

Công nghệ nhận dạng và tổng hợp tiếng nói

Trịnh Văn Loan – ĐHBK Hà Nội

1

Tài liệu tham khảo

- ▶ An Introduction to Digital Speech Processing
Lawrence R. Rabiner, Ronald W. Schafer, Now, 2007
- ▶ Digital Processing of Speech Signals
Lawrence R. Rabiner, Ronald W. Schafer, Prentice-Hall, 1978
- ▶ Discrete-Time Processing of Speech Signals
John R. Deller, John G. Proakis, Hansen John H. L., IEEE Press, 2000
- ▶ Fundamentals of Speech Recognition
Lawrence Rabiner, Bing-Hwang Juang, Pearson College Div, 1993
- ▶ Automatic Speech Recognition: A Deep Learning Approach (Signals and Communication Technology)
Dong Yu, Li Deng, Springer, 2015
- ▶ Text-to-Speech Synthesis
Paul Taylor, Cambridge University Press, 2009
- ▶ Improvements of Vietnamese Hidden Markov Model based speech synthesis
Duy Khanh Ninh, LAP LAMBERT Academic Publishing, 2020
- ▶ *Nguyễn Hữu Quỳnh, Hà Nội, 1994*
- ▶ *Dẫn luận Ngôn ngữ học*
Nguyễn Thiện Giáp, Đoàn Thiện Thuật, Nguyễn Minh Thuyết, Hà Nội, 1994

2

1. Các khái niệm cơ bản

- ▶ Xử lý tiếng nói là gì ?
- ▶ Xử lý tiếng nói bao hàm các lĩnh vực:
 - ▶ Nhận dạng tiếng nói
 - ▶ Nhận dạng người nói
 - ▶ Mã hóa và giải mã tiếng nói
 - ▶ Tổng hợp tiếng nói
 - ▶ Tăng cường chất lượng tín hiệu tiếng nói

3

- ▶ Các ứng dụng của xử lý tiếng nói
 - ▶ Tương tác người - máy
 - ▶ Viễn thông
 - ▶ Các công nghệ trợ giúp (khuyết tật, khiếm thị, học ngôn ngữ)
 - ▶ Khai thác dữ liệu âm thanh
 - ▶ An ninh, bảo mật
- ▶ Các lĩnh vực khoa học liên quan
 - ▶ Xử lý tín hiệu số
 - ▶ Xử lý ngôn ngữ tự nhiên
 - ▶ Học máy
 - ▶ Ngữ âm học
 - ▶ Tương tác người máy
 - ▶ Tâm lý học cảm thụ

4

1.2 Các khái niệm cơ bản về tiếng nói

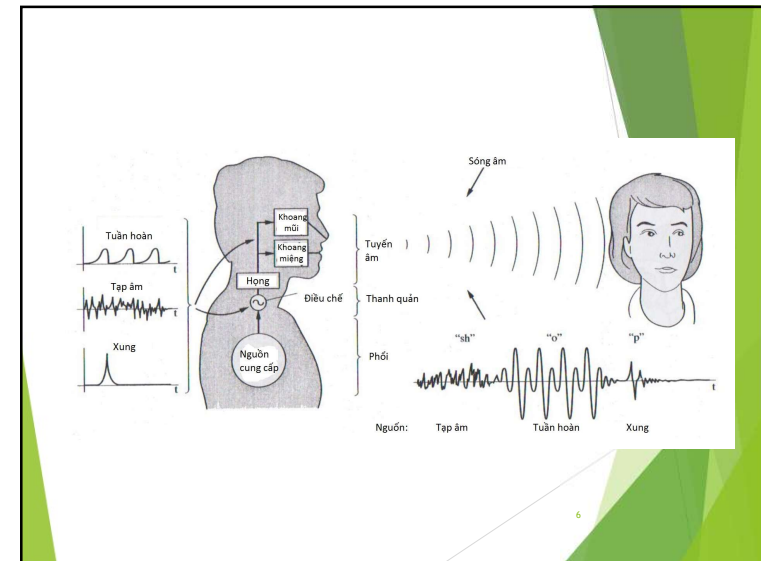
► Phân biệt tiếng nói và âm thanh

Tiếng nói được phân biệt với các âm thanh khác bởi các đặc tính âm học có nguồn gốc từ cơ chế tạo tiếng nói.

► Có các nguồn âm

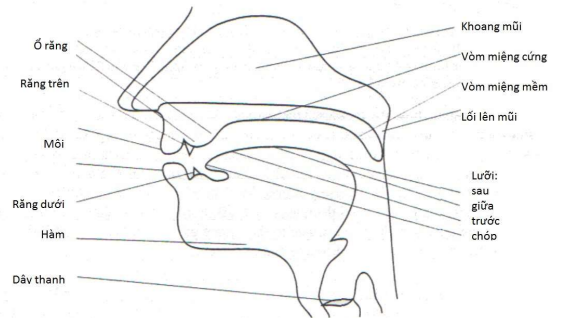
- Tuần hoàn (dây thanh rung)
- Tụp âm (dây thanh không rung)
- Xung

