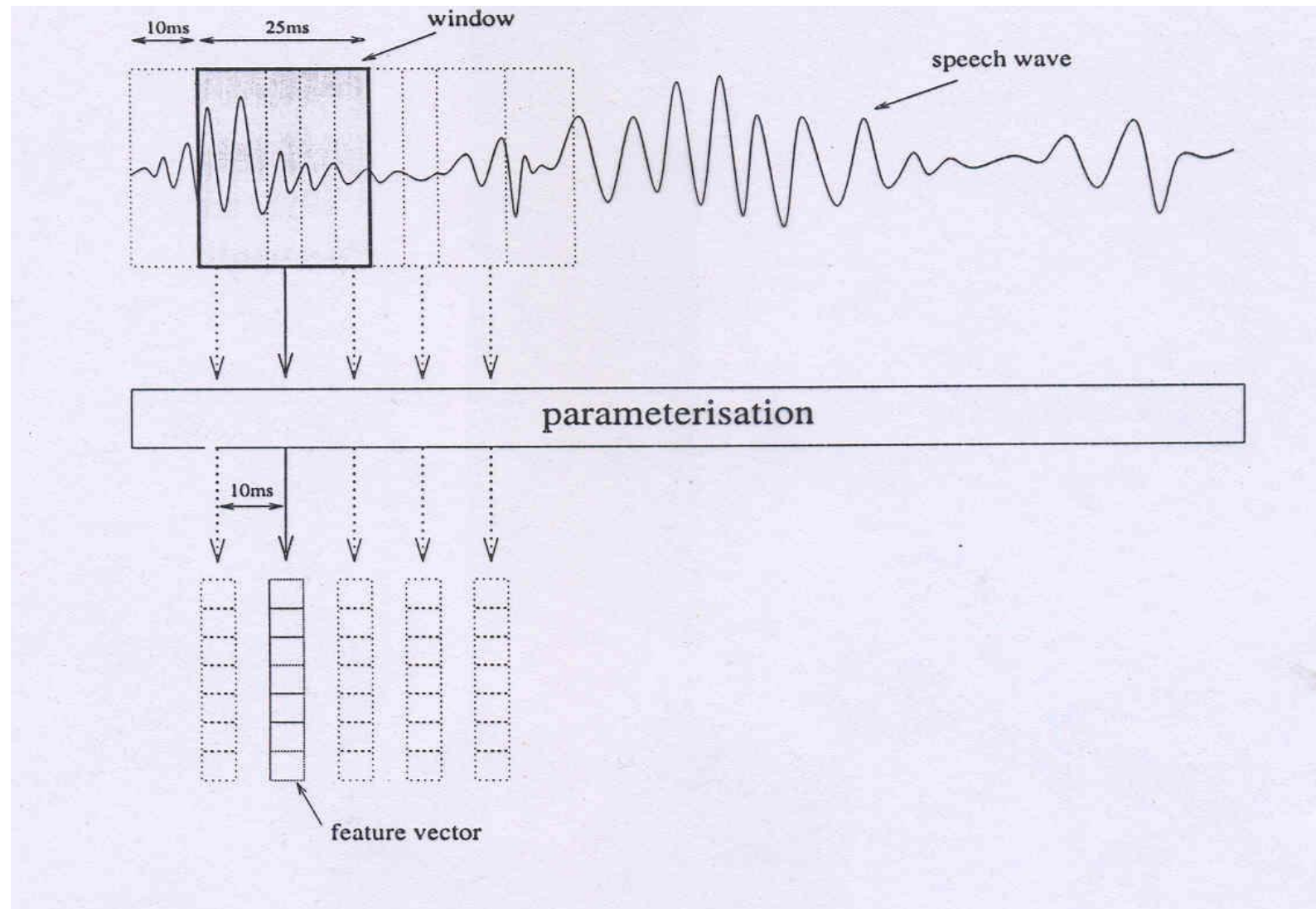