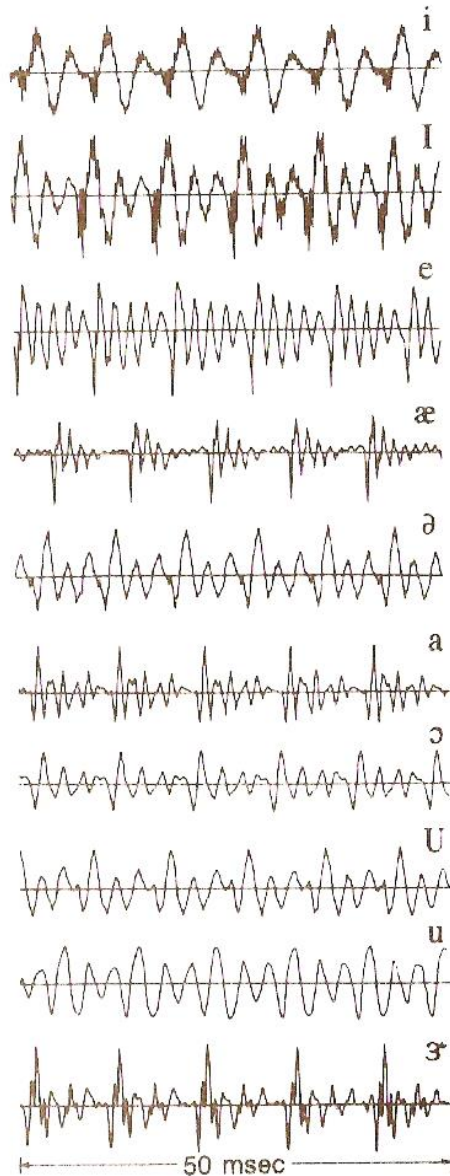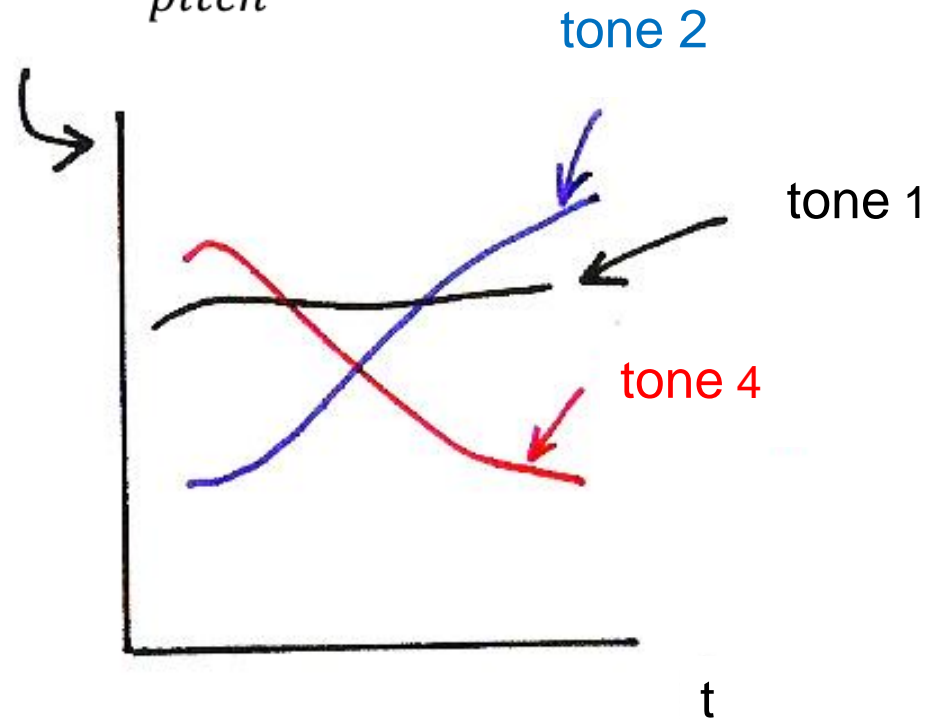# 7.0 Speech Signals and Front-end Processing

**References**: 1.  3.3, 3.4 of Becchetti

3.  9.3 of Huang

# Waveform plots of typical vowel sounds - Voiced（濁音）
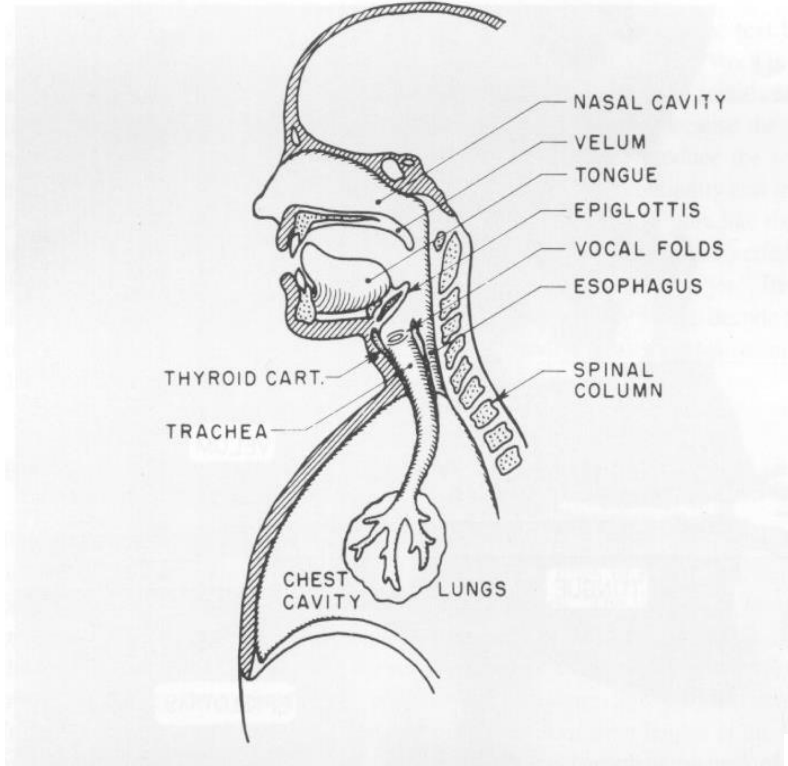


i
I
e
æ
ə
a
ɔ
U
u
ɝ

⊢——— 50 msec ———⊣

$$F_0 = \frac{1}{pitch} \ (音高)$$

tone 2

tone 1

tone 4

t

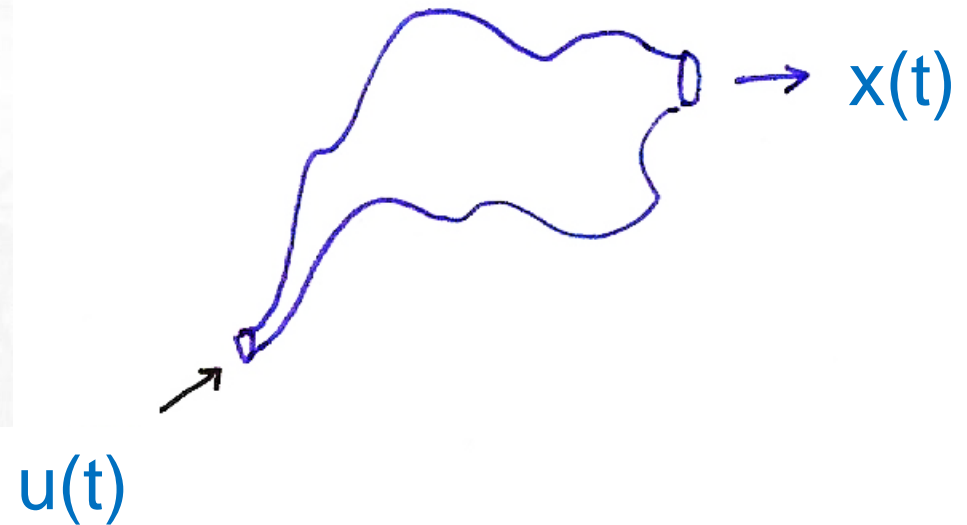# Speech Production and Source Model

• Human vocal mechanism



- NASAL CAVITY
- VELUM
- TONGUE
- EPIGLOTTIS
- VOCAL FOLDS
- ESOPHAGUS
- THYROID CART.
- SPINAL COLUMN
- TRACHEA
- CHEST CAVITY
- LUNGS

• Speech Source Model



$u(t)$   $x(t)$

Vocal tract

excitation

x(t)

u(t)