5



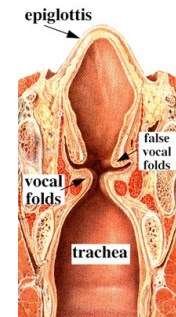
6

Bộ máy phát âm



7

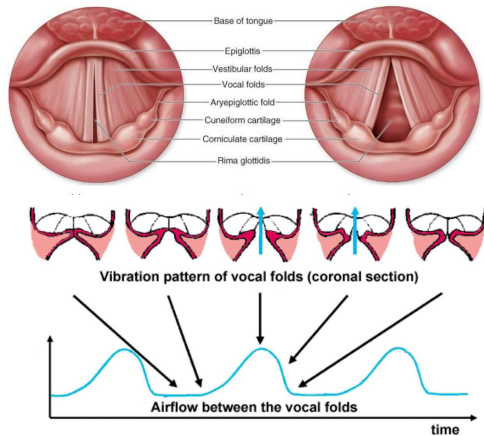
Bộ máy phát âm



NASAL CAVITY: Khoang mũi
 SOFT PALATE: Vòm miệng mềm
 EPIGLOTTIS: Nắp thanh quản
 VOCAL FOLDS (CORDS): Dây thanh
 OESOPHAGUS: Thực quản
 TRACHEA: Khí quản
 PHARYNX: Họng

8

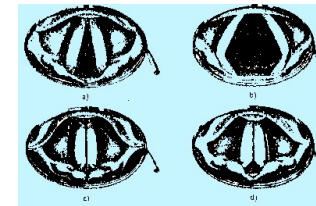
Dây thanh và Thanh môn



9

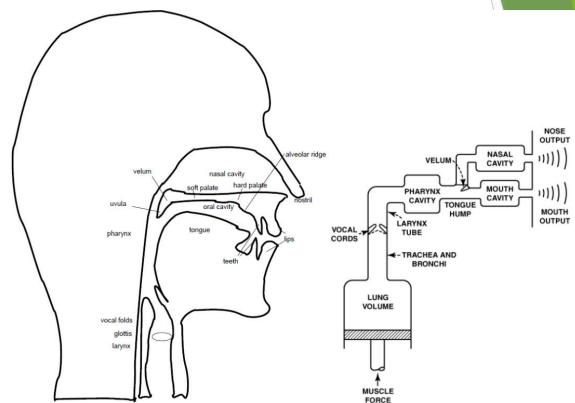
Thanh môn

- Ở các vị trí hít, thở, phát âm, nói thì thào



10

Sơ đồ khối bộ máy phát âm



11

Hệ thống thính giác

- Hệ thống thính giác có 2 thành phần quan trọng:
 - Cơ quan thính giác ngoại vi (tai)
 - Biến đổi áp suất âm thanh thành dao động cơ học kích thích tế bào thần kinh
 - Hệ thống thần kinh thính giác (não)
 - Trích xuất các thông tin cảm nhận được ở mức độ khác nhau

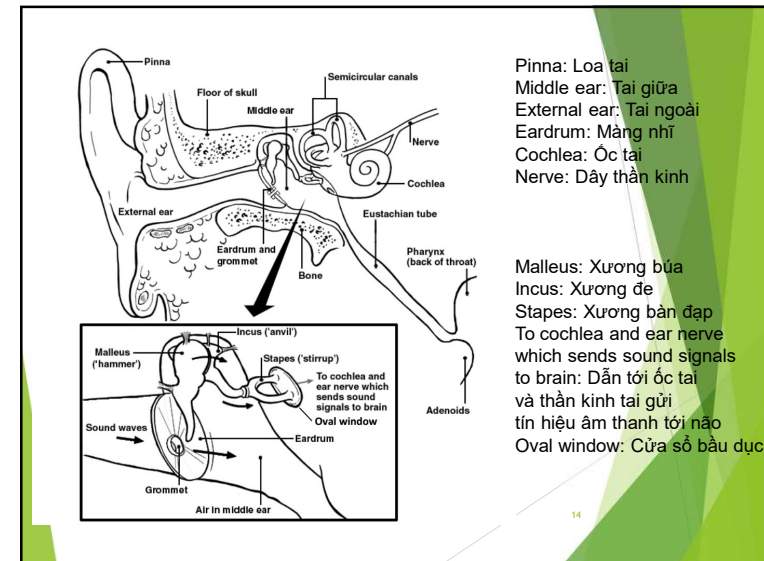
12

Hệ thống thính giác

- ▶ Tai có thể được phân chia
 - ▶ Tai ngoài:
 - ▶ Bao gồm loa tai, ống tai ngoài và màng nhĩ
 - ▶ Biến đổi áp suất âm thanh thành rung động
 - ▶ Tai giữa
 - ▶ Gồm các xương: xương búa, xương đe và xương bàn đạp
 - ▶ Vận chuyển rung động màng nhĩ vào tai trong
 - ▶ Tai trong:
 - ▶ Gồm ốc tai
 - ▶ Biến đổi các rung động thành các xung kích thích màng đáy
 - ▶ Màng đáy có thể được mô hình hóa như băng bộ lọc