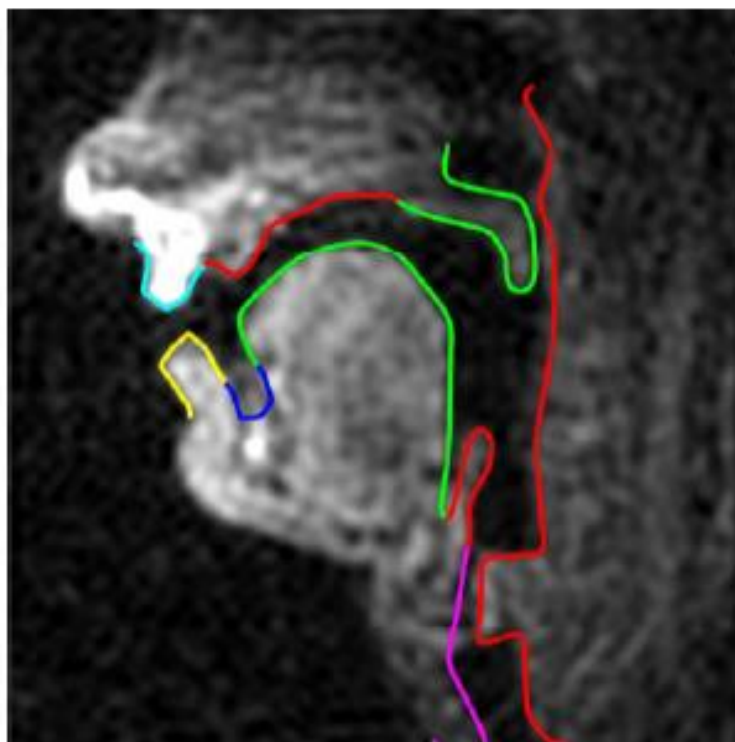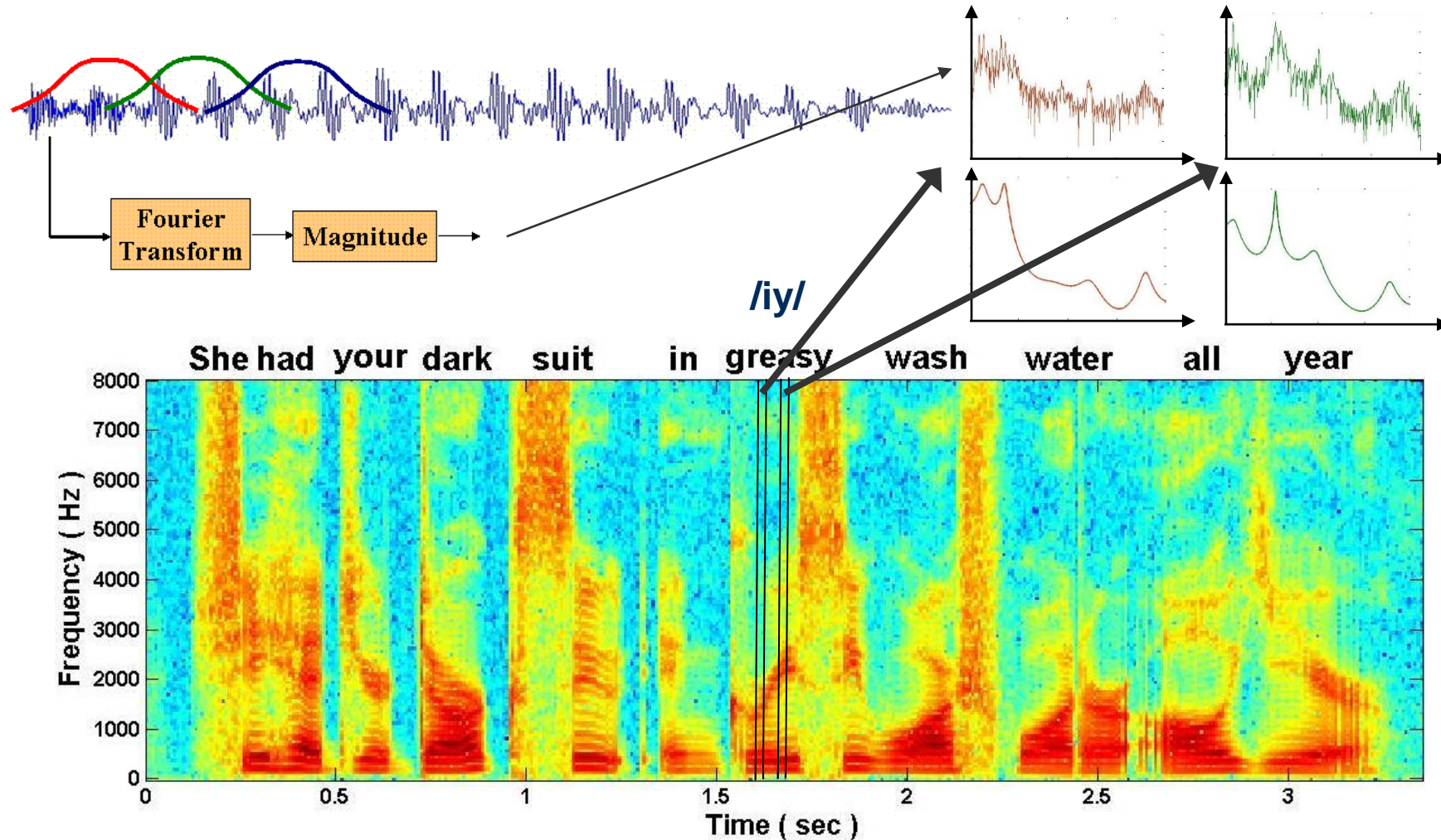