# **Voiced and Unvoiced Speech**

# Waveform plots of typical consonant sounds

Unvoiced （清音）　　　Voiced （浊音）

# Waveform plot of a sentence



SHOULD WE CHASE

š

U

d                          w

i

č

e$^y$

s

——100 msec——

# Time and Frequency Domains

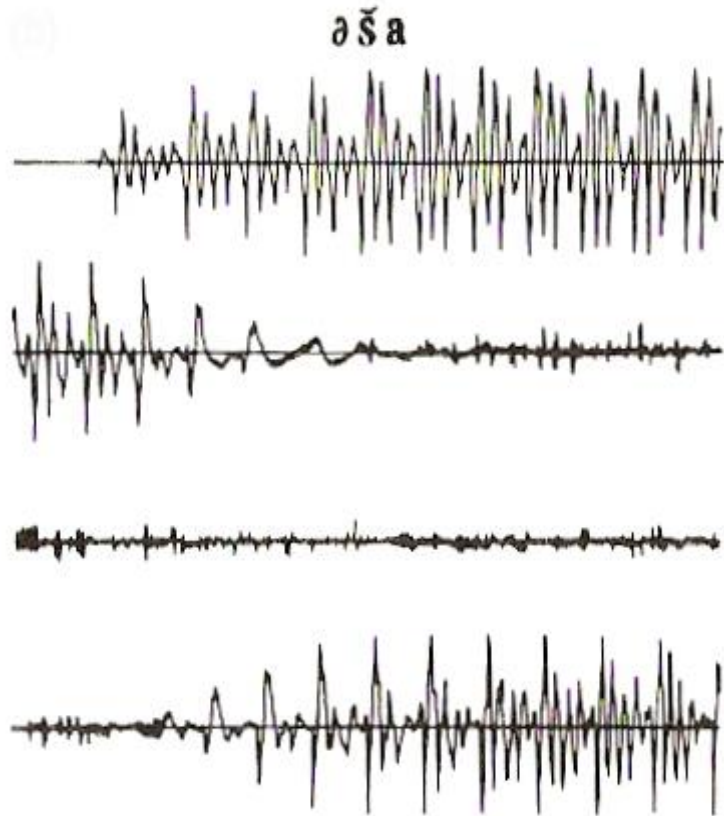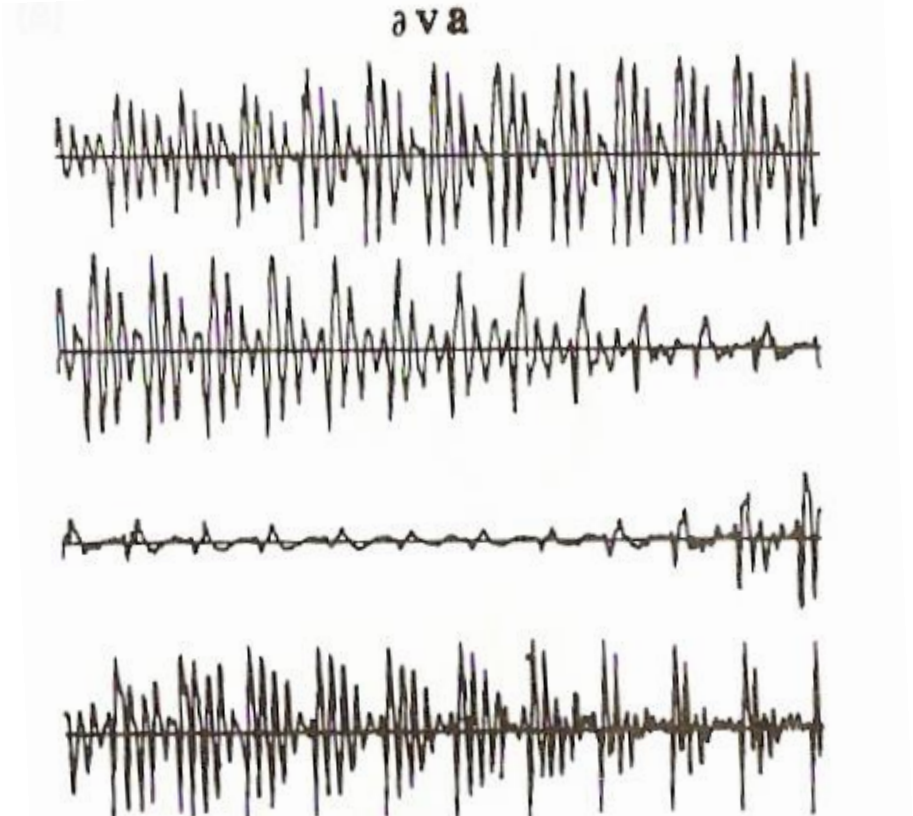$x[n]$  $x(t)$

$t_0$

$x(t_0)$

$t, n$

time domain

$X[k]$

$A_2 e^{j\phi_2}$

$X(\omega)$

$A_1 e^{j\phi_1}$

$\omega, f, k$

$\omega_2$

$\omega_1$

$e^{j\omega_1 t}$

1-1 mapping
Fourier Transform
Fast Fourier Transform (FFT)

frequency domain

Im

$A$ $\phi$

Re

$Re\{e^{j\omega_1 t}\} = \cos(\omega_1 t)$

$Re\{(A_1\, e^{j\phi_1})\, (e^{j\omega_1 t})\} = A_1 \cos(\omega_1 t + \phi_1)$

$\vec{X} = a_1 \vec{i} + a_2 \vec{j} + a_3 \vec{k}$

$\vec{X} = \sum_i a_i\, x_i$

$x(t) = \sum_i a_i\, x_i\,(t)$

$t$

$0$

# Frequency domain spectra of speech signals

Voiced

Unvoiced

# Frequency Domain

Voiced

$$F_0 = \frac{1}{pitch}$$

Formant
Structure

excitation

$F_1$  $F_2$  $F_3$  $F_4$

formant frequencies

Unvoiced

Formant
Structure

excitation

$F_1$  $F_2$  $F_3$  $F_4$

formant frequencies

# Input/Output Relationship for Time/Frequency Domains

$u(t)$

$u[n]$

$\mathscr{F}$

$g(t)$
$g[n]$

$x(t)$

$x[n]$

$\mathscr{F}$

$U(\omega)$

$U[k]$

$\Downarrow \mathscr{F}$

$G(\omega)$

$G[k]$
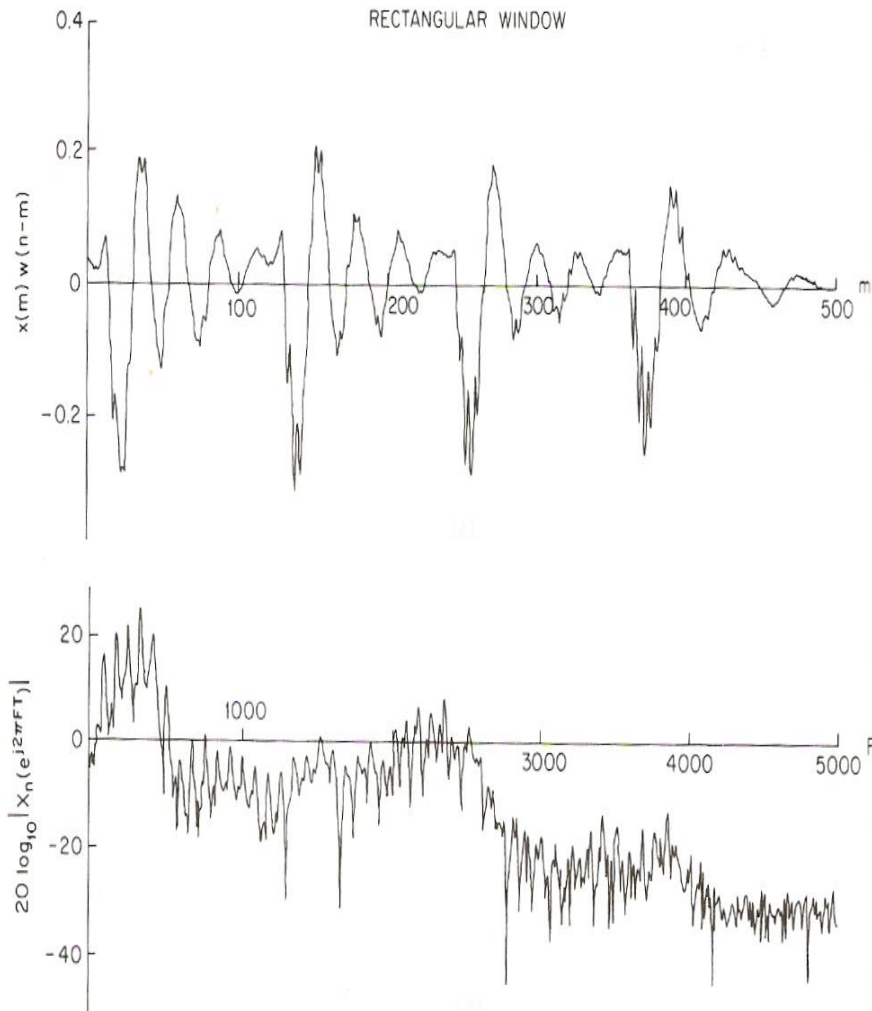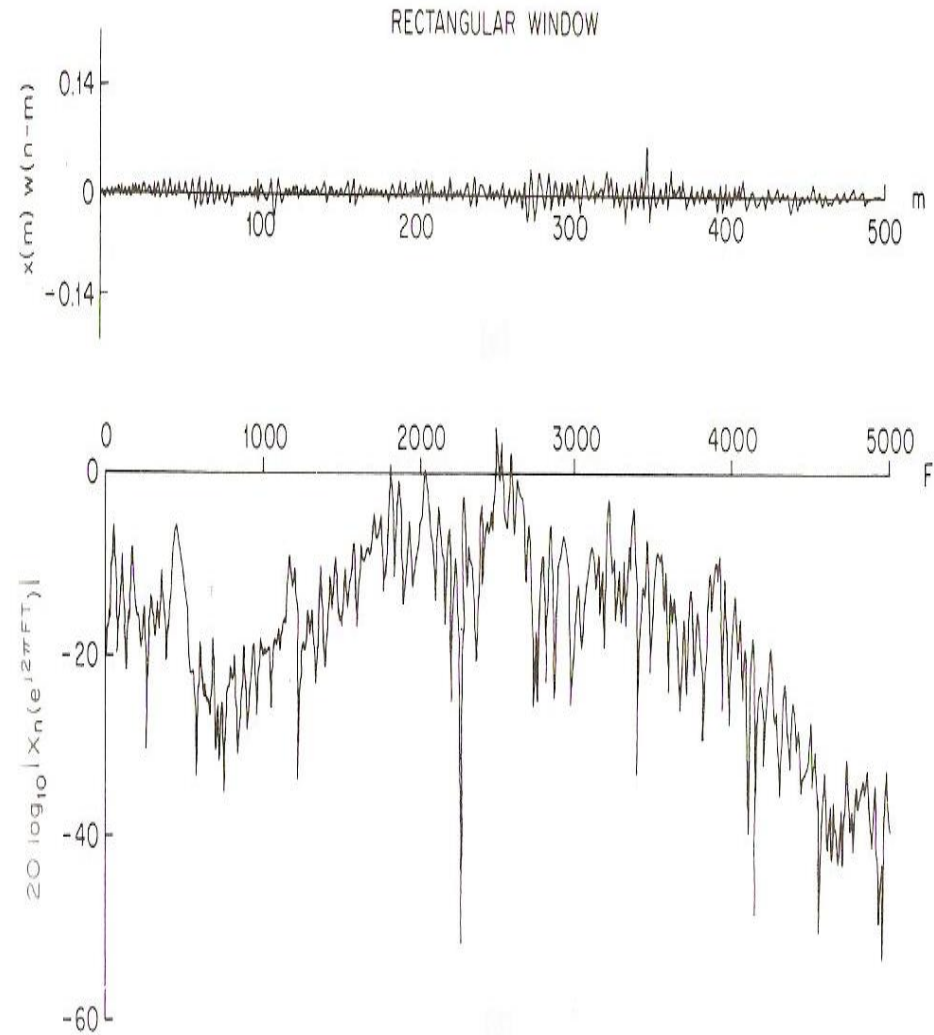
$X(\omega)$