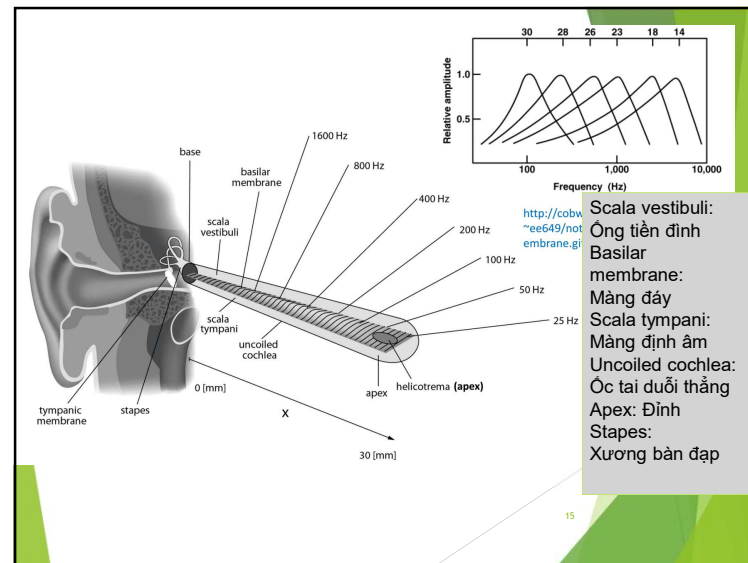
13

13



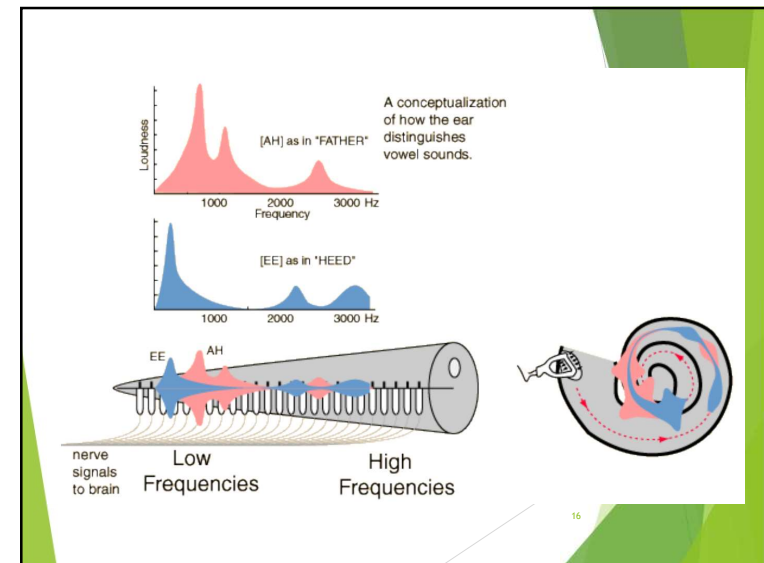
14

14



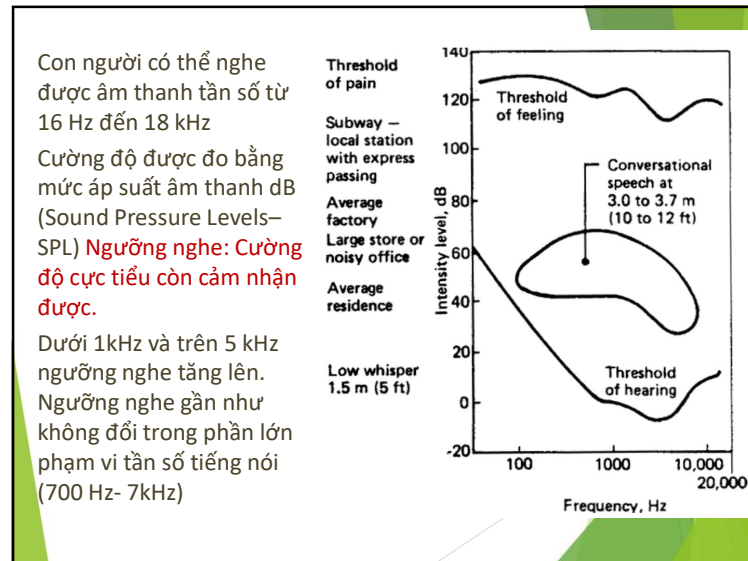
15

15



16

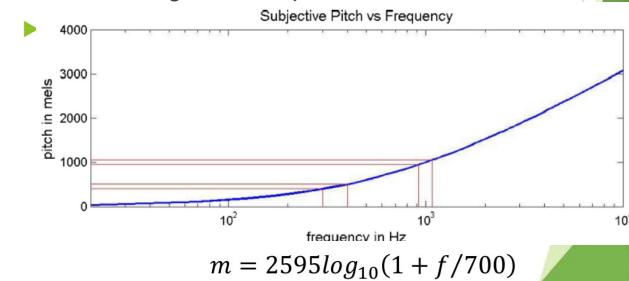
16



17

Cảm nhận cao độ (pitch)

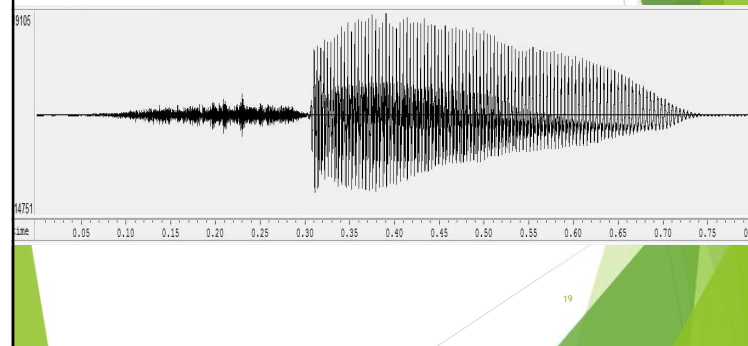
- ▶ Cao độ là F0 (tần số cơ bản) được cảm nhận, mang tính chủ quan