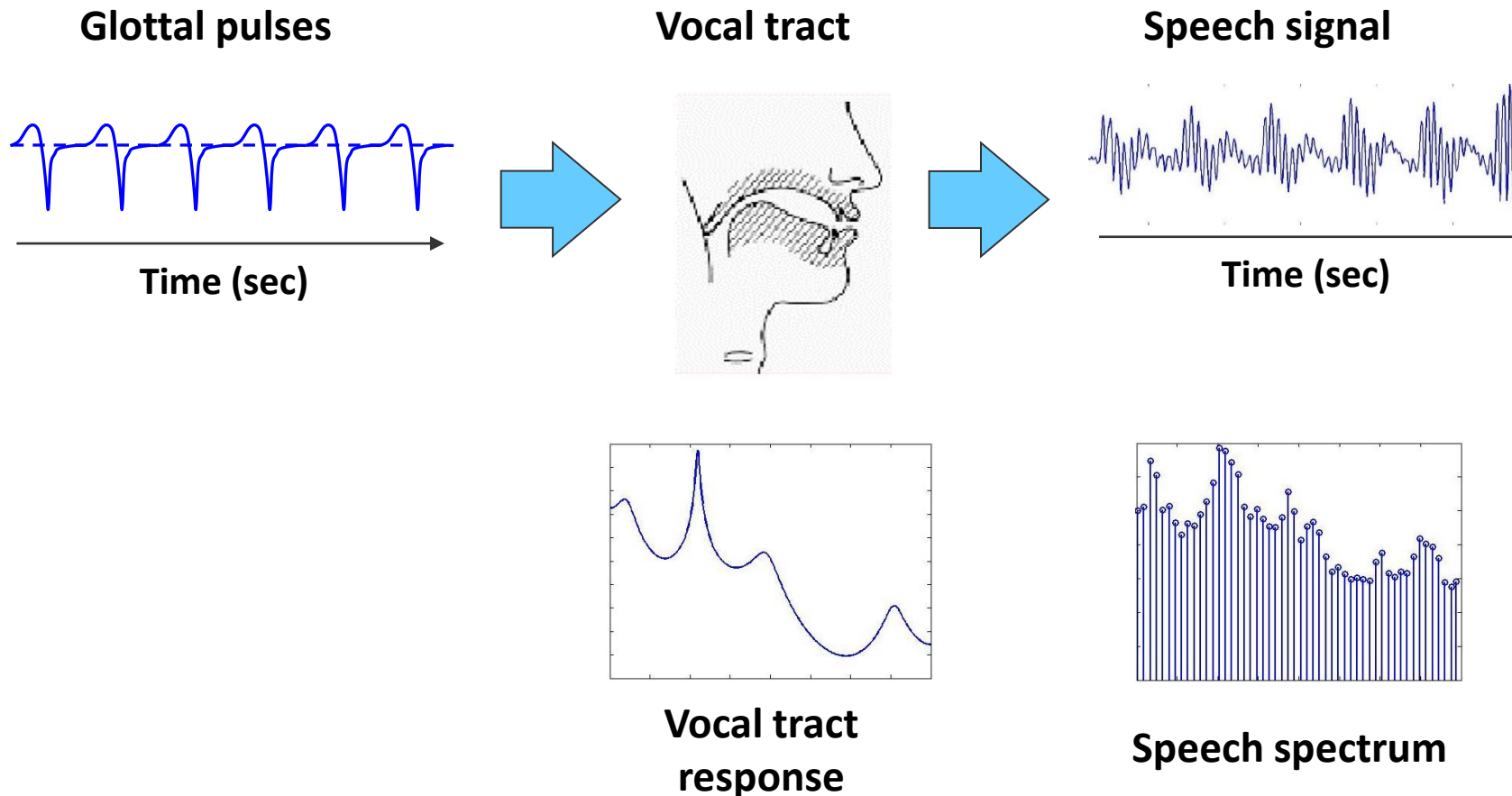square wave