$X[k]$

$\omega, k$

excitation     formant structure

$$x(t) = u(t) * g(t) = \int_{\tau} u(\tau) g(t - \tau) d\tau$$

$$x[n] = u[n] * g[n] = \sum_{k} u[k] g[n - k]$$

time domain: convolution

$$X(\omega) = U(\omega) G(\omega)$$

$$X[k] = U[k] G[k]$$

frequency domain: product

$g(t), G(\omega)$: Formant structure: differences between phonemes

$u(t), U(\omega)$: excitation

# Spectrogram

# Spectrogram



SHOULD WE CHASE

# Formant Frequencies

# Formant frequency contours



He will allow a rare lie.

Reference: 6.1 of Huang, or 2.2, 2.3 of Rabiner and Juang

# Speech Signals

- **Voiced/unvoiced**　　濁音、清音
- **Pitch/tone**　　　　音高、聲調
- **Vocal tract**　　　聲道
- **Frequency domain/formant frequency**
- **Spectrogram representation**
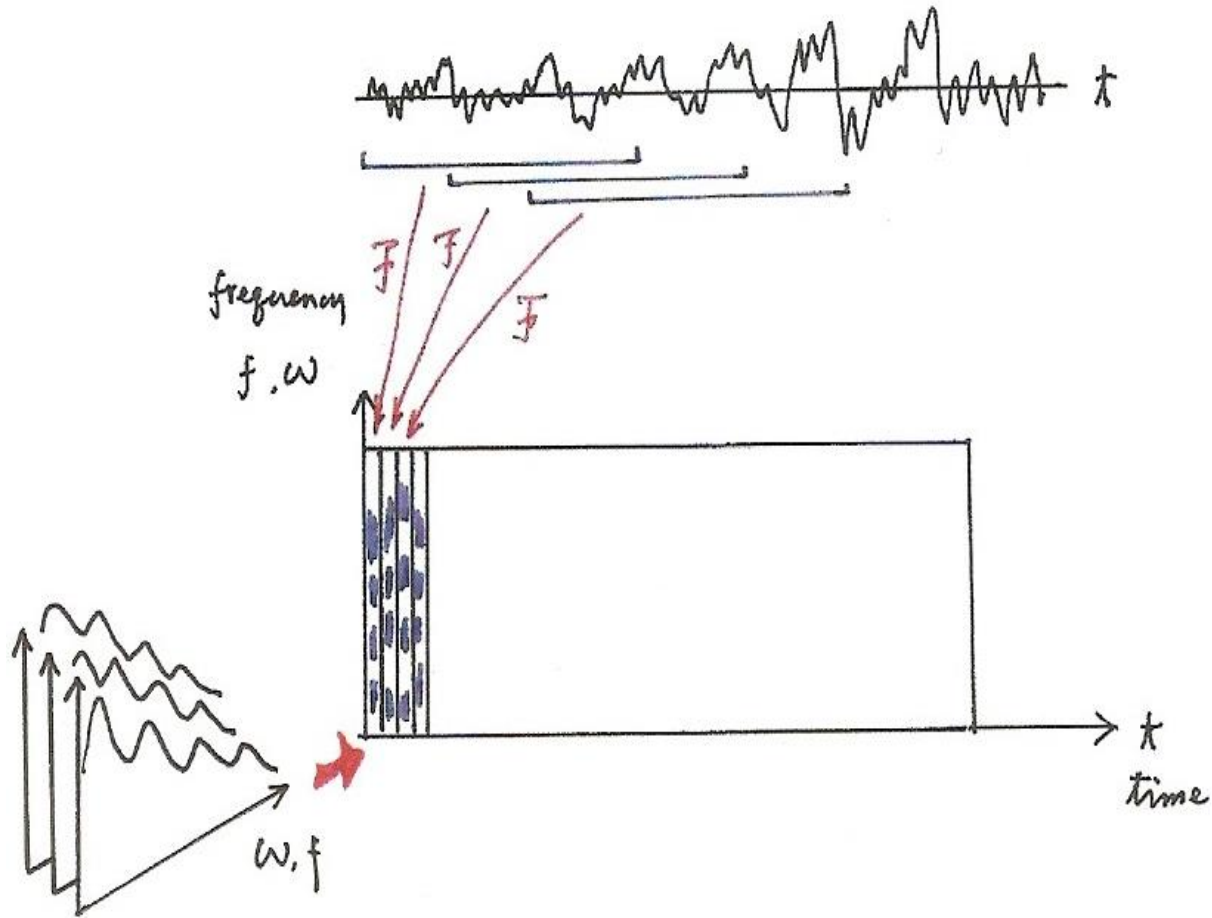- **Speech Source Model**



Ex  　　　　　G(ω),G(z), g[n]

u[n]　　　　　　　　x[n]

Excitation Generator → U(ω) U(z) → Vocal Tract Model → x[n]

x[n]=u[n]*g[n]
X(ω)=U(ω)G(ω)
X(z)=U(z)G(z)

parameters　　　　parameters

- digitization and transmission of the parameters will be adequate
- at receiver the parameters can produce x[n] with the model
- much less parameters with much slower variation in time lead to much less bits required
- the key for low bit rate speech coding

# Speech Source Model

x(t)

t

a[n]

n

0   1   2

# Speech Source Model

- ## Sophisticated model for speech production



- ## Simplified model for speech production

# Simplified Speech Source Model



**Excitation**

**Vocal Tract Model**

$$G(z) = \cfrac{1}{1 - \displaystyle\sum_{k=1}^{P} a_k z^{-k}}$$

– Excitation parameters

   v/u : voiced/ unvoiced

   N : pitch for voiced

   G : signal gain

   $\rightarrow$ excitation signal u[n]

– Vocal Tract parameters

   $\{a_k\}$ : LPC coefficients

   $\rightarrow$formant structure of speech signals

– A good approximation, though not precise enough

Reference: 3.3.1-3.3.6 of Rabiner and Juang, or 6.3 of Huang

# Speech Source Model



$$u[n] \rightarrow \boxed{G(z) = \cfrac{1}{1 - \sum_{k=1}^{P} a_k \, z^{-k}}} \rightarrow x[n]$$

$$x[n] - \sum_{k=1}^{P} a_k \, x[n-k] = u[n]$$

# Feature Extraction - MFCC

- **Mel-Frequency Cepstral Coefficients (MFCC)**
  - Most widely used in the speech recognition
  - Has generally obtained a better accuracy at relatively low computational complexity
  - The process of MFCC extraction :

Speech signal $x(n)$ → **Pre-emphasis** → $x'(n)$ → $\otimes$ → $x_t(n)$ → **DFT** → $X_t(k)$ → **Mel filter-bank**

Window

**Mel filter-bank** → $Y_t(m)$ → **Log(| |$^2$)**

**energy** → $e_t$

**Log(| |$^2$)** → $Y_t'(m)$ → **IDFT** → $y_t(j)$ → **derivatives**

**MFCC**

$$\mathbf{y}_t = \left\{ \begin{array}{c} y_t(j), e_t \\ \Delta\{y_t(j)\}, \Delta\{e_t\} \\ \Delta^2\{y_t(j)\}, \Delta^2\{e_t\} \end{array} \right\}$$

# Pre-emphasis

- The process of Pre-emphasis :
  - a high-pass filter

Speech signal $x(n)$

$H(z)=1-a \cdot z^{-1} \quad 0<a\leq1$

$x'(n)=x(n)-ax(n-1)$



Magnitude (dB)