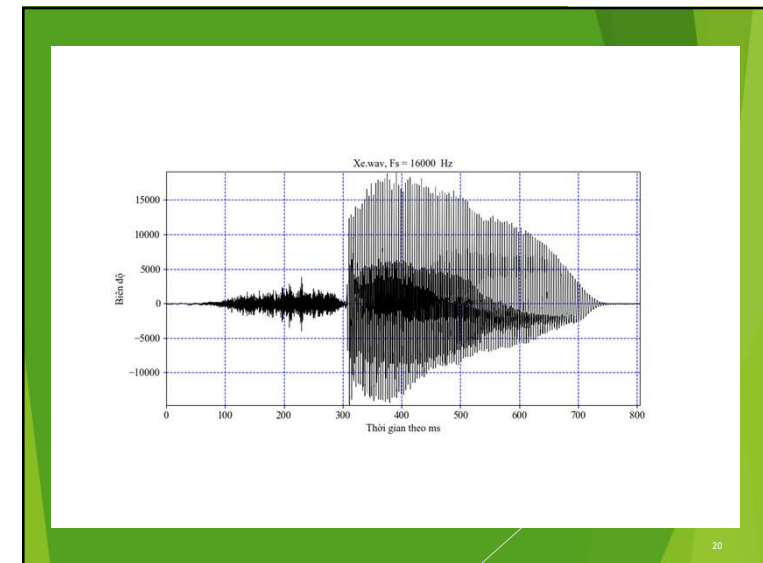
18

1.3 Một số phương pháp biểu diễn tín hiệu tiếng nói

- ▶ Dạng sóng theo thời gian



19



20

```

import numpy as np
import scipy.io.wavfile as wf
import matplotlib.pyplot as plt

filename = "Xe.wav"
rate, data = wf.read(filename)
data = data.astype(np.int32)
# Timing axe
Time = np.linspace(0, 1000 * len(data) / rate, num=len(data))
fig = plt.figure(figsize=(10, 5))

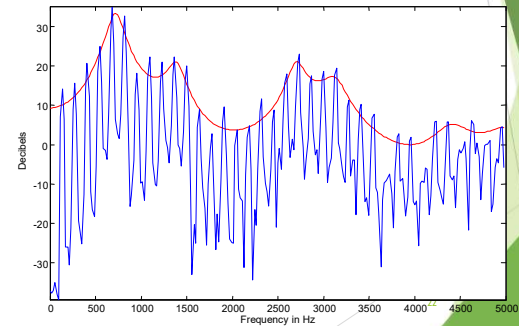
# Hạn biên theo chiều X và Y
plt.xlim(Time[0], Time[-1])
plt.ylim(np.min(data), np.max(data))
plt.yticks(np.min(data), np.max(data))
plt.xticks(fontsize = 12, family = 'Times New Roman')
plt.yticks(fontsize = 12, family = 'Times New Roman')
plt.grid(color='b', linestyle='dashed')
plt.xlabel(u'Thời gian theo ms', fontsize = 12, family = 'Times New Roman')
plt.ylabel(u'Biên độ', fontsize = 12, family = 'Times New Roman')
plt.title(filename + ", Fs = " + str(rate) + " Hz", fontsize = 12, family = 'Times New Roman')
# Draw wave form
plt.plot(Time, data, color='k', linewidth = 0.5)
plt.show()