sawtooth wave

triangle wave

semicircle

# Spectrogram

- Speech is a continuous evolution of the vocal tract
- Spectrogram shows time-frequency evolution
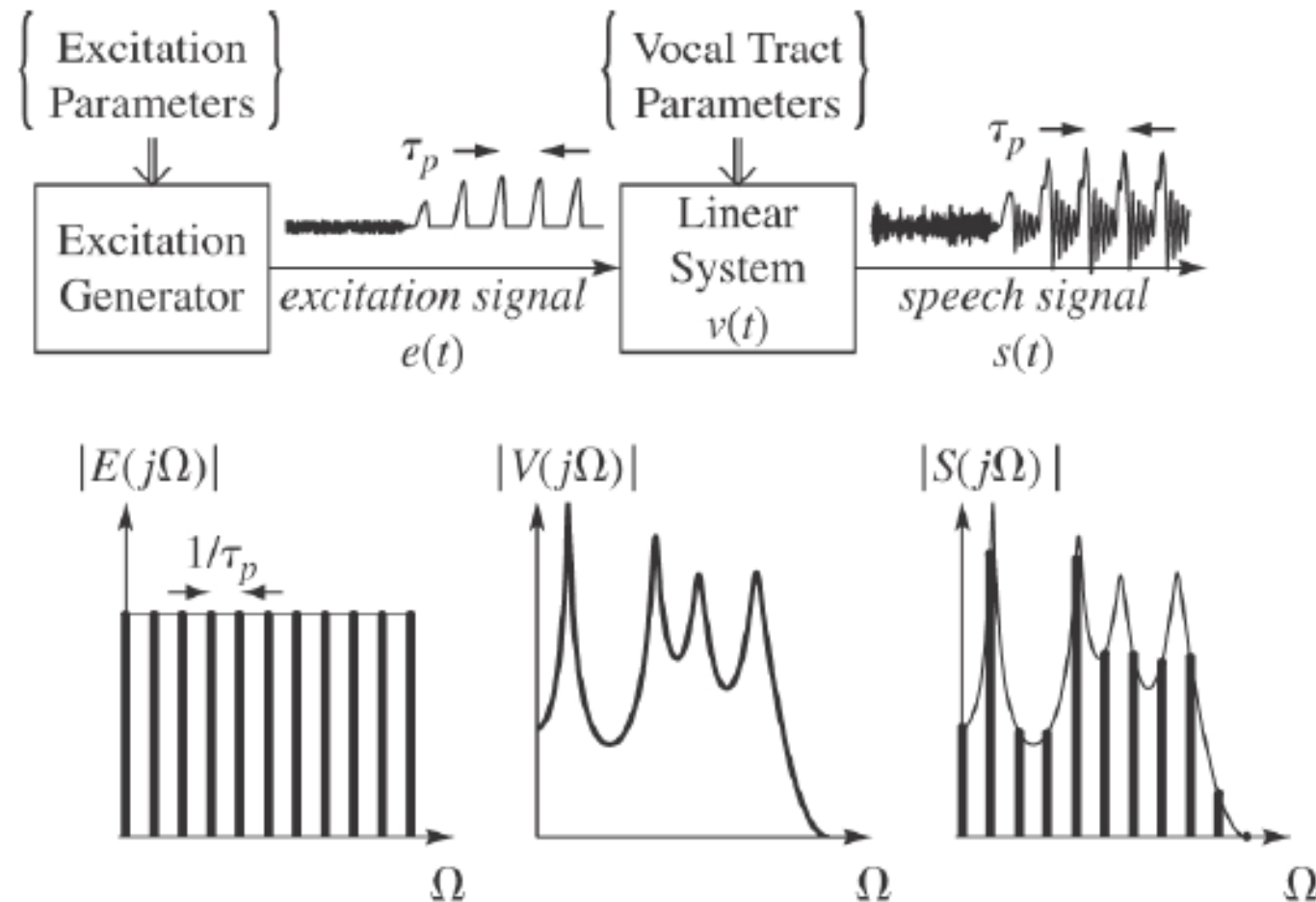- Represented as a time-series of short-time spectra

# Auto-regressive Model and Speech Spectrum

# Short-time Analysis and Parameterization

# MRI of Speech (Prof. Shri Narayanan, USC)



(a)

(b)

# Spectrogram

❑ Speech is a continuous evolution of the vocal tract

❑ Spectrogram shows time-frequency evolution

❑ Represented as a time-series of short-time spectra

## Source-Filter Model

❑ Features based on speech production model: Source-filter interaction
 – Anatomical structure (vocal tract / glottis) conveyed in speech spectrum

**Glottal pulses**



Time (sec)

**Vocal tract**



**Speech signal**



Time (sec)



**Vocal tract response**



**Speech spectrum**

# Source-System Model of Speech Production

# Linear Prediction based Speech Production Model



$$\frac{1}{A(z)} = \frac{1}{1 - \sum_{i=1}^{p} a_i z^{-i}}$$

- Vocal Tract • $\Longleftrightarrow$ • H(z) (LPC Filter)
- Air • $\Longleftrightarrow$ • u(n) (Innovations)
- Vocal Cord Vibration • $\Longleftrightarrow$ • V (voiced)
- Vocal Cord Vibration Period • $\Longleftrightarrow$ • T (pitch period)
- Fricatives and Plosives • $\Longleftrightarrow$ • UV (unvoiced)
- Air Volume • $\Longleftrightarrow$ • G (gain)