$a_{pre}=-0.4$

$a_{pre}=-0.95$

$a_{pre}=-1.0$

Frequency (Log - Hz)

$X(\omega)$

$\omega$

# Why pre-emphasis?

- **Reason :**
  - Voiced sections of the speech signal naturally have a negative spectral slope (attenuation) of approximately 20 dB per decade due to the physiological characteristics of the speech production system
  - High frequency formants have small amplitude with respect to low frequency formants. A pre-emphasis of high frequencies is therefore helpful to obtain similar amplitude for all formants

# Why Windowing?

- **Why dividing the speech signal into successive and overlapping frames?**
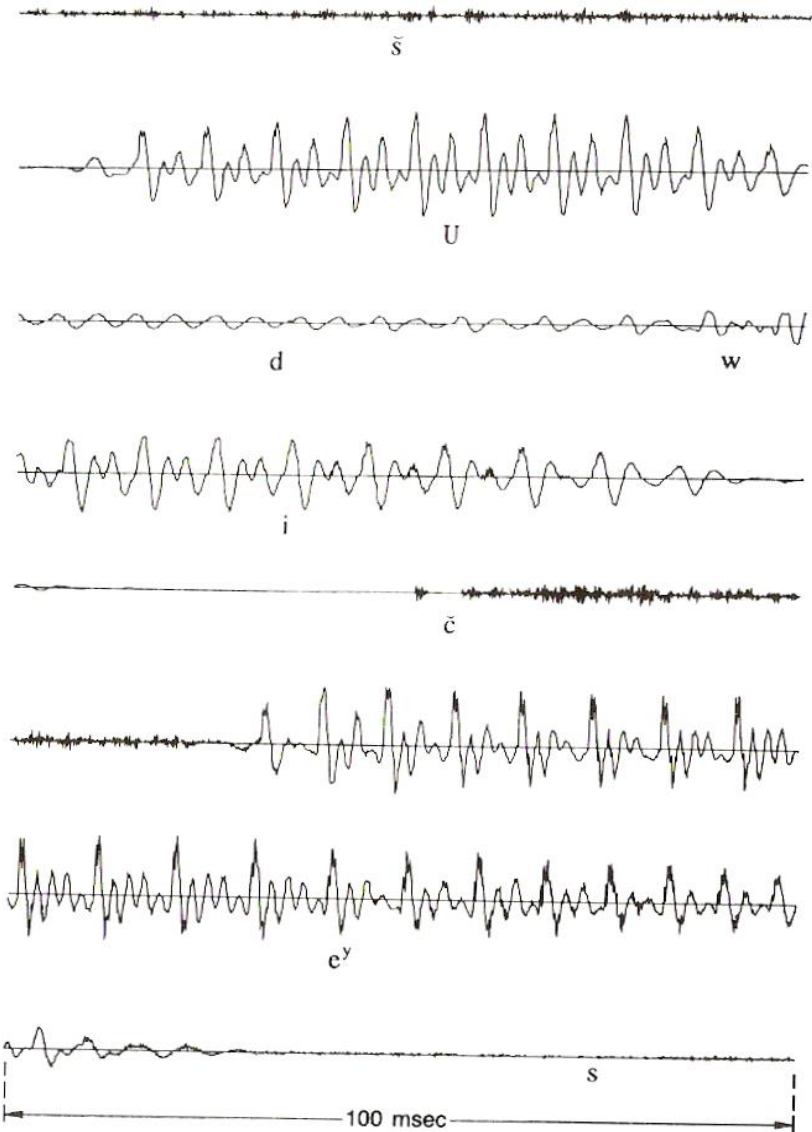  - Voice signals change their characteristics from time to time. The characteristics remain unchanged only in short time intervals (short-time stationary, short-time Fourier transform)
- **Frames**
  - **Frame Length :** the length of time over which a set of parameters can be obtained and is valid. Frame length ranges between **20 ~ 10** ms
  - **Frame Shift:** the length of time between successive parameter calculations
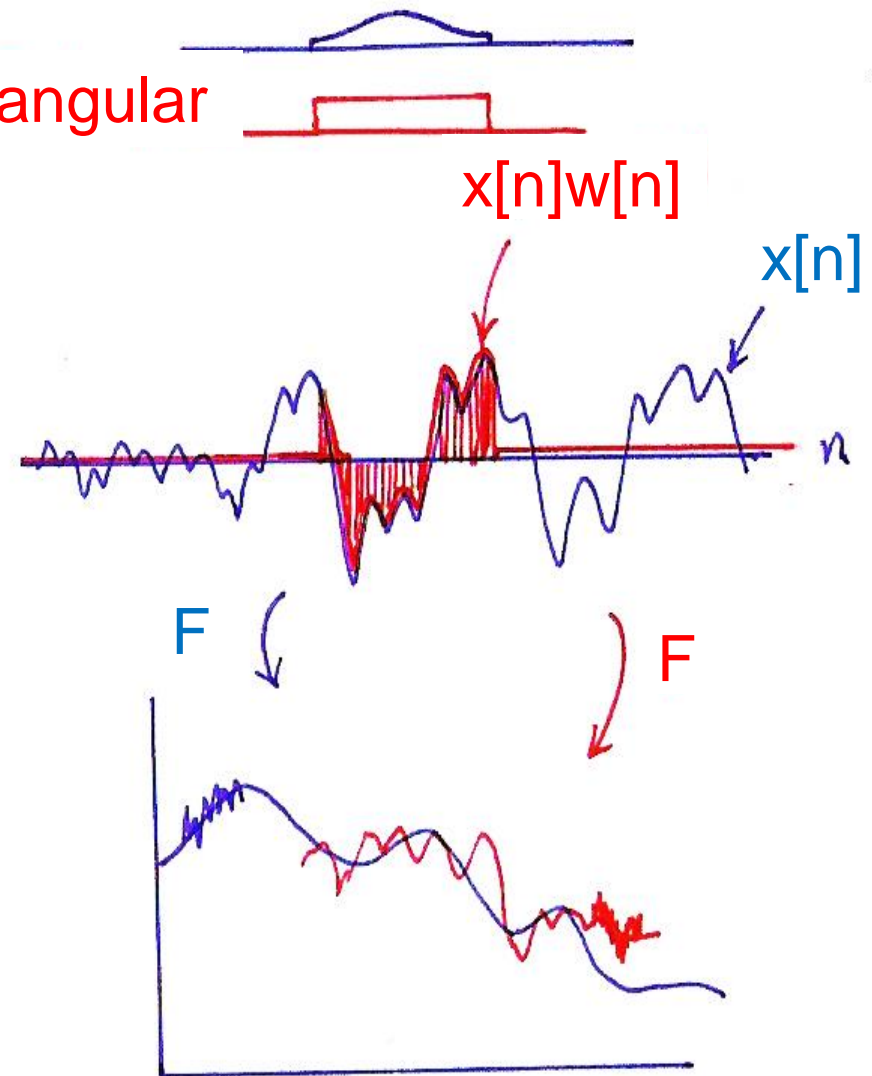  - **Frame Rate:** number of frames per second
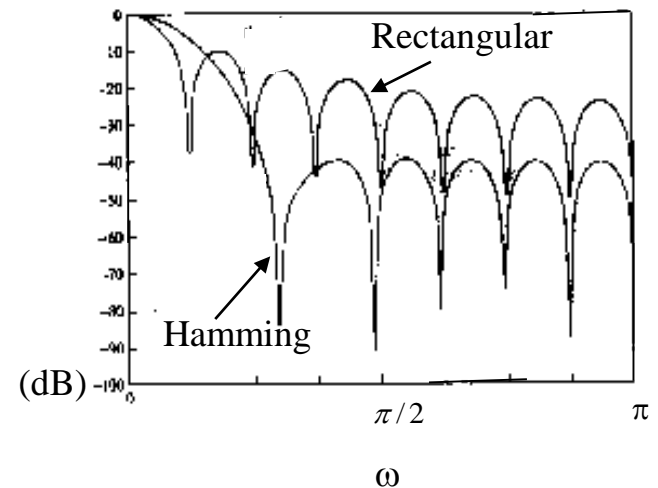
# Waveform plot of a sentence