```

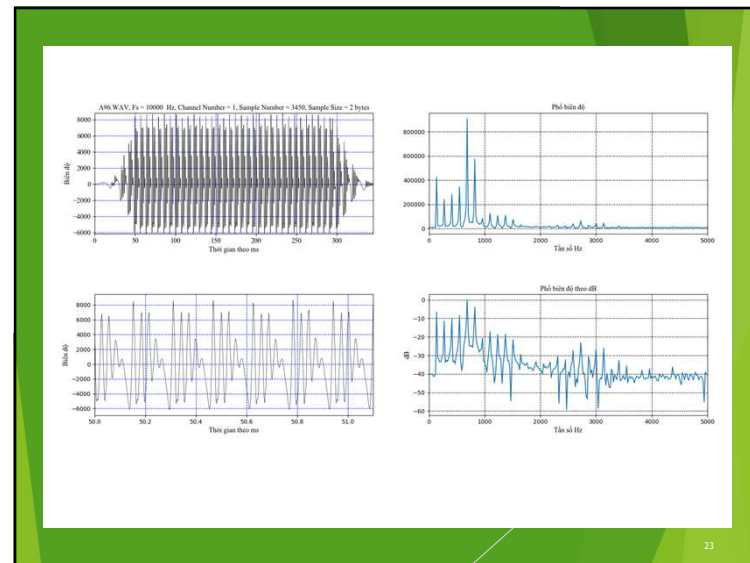
21

21

► Phổ tín hiệu tiếng nói



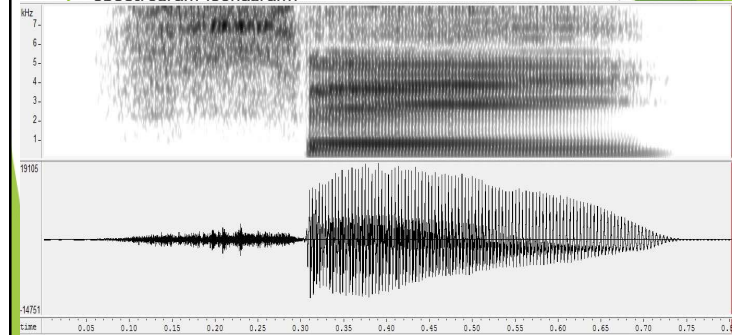
22



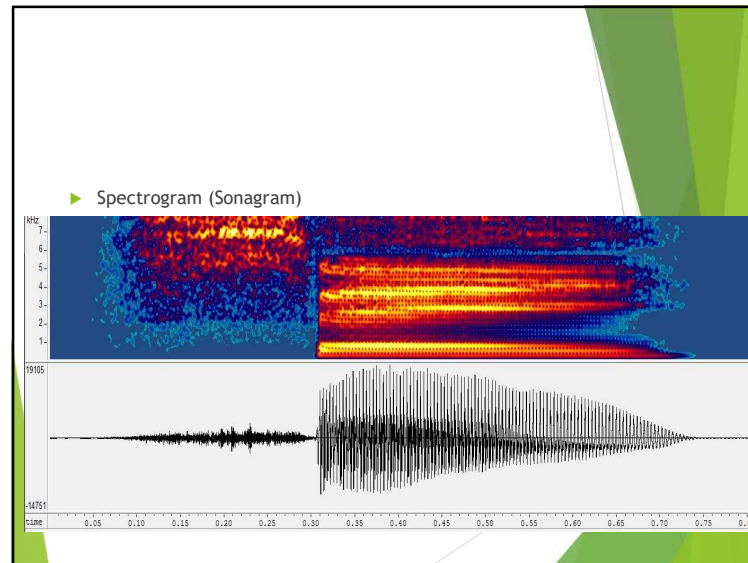
23

23

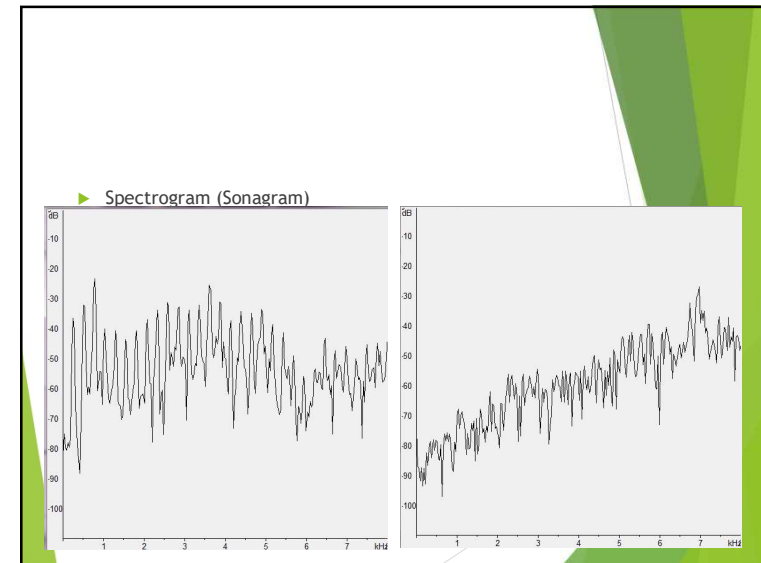
► Spectrogram (Sonagram)