# LP Analysis: Envelope (Filter) & Excitation (Source)

**Speech Signal S(n)**

**S(w) with LP spectral envelope superimposed**

LP Envelope

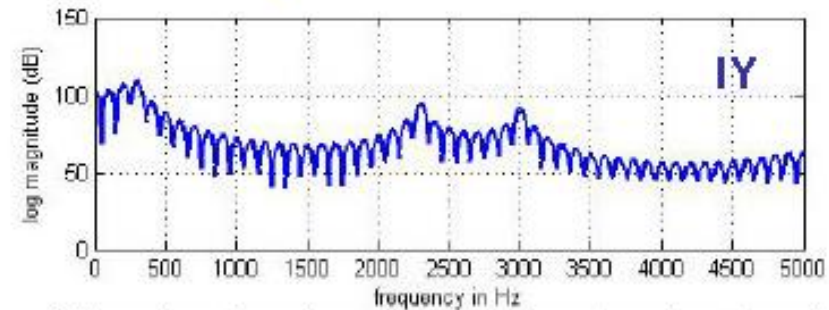**Excitation Signal E(n)**

Pitch Period

# Spectral slices

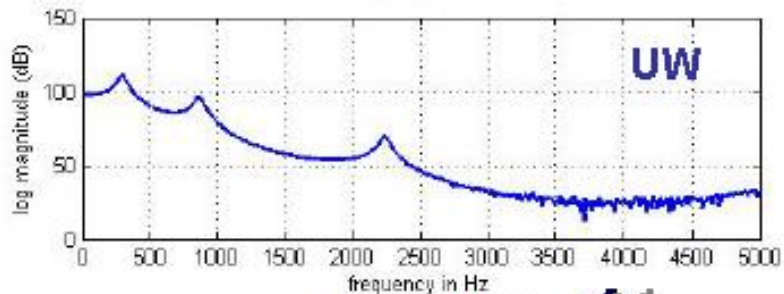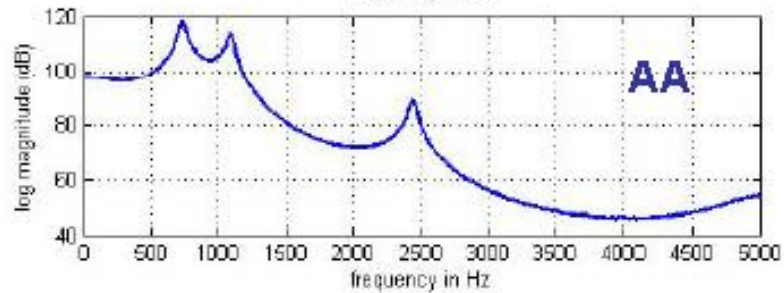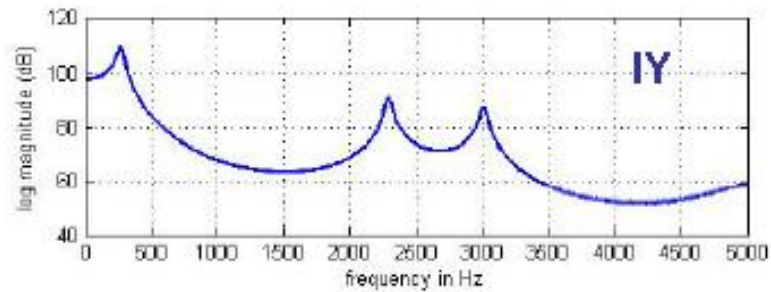## Spectral Envelopes