SHOULD WE CHASE

š

U

d                                    w

i

č

eʸ

s

|←——————— 100 msec ———————→|

Hamming

Rectangular

x[n]w[n]

x[n]

n

F

F

# Effect of Windowing (1)

- **Windowing :**
  - $x_t(n)=w(n) \bullet x'(n)$, $w(n)$: the shape of the window (product in time domain)
    - $X_t(\omega)=W(\omega)*X'(\omega)$, *: convolution (convolution in frequency domain)
  - Rectangular window ($w(n)=1$ for $0 \le n \le L\text{-}1$):
    - simply extract a segment of the signal
    - whose frequency response has high side lobes
  - *Main lobe* : spreads out the narrow band power of the signal (that around the formant frequency) in a wider frequency range, and thus reduces the local frequency resolution in formant allocation
  - *Side lobe* : swap energy from different and distant frequencies

excitation    formant structure

$$x(t) = u(t) * g(t) = \int_{\tau} u(\tau) g(t-\tau) d\tau$$

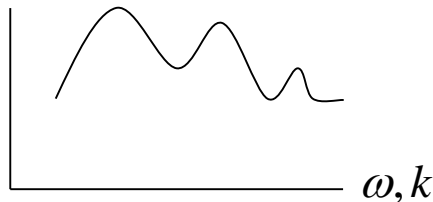$$x[n] = u[n] * g[n] = \sum_{k} u[k] g[n-k]$$

time domain: convolution

$$X(\omega) = U(\omega) G(\omega)$$

$$X[k] = U[k] G[k]$$

frequency domain: product

$u(t)$
$u[n]$

$g(t)$
$g[n]$

$x(t)$
$x[n]$

$\mathcal{F}$

$U(\omega)$
$U[k]$

$\mathcal{F}$

$G(\omega)$
$G[k]$

$X(\omega)$
$X[k]$

$\mathcal{F}$

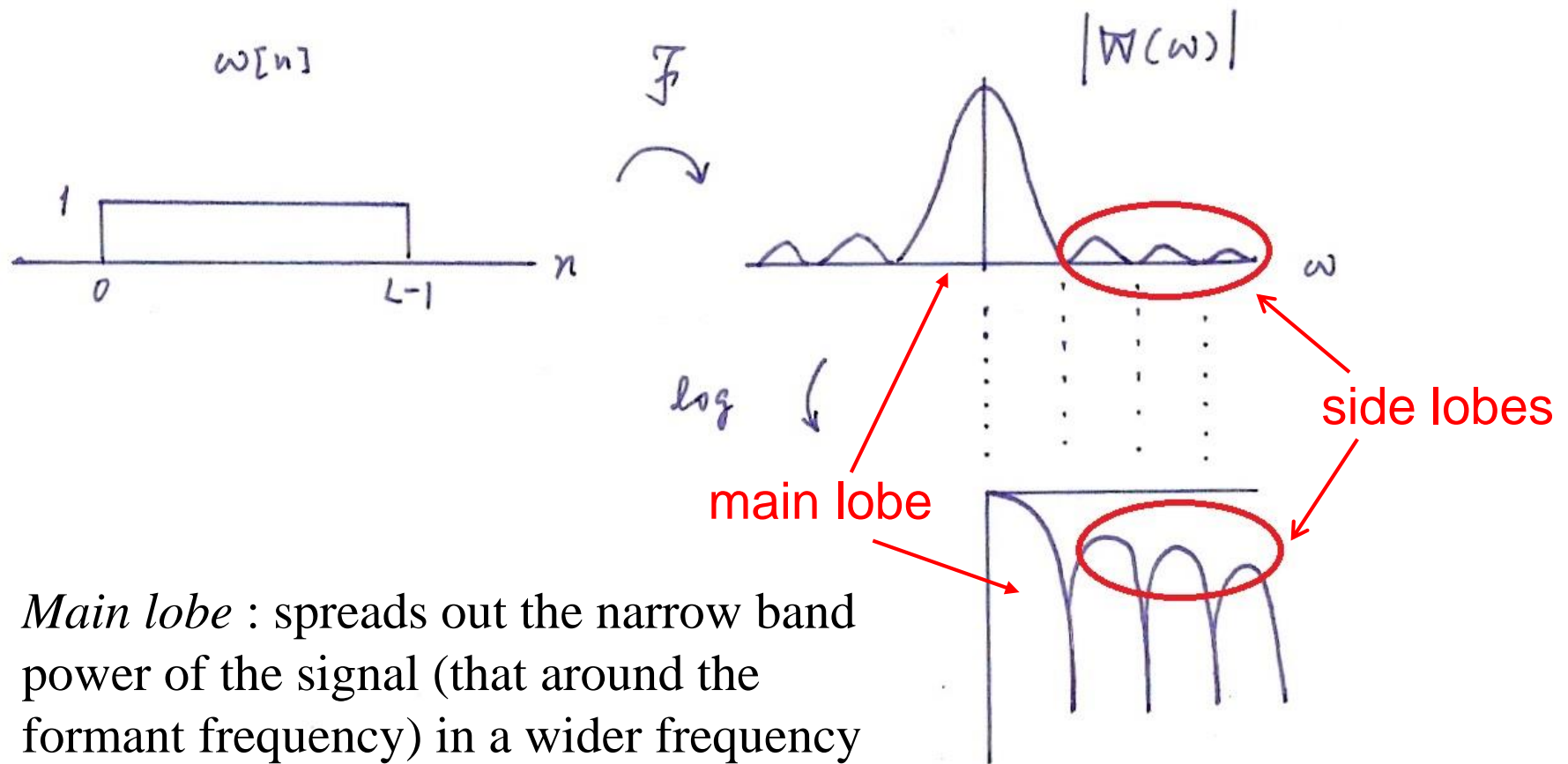$\omega, k$

$g(t), G(\omega)$: Formant structure: differences between phonemes

$u(t), U(\omega)$: excitation

# **Windowing**



- *Main lobe* : spreads out the narrow band power of the signal (that around the formant frequency) in a wider frequency range, and thus reduces the local frequency resolution in formant allocation
- *Side lobe* : swap energy from different and distant frequencies
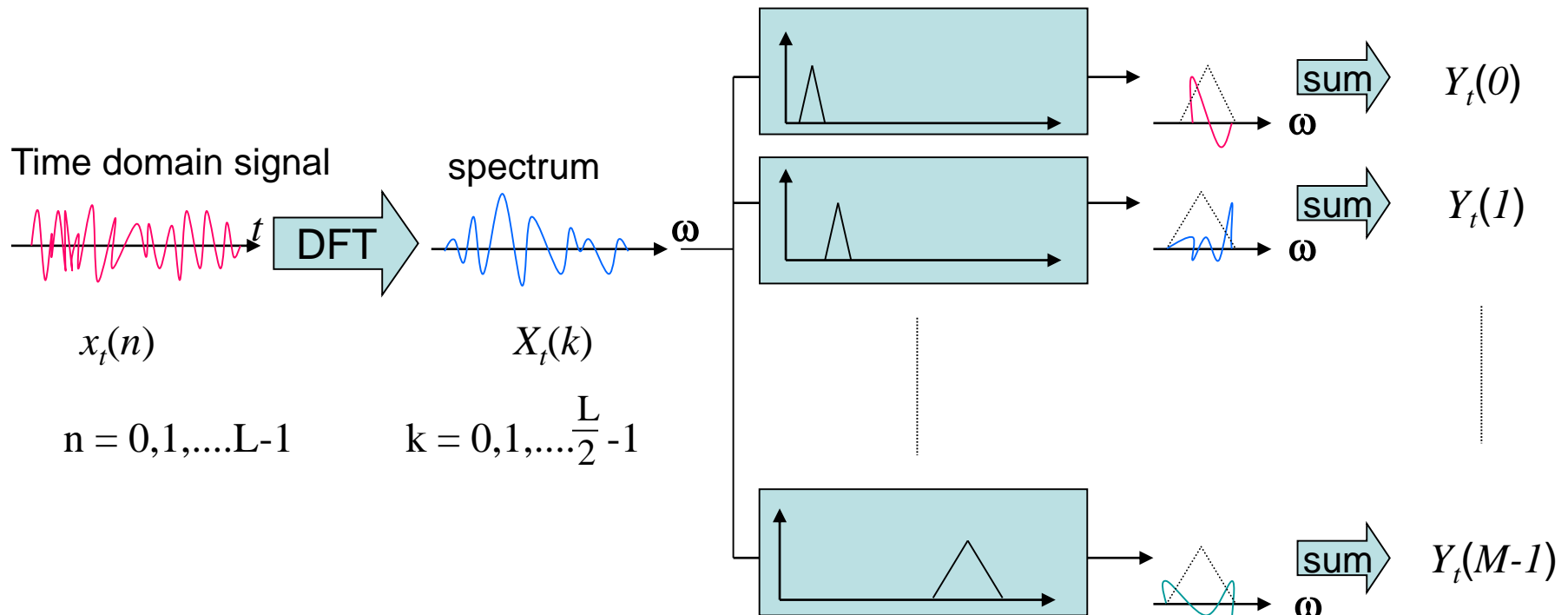
# Effect of Windowing (2)