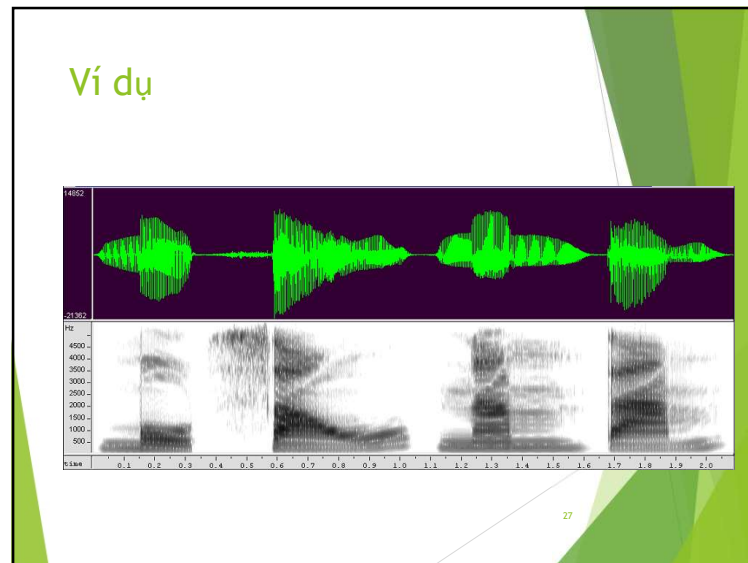
24



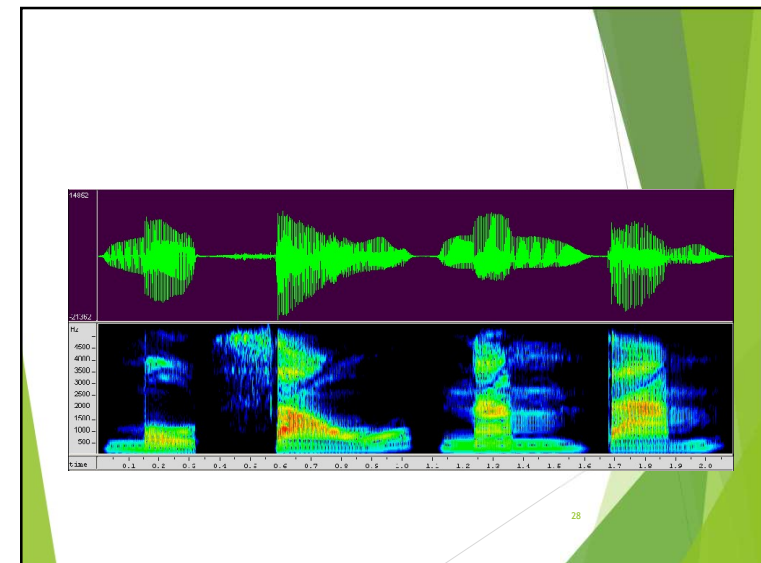
25



26

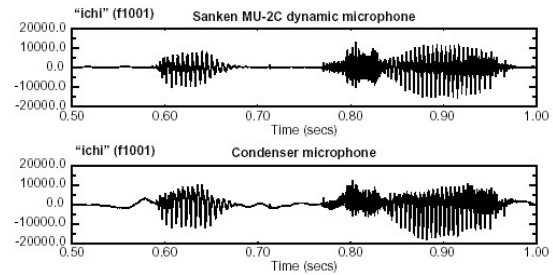


27



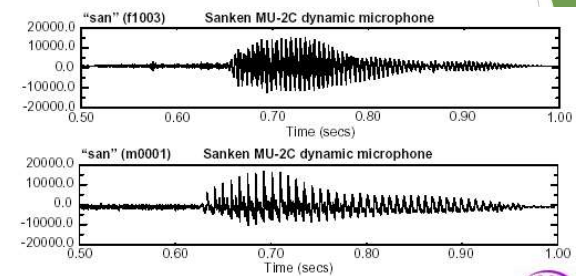
28

► Tín hiệu tiếng nói thu bằng micro khác loại



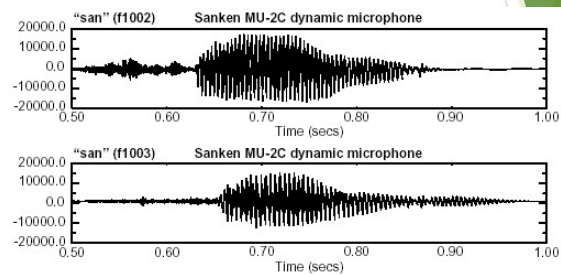
29

► Hai giọng khác nhau cho cùng một âm



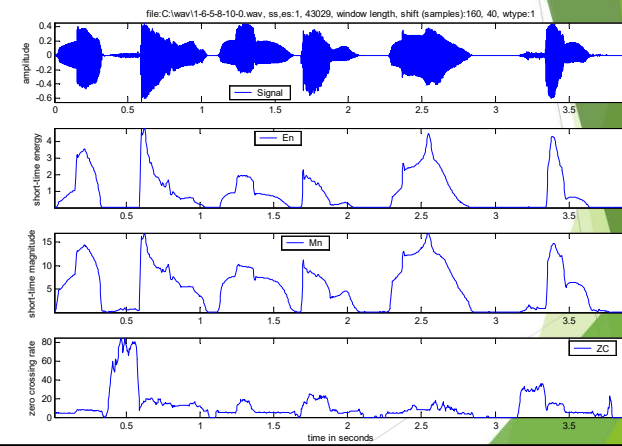
30

► Cùng người nói, cùng một âm



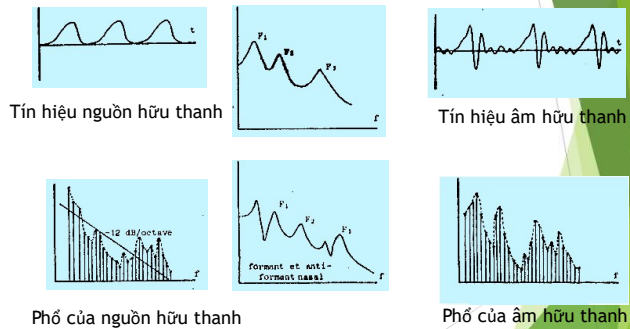
31

Năng lượng, tỷ lệ biến thiên qua giá trị không



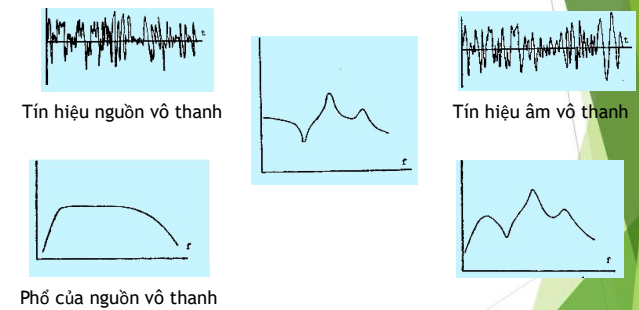
32

Tạo âm hữu thanh Formant và antiformant



33

Tạo âm vô thanh



34

1.4 các đặc điểm cơ bản ngữ âm tiếng Việt

- Đơn âm tiết
- Có thanh điệu (6), biến đổi thanh điệu kèm theo biến đổi nghĩa
- Không biến đổi hình thái

35

- Hệ thống âm vị: 14 nguyên âm (11 nguyên âm đơn, 3 nguyên âm đôi, 22 phụ âm)

1	i, y	ý chí
2	ê	ê chế
3	e	e dề
4	a	a ha
5	ă	mắt
6	ơ	bơ phờ
7	â	ân cần
8	ư	từ từ
9	ô	ótó
10	o	ơ ro
11	u	lù mù

1	ia, yê, ya, iê (đọc ia, yê)	kia kia, yêu kiểu, khuya, tiên tiên
2	ua, uô (đọc ua)	tua rua, luôn
3	ưạ, ươ (đọc ưạ)	lừa thua, lướt

36

- Hệ thống âm vị: 22 phụ âm

1	b	<i>bồng bênh</i>	12	tr	<i>trông</i>
2	p	<i>óp ép</i>	13	s	<i>sinh viên</i>
3	v	<i>vấn vơ</i>	14	r	<i>rừng</i>
4	ph	<i>phôi pha</i>	15	ch	<i>chông</i>
5	m	<i>mơ màng</i>	16	nh	<i>nhọc</i>
6	đ	<i>đất đai</i>	17	ng,ngh	<i>ngô nghê</i>
7	t	<i>tin tưởng</i>	18	c,k,q	<i>con,kẹt,qua</i>
8	th	<i>thơ thần</i>	19	kh	<i>khúc</i>
9	d,gi	<i>duyên, giữ</i>	20	g,gh	<i>gỗ ghê</i>
10	n	<i>nóng</i>	21	h	<i>hả hê</i>
11	l	<i>long lanh</i>	22	x	<i>xa xôi</i>