Vowel / iy / LP Spectrum

# Canonic Vowel Spectra


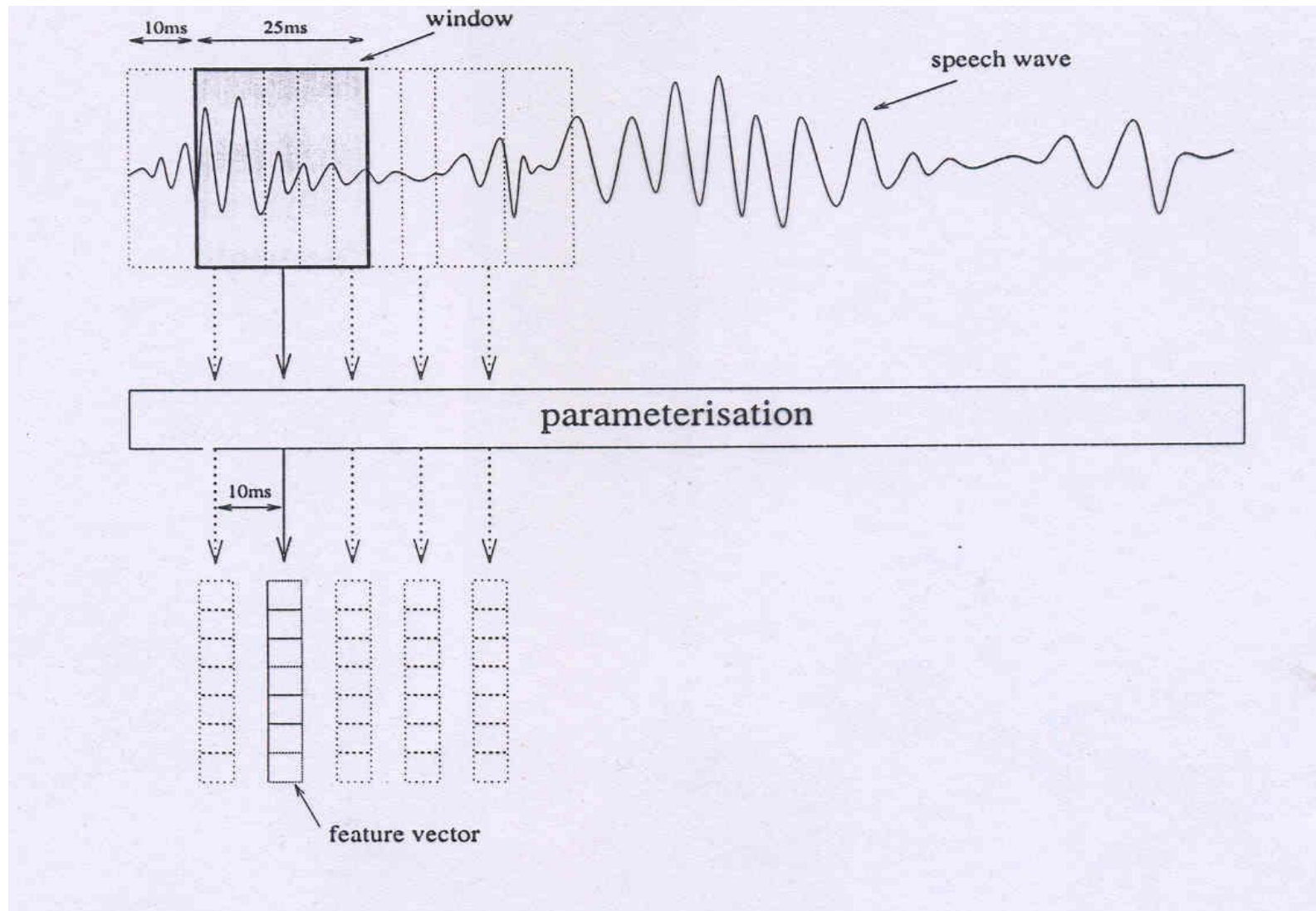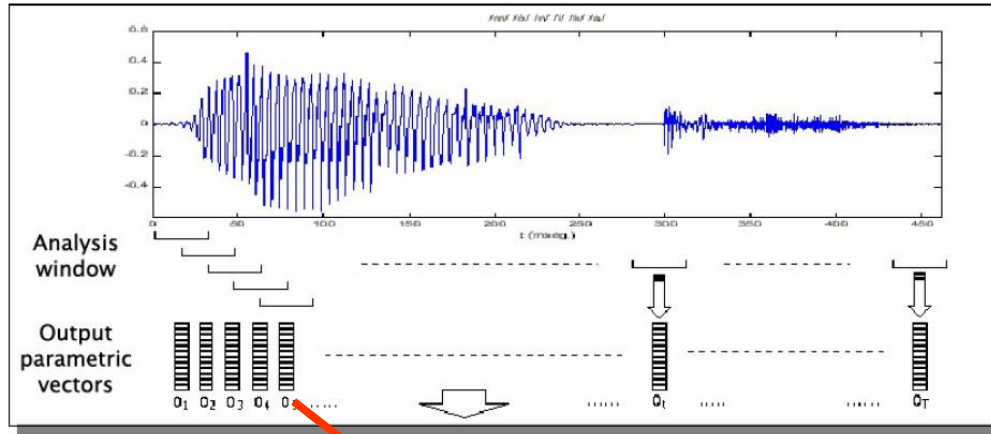
10 Hz    33 Hz    100 Hz