- **Windowing (Cont.):**
  - For a designed window, we wish that
    - the main lobe is as narrow as possible
    - the side lobe is as low as possible
      - However, it is impossible to achieve both simultaneously. Some trade-off is needed
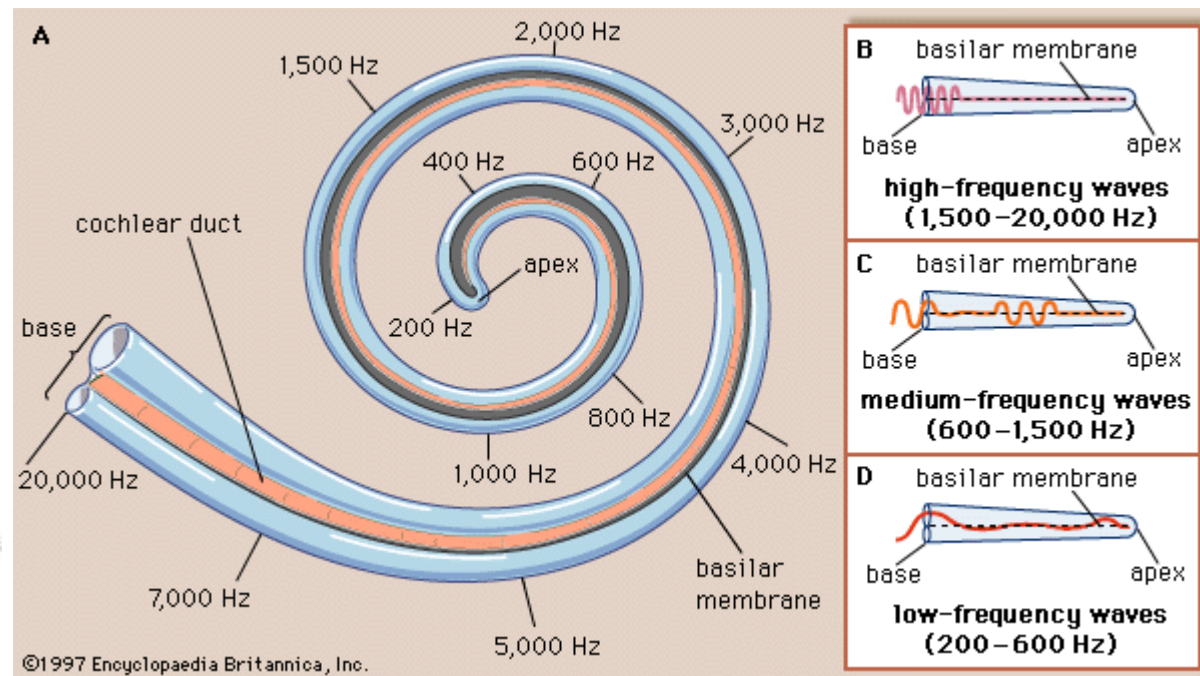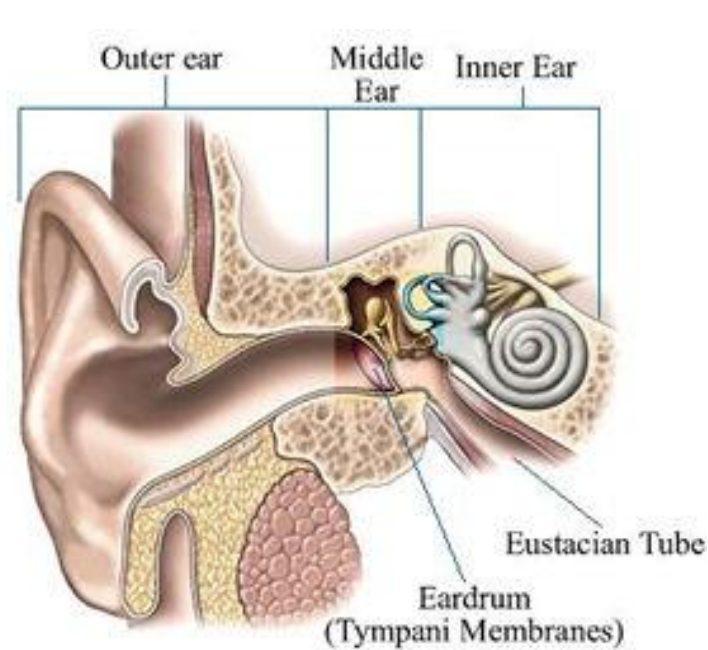  - The most widely used window shape is the Hamming window

$$w(n) = \begin{cases} 0.54 - 0.46\cos\left(\dfrac{2\pi n}{L-1}\right), & n = 0,1,......,L-1 \\ 0 & \text{otherwise} \end{cases}$$
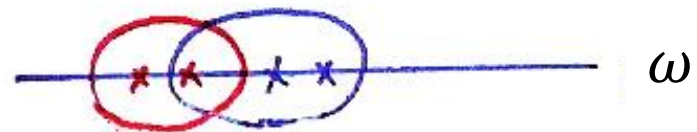
# DFT and Mel-filter-bank Processing

- **For each frame of signal (*L* points, e.g., L=512),**
  - the Discrete Fourier Transform (DFT) is first performed to obtain its spectrum (*L* points, for example *L*=512)
  - The bank of filters based on Mel scale is then applied, and each filter output is the sum of its filtered spectral components (*M* filters, and thus *M* outputs, for example *M*=24)

Time domain signal $\;\xrightarrow{\text{DFT}}\;$ spectrum

$x_t(n)$   $X_t(k)$

$n = 0,1,....L-1$   $k = 0,1,....\dfrac{L}{2}-1$

sum $\;\;\; Y_t(0)$

sum $\;\;\; Y_t(1)$

sum $\;\;\; Y_t(M-1)$

# Peripheral Processing for Human Perception

# Mel-scale Filter Bank



$\omega$

$\log \omega$

$X(\omega)$

$\omega$

$\omega$

# Why Filter-bank Processing?

- **The filter-bank processing simulates human ear perception**
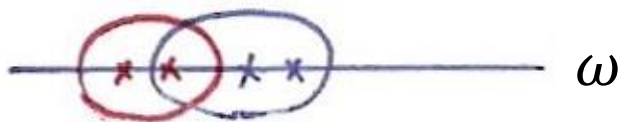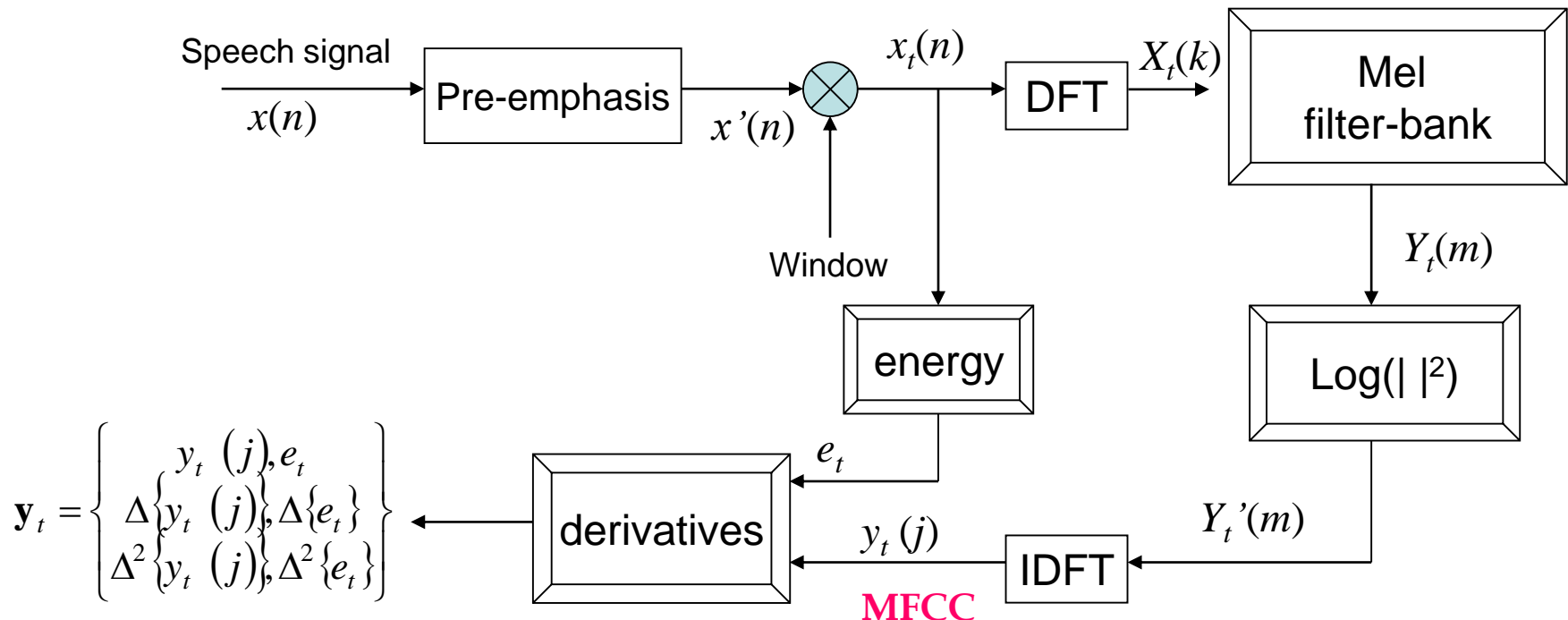  - Frequencies of a complex sound within a certain frequency band cannot be individually identified.
  - When one of the components of this sound falls outside this frequency band, it can be individually distinguished.
  - This frequency band is referred to as the critical band.
  - These critical bands somehow overlap with each other.
  - The critical bands are roughly distributed linearly in the logarithm frequency scale (including the center frequencies and the bandwidths), specially at higher frequencies.
  - Human perception for pitch of signals is proportional to the *logarithm* of the frequencies (relative ratios between the frequencies)
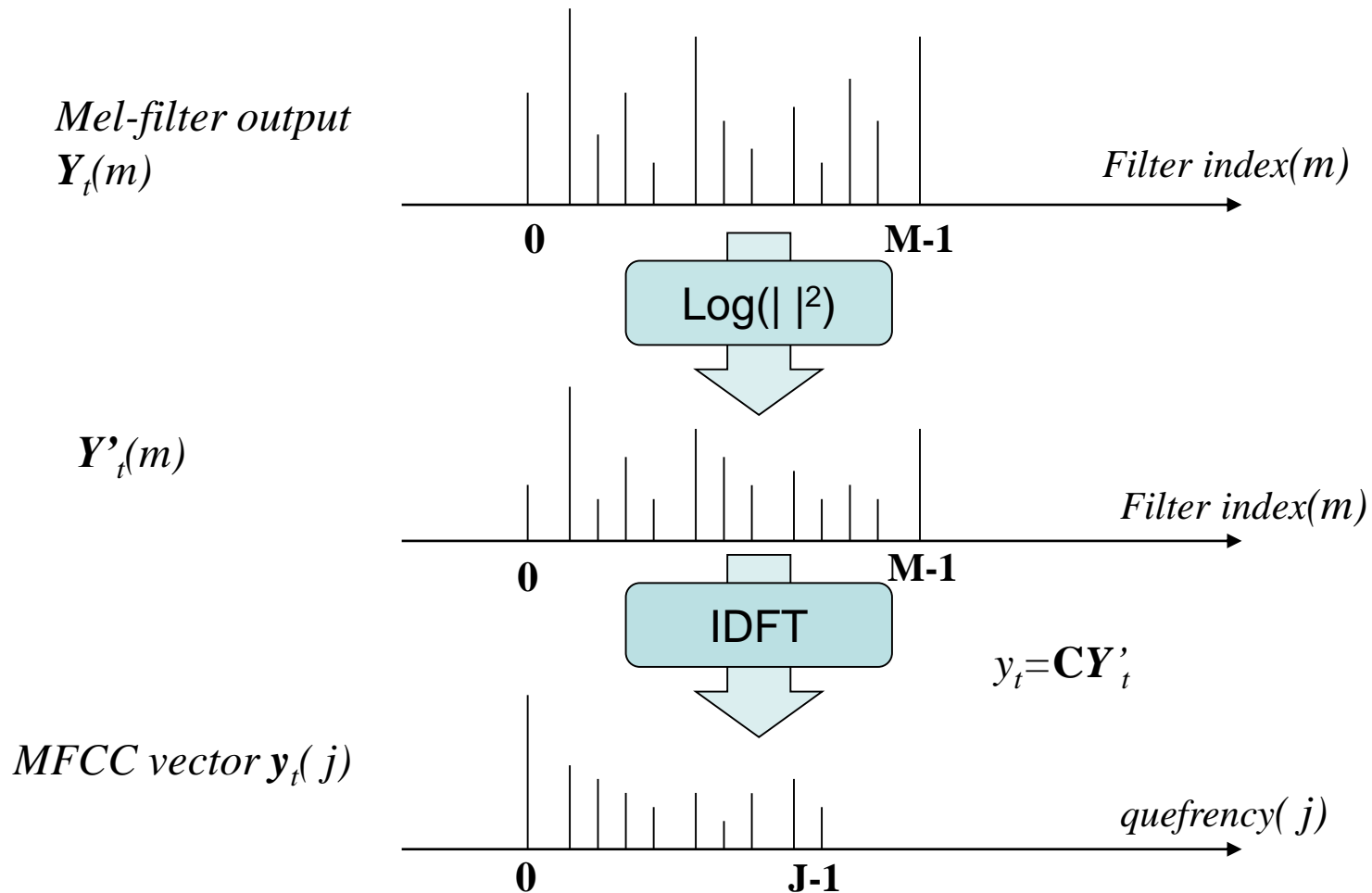
$\omega$