37

- Phân loại nguyên âm theo độ nâng của lưỡi và chuyển động của lưỡi

Độ nâng	cao	trung bình	thấp
Hàng			
trước	i e	e	
giữa	ư	ơ â	a ă
sau	u ô	o	

38

- Phân loại nguyên âm theo độ mở của miệng và chuyển động của lưỡi

Độ mở	Hàng	<i>hàng trước</i>	<i>hàng sau không tròn môi</i>	<i>hàng sau tròn môi</i>
hẹp	i	ia,yê,ya,iê	ư ơ	u ua
hơi hẹp	ê		ơ â	ô
hơi rộng	e			o
rộng			a ă	

39

- Âm tắc: tiếng nổ, phát sinh do luồng khí từ phổi đi ra bị cản trở hoàn toàn, phải phá vỡ sự cản trở đó để thoát ra.
- Âm xát: tiếng cọ xát, phát sinh do luồng không khí đi ra bị cản trở không hoàn toàn (chỉ bị khó khăn), phải lách qua một khe hở nhỏ và trong khi thoát ra như vậy phải cọ xát vào thành của bộ máy phát âm.
- Phụ âm bên: đầu lưỡi tiếp xúc với lợi chặn lối thoát của không khí, buộc nó phải lách qua khe hở ở hai bên cạnh lưỡi tiếp giáp với má mà ra ngoài tạo nên tiếng xát nhẹ (l).
- Luồng không khí thoát ra ngoài bị cản trở, tạo nên tiếng xát hay tiếng nổ, dạng tín hiệu không tuần hoàn gọi là tiếng động (ồn).
- Trong khi phát âm một số phụ âm, dây thanh cũng hoạt động đồng thời tạo nên tiếng thanh.
- Phụ âm có tỉ lệ tiếng động lớn hơn gọi là phụ âm ồn.
- Phụ âm có tỉ lệ tiếng thanh lớn hơn gọi là phụ âm vang.

40

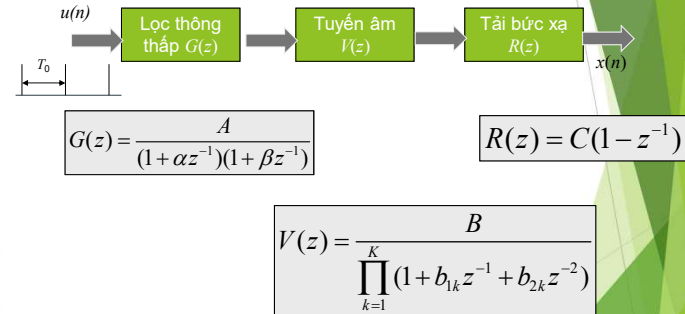
- Phân loại phụ âm theo tắc hay xát, hữu thanh hay vô thanh, mũi hóa

Vị trí cấu âm				Đầu lưỡi				
Phương thức cấu âm				Môi	Răng	Vòm miệng	Mặt lưỡi	Cuối lưỡi
Tắc	Ôn	Bật hơi	Vô thanh	p	t	tr	ch	c, k, qu
		Không bật hơi	Hữu thanh	b	đ			
	Xát	Vang mũi	m	n		nh	ng, ngh	
		Vô thanh	ph	x	s		kh	h
Xát	Ôn	Hữu thanh	v	d, gi	r		g	
		Vang bên		l				

41

41

1.5 Mô hình tạo tiếng nói



42

Mô hình toàn điểm cực (AR)

$$T(z) = G(z)V(z)R(z) = \frac{\sigma}{A(z)}$$

- $A(z)$: Hàm truyền đạt của bộ lọc đảo

$$T(z) = \frac{\sigma}{A(z)}$$

$$A(z) = 1 + \sum_{i=1}^{2K+1} a_i z^{-i} \quad A(z) = \sum_{i=0}^p a_i z^{-i} \quad a_0 = 1$$

$$x(n) + \sum_{i=1}^p a_i x(n-i) = \sigma u(n) \quad P = 2K+1$$

43

43

Mô hình ARMA (Autoregressive Moving Average)

$$T(z) = \frac{\sigma_1}{A_1(z)} + \frac{\sigma_2}{A_2(z)} = \sigma \frac{C(z)}{A(z)}$$

$$C(z) = \sum_{i=0}^q c_i z^{-i} \quad c_0 = 1$$

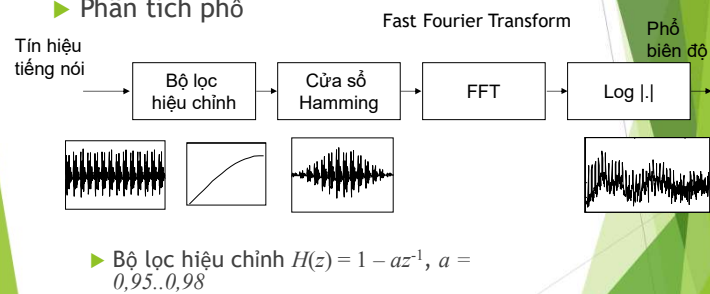
$$x(n) + \sum_{i=1}^p a_i x(n-i) = \sigma \sum_{i=0}^q c_i u(n-i)$$

44

44

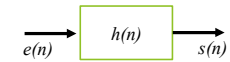
1.6 Các kỹ thuật cơ bản xử lý tín hiệu tiếng nói

► Phân tích phổ



45

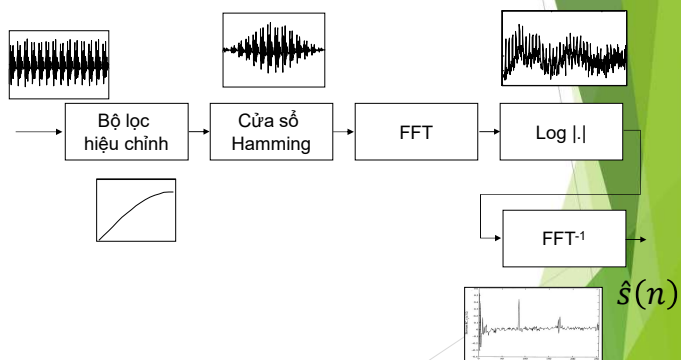
Xử lý đồng hình (homomorphic)



- $s(n) = h(n) * e(n) \rightarrow S(\omega) = H(\omega)E(\omega)$
- $\log S(\omega) = \log H(\omega) + \log E(\omega)$
- $\mathbb{F}^{-1}\{\log S(\omega)\} = \mathbb{F}^{-1}\{\log H(\omega)\} + \mathbb{F}^{-1}\{\log E(\omega)\}$
- $\hat{s}(n) = \hat{h}(n) + \hat{e}(n)$

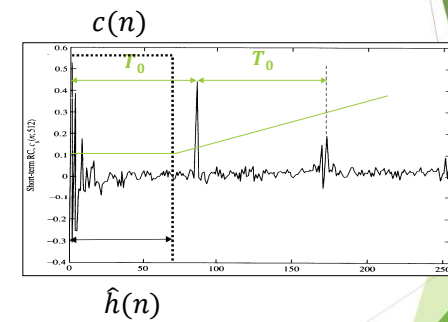
46

Sơ đồ khối xử lý đồng hình



47

Ví dụ



48

Tiên đoán tuyến tính (Linear Prediction Coding)

- Mô hình AR

Tiên đoán

Sai số tiên đoán

Sai số bình phương toàn phần

Tối thiểu hóa sai số

$$x(n) + \sum_{i=1}^p a_i x(n-i) = \sigma u(n)$$

$$\hat{x}(n) = - \sum_{i=1}^p \hat{a}_i x(n-i)$$

$$e(n) = x(n) - \hat{x}(n)$$

$$E = \sum_n e^2(n)$$

$$\frac{\partial E}{\partial \hat{a}_i} = 0, i = 1, 2, \dots, p$$

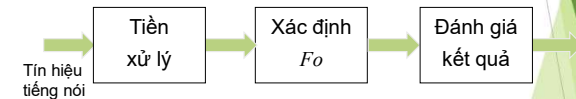
49

Xác định tần số cơ bản

- Giá trị F_0 phụ thuộc vào giới tính và lứa tuổi

- Giọng nam: 80..250 Hz

- Giọng nữ: 150..500 Hz



50

Một số phương pháp xác định F_0

- Dựa vào hàm tự tương quan
- Dựa vào hàm vi sai biên độ trung bình
- Xử lý đồng hình

51

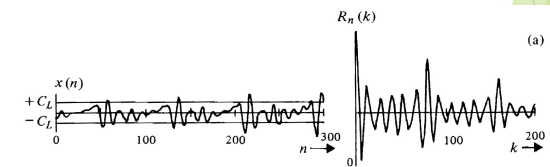
Dựa vào hàm tự tương quan

- Tính hàm tự tương quan $R(k)$ của tín hiệu tiếng nói $x(n)$

$$R(k) = \sum_{n=0}^{N-1-k} x(n)x(n+k) \quad k = 0, 1, \dots, K$$

$$Fs = 10 \text{ kHz}, N = 300, K = 150.$$

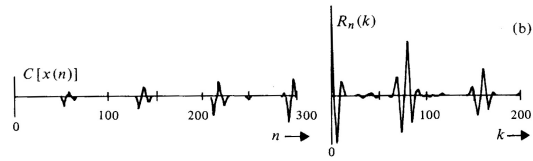
Tìm cực đại trong khoảng $(0, K)$



52

Phương pháp tự tương quan có cải tiến

- ▶ Hạn chế, loại bỏ $|x| < C_L$



53

Dựa vào hàm vi sai biên độ trung bình (AMDF- Average Magnitude Difference Function)

$$D(k) = \sum_{m=0}^{N-1} |x(n+m) - x(n+m-k)| \quad k = 0, 1, \dots, K$$

$$D(iP) = 0, \quad i = 0, 1, \dots, \quad \frac{1}{N} \sum_{n=0}^{N-1} |u(n)| \leq \left[\frac{1}{N} \sum_{n=0}^{N-1} u^2(n) \right]^{1/2}$$

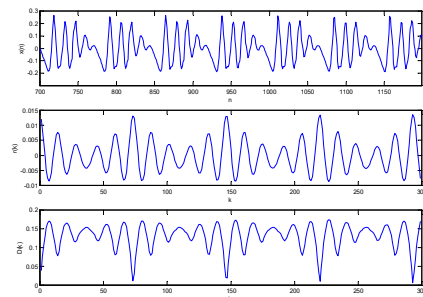
$$D(k) = \lambda \left\{ \frac{1}{N} \sum_{m=0}^{N-1} [x(n+m) - x(n+m-k)]^2 \right\}^{1/2}$$

$$= \lambda \left\{ \frac{1}{N} [2r(0) - 2r(k)] \right\}^{1/2} \quad k = 0, 1, \dots, K$$

với $\lambda < 1$

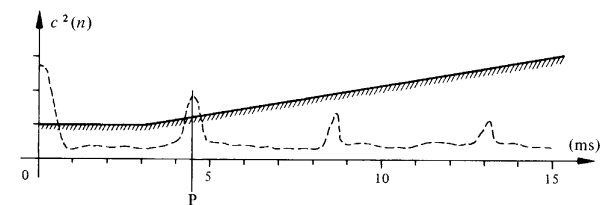
54

Ví dụ



55

Xử lý đồng hình