**100 Hz Fundamental**

54

# Short-time Analysis and Parameterization
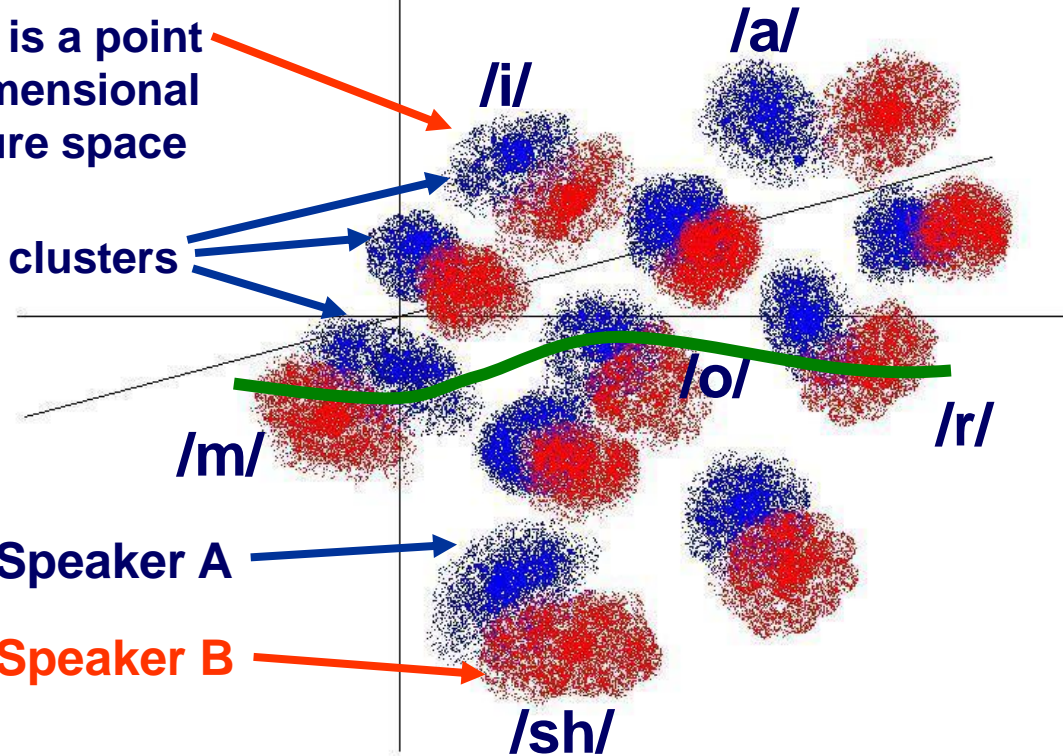
**Feature Space**

Bag of vectors representation
of speaker's acoustic space

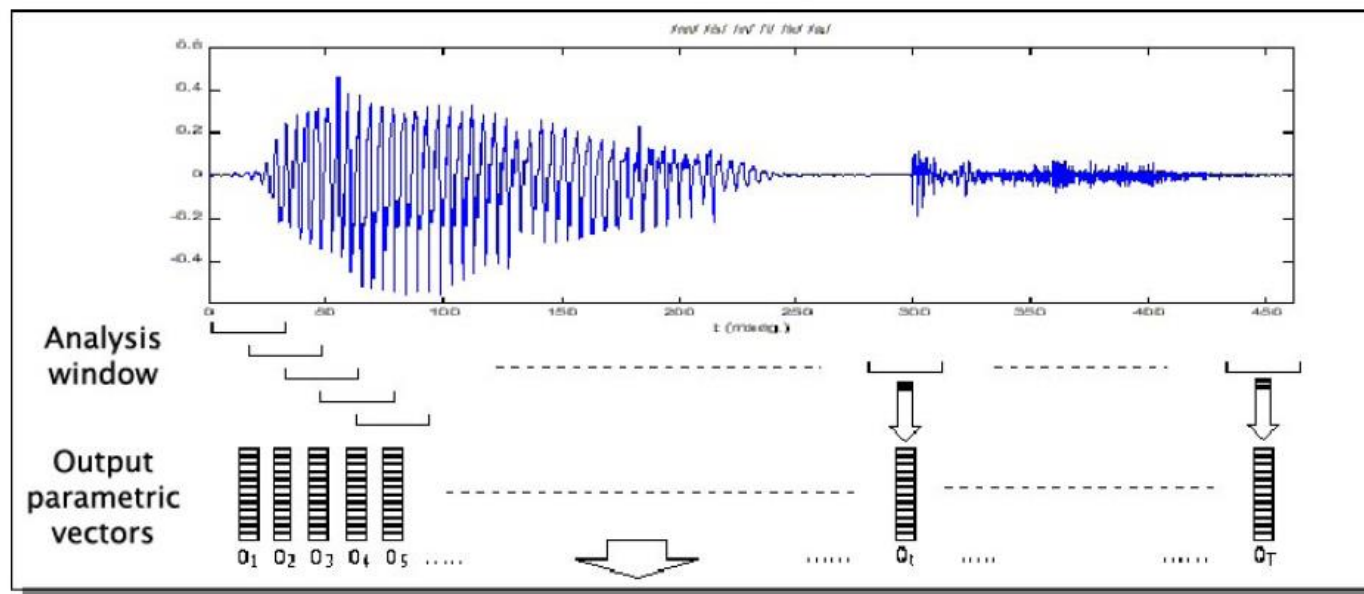Each vector is a point
in the 13-dimensional
MFCC feature space

Phone clusters

Speaker A

Speaker B

/a/

/i/

/o/

/r/

/m/

/sh/

**Feature Space**



Analysis window

Output parametric vectors

$o_1$  $o_2$  $o_3$  $o_4$  $o_5$  .....        ......  $o_t$   .....        ......  $o_T$