$\log \omega$

# Feature Extraction - MFCC

- **Mel-Frequency Cepstral Coefficients (MFCC)**
  - Most widely used in the speech recognition
  - Has generally obtained a better accuracy at relatively low computational complexity
  - The process of MFCC extraction :

# Logarithmic Operation and IDFT

- **The final process of MFCC evaluation : logarithm operation and IDFT**

*Mel-filter output* $Y_t(m)$

*Filter index(m)*

0    M-1

Log(| |²)

$Y'_t(m)$

*Filter index(m)*

0    M-1

IDFT

$y_t = \mathbf{C}Y'_t$

*MFCC vector* $y_t(j)$

*quefrency( j)*

0    J-1

# Why Log Energy Computation?

- **Using the magnitude (or energy) only**
  - Phase information is not very helpful in speech recognition
    - Replacing the phase part of the original speech signal with continuous random phase usually won't be perceived by human ears

- **Using the Logarithmic operation**
  - Human perception sensitivity is proportional to signal energy in logarithmic scale (relative ratios between signal energy values)
  - The logarithm compresses larger values while expands smaller values, which is a characteristic of the human hearing system
  - The dynamic compression also makes feature extraction less sensitive to variations in signal dynamics
  - To make a convolved noisy process additive
    - Speech signal $x(n)$, excitation $u(n)$ and the impulse response of vocal tract $g(n)$

$$x(n)=u(n)*g(n) \Rightarrow X(\omega)=U(\omega)G(\omega)$$
$$\Rightarrow |X(\omega)|=|U(\omega)||G(\omega)| \Rightarrow \log|X(\omega)|=\log|U(\omega)|+\log|G(\omega)|$$

# Why Inverse DFT?

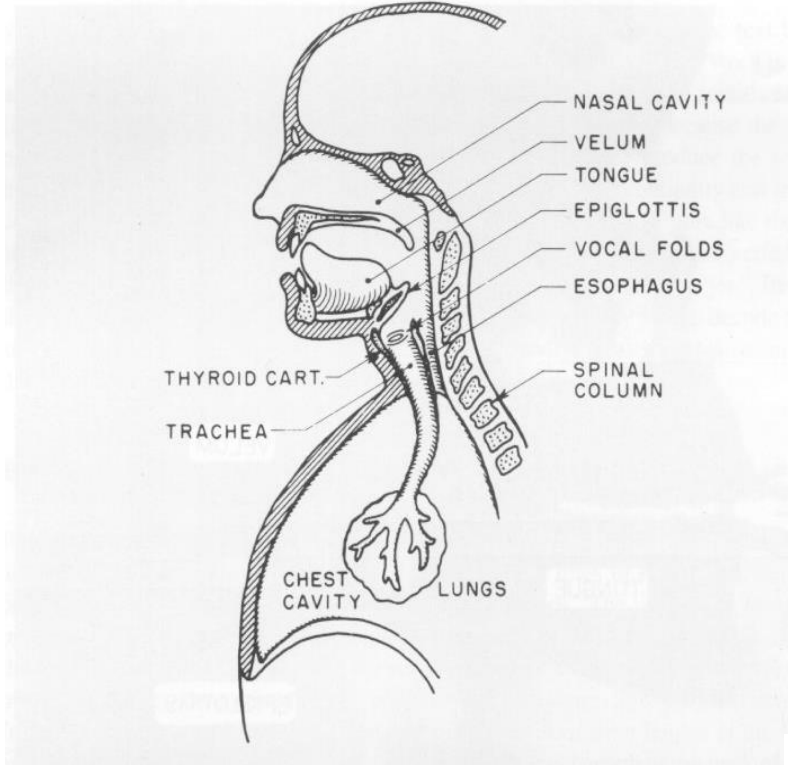- **Final procedure for MFCC : performing the inverse DFT on the log-spectral power**

$$y_t(j) = \sum_{m=0}^{M-1} \log\left(|Y_t(m)|^2\right)\cos\left[j\left(m-\frac{1}{2}\right)\frac{\pi}{M}\right], \quad j = 0,1,....,J-1 < M$$
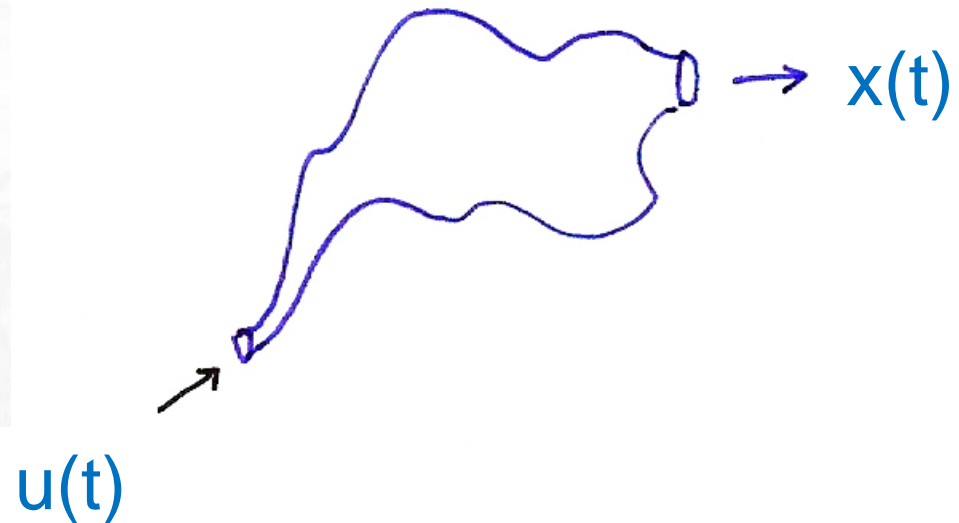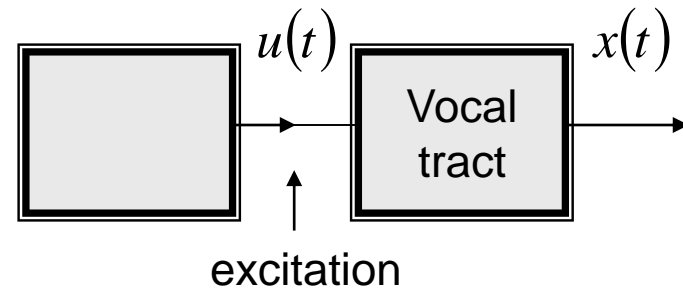
- **Advantages :**
  - Since the log-power spectrum is real and symmetric, the inverse DFT reduces to a Discrete Cosine Transform (DCT). The DCT has the property to produce highly uncorrelated features $y_t$
    - diagonal rather than full covariance matrices can be used in the Gaussian distributions in many cases
  - Easier to remove the interference of excitation on formant structures
    - the phoneme for a segment of speech signal is primarily based on the formant structure (or vocal tract shape)
    - on the frequency scale the formant structure changes slowly over frequency, while the excitation changes much faster

• Human vocal  mechanism

• Speech Source Model

NASAL CAVITY
VELUM
TONGUE
EPIGLOTTIS
VOCAL FOLDS
ESOPHAGUS
THYROID CART.
SPINAL COLUMN
TRACHEA
CHEST CAVITY
LUNGS

$u(t)$   $x(t)$

Vocal tract

excitation

x(t)

u(t)

Voiced

Unvoiced

excitation    formant structure

$$x(t) = u(t) * g(t) = \int_{\tau} u(\tau)g(t-\tau)d\tau$$

$$x[n] = u[n] * g[n] = \sum_{k} u[k]g[n-k]$$

time domain: convolution

$u(t)$
$u[n]$    $g(t)$    $x(t)$
         $g[n]$    $x[n]$

$\mathcal{F}$

$U(\omega)$    $\mathcal{F}$    $X(\omega)$
$U[k]$    $G(\omega)$    $X[k]$
         $G[k]$

$\omega, k$

$$X(\omega) = U(\omega)G(\omega)$$

$$X[k] = U[k]G[k]$$

frequency domain: product

$g(t), G(\omega)$: Formant structure: differences between phonemes

$u(t), U(\omega)$: excitation

# Logarithmic Operation

$u[n]$ $\rightarrow$ $\boxed{g[n]}$ $\rightarrow$ $x[n]= u[n]*g[n]$

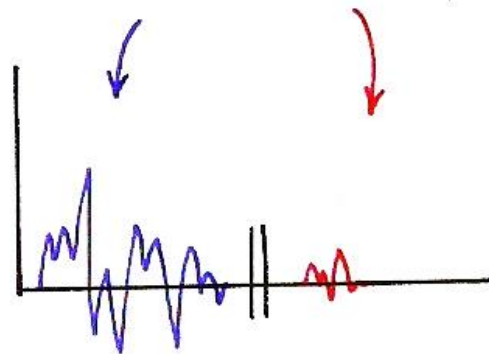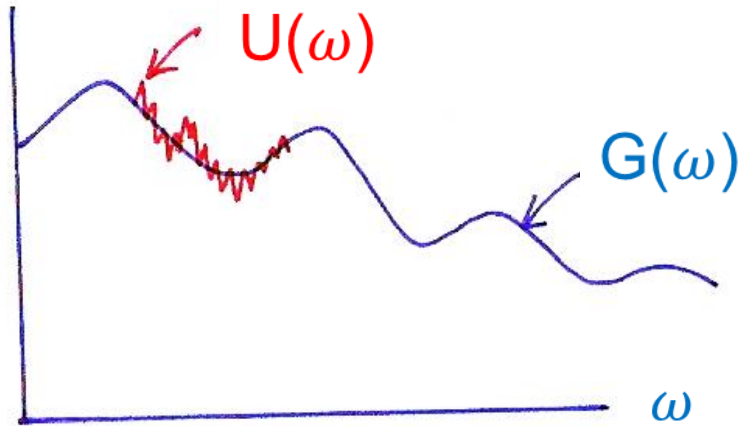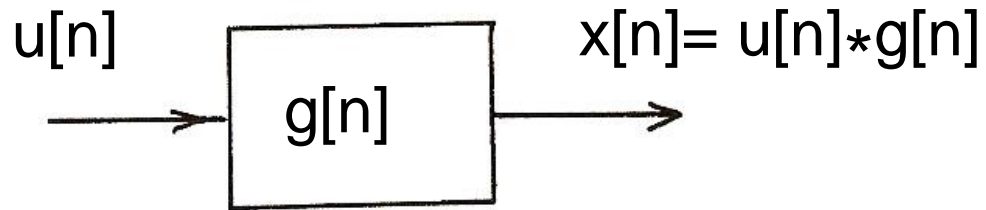$U(\omega)$

$G(\omega)$

$\omega$

# Derivatives

- **Derivative operation : to obtain the change of the feature vectors with time**



MFCC stream $y_t(j)$

$$\Delta y_t(j) = \frac{\sum\limits_{m=-p}^{p} m \bullet y_{t-m}(j)}{\sum\limits_{m=-p}^{p} m^2}$$

$\Delta$MFCC stream $\Delta y_t(j)$

$$\Delta^2 y_t(j) = \frac{\sum\limits_{m=-p}^{p} m \bullet \Delta y_{t-m}(j)}{\sum\limits_{m=-p}^{p} m^2}$$

$\Delta^2$ MFCC stream $\Delta^2 y_t(j)$

# Linear Regression

$(x_i, y_i)$

$y = ax + b$

$$\sum_i \left(ax_i + b - y_i\right)^2 = \min$$

find a, b

# Why Delta Coefficients?

- **To capture the dynamic characters of the speech signal**
  - Such information carries relevant information for speech recognition
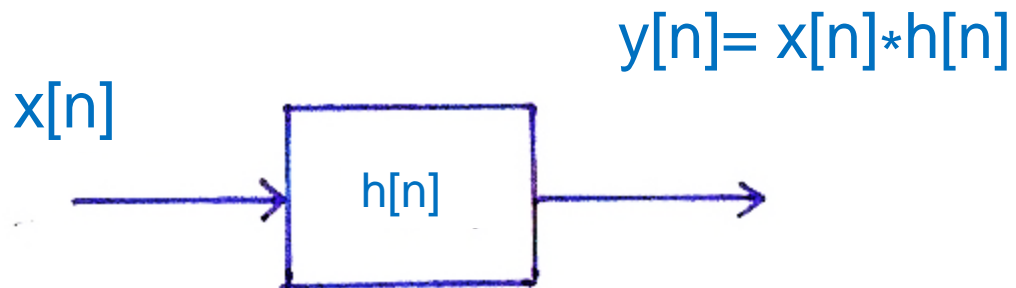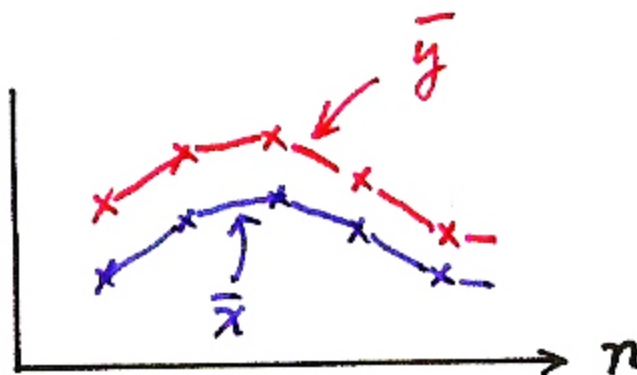  - The value of $p$ should be properly chosen
    - The dynamic characters may not be properly extracted if $p$ is too small
    - Too large p may imply frames too far away
- **To cancel the DC part (channel distortion or convolutional noise) of the MFCC features**
  - Assume, for clean speech, an MFCC parameter stream for an utterance is $\{\mathbf{y}(t\text{-}N), \mathbf{y}(t\text{-}N\text{+}1),\text{…......,}\mathbf{y}(t), \mathbf{y}(t\text{+}1), \mathbf{y}(t\text{+}2), \text{……}\}$,

    $\mathbf{y}(t)$ is an MFCC parameter at time $t$,
    while after channel distortion, the MFCC stream becomes
    $\{\mathbf{y}(t\text{-}N)\text{+}h, \mathbf{y}(t\text{-}N\text{+}1)\text{+}h,\text{…......,}\mathbf{y}(t)\text{+}h, \mathbf{y}(t\text{+}1)\text{+}h, \mathbf{y}(t\text{+}2)\text{+}h, \text{……}\}$
    the channel effect $h$ is eliminated in the delta (difference) coefficients

# Convolutional Noise

$$y[n] = x[n] * h[n]$$
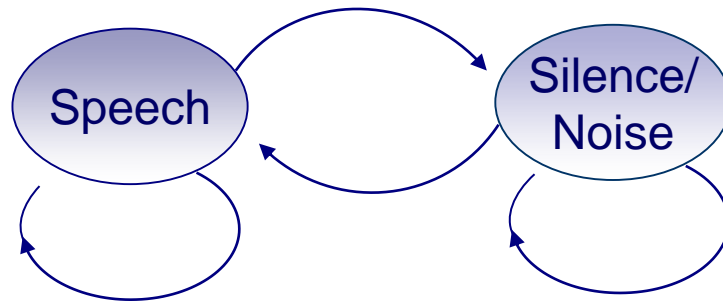
x[n]

h[n]

MFCC

$$\overline{y} = \overline{x} + \overline{h}$$

# End-point Detection

- **Push (and Hold) to Talk/Continuously Listening**
- **Adaptive Energy Threshold**
- **Low Rejection Rate**
  - false acceptance may be rescued
- **Vocabulary Words Preceded and Followed by a Silence/Noise Model**
- **Two-class Pattern Classifier**



  - Gaussian density functions used to model the two classes
  - log-energy, delta log-energy as the feature parameters
  - dynamically adapted parameters

# End-point Detection