**One feature vector every 10 ms**

/s/   /o/   /i/   $o_T$

$o_1$   $o_2$   a   b   $o_n$   c   d

/n/

cd: Phone segment

/s/   /o/   /i/

$o_1$ $o_2$ a b $o_n$ $o_T$

c /n/ d

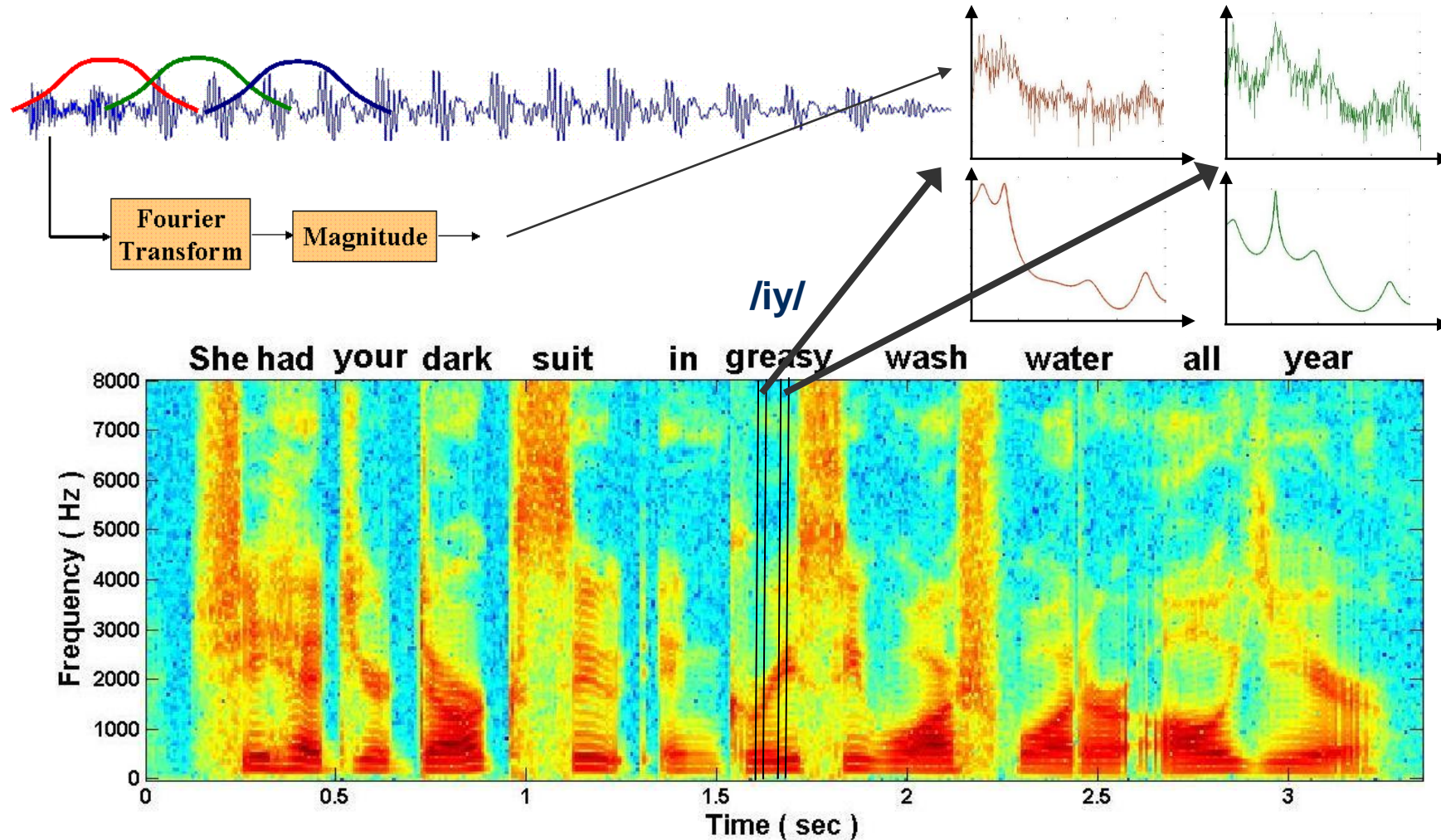cd: Phone segment

**SPEECH RECOGNITION ALGORITHMS**
- ❑ **TAKE THIS FEATURE VECTOR SEQUENCE**
- ❑ **AS INPUT AND DETERMINE "WHAT HAS BEEN SAID"**
- ❑ **e.g. SEQUENCE OF PHONES / SEQUENCE OF WORDS etc.**

# Spectrogram

❑ Speech is a continuous evolution of the vocal tract

❑ Spectrogram shows time-frequency evolution

❑ Represented as a time-series of short-time spectra

Thank you !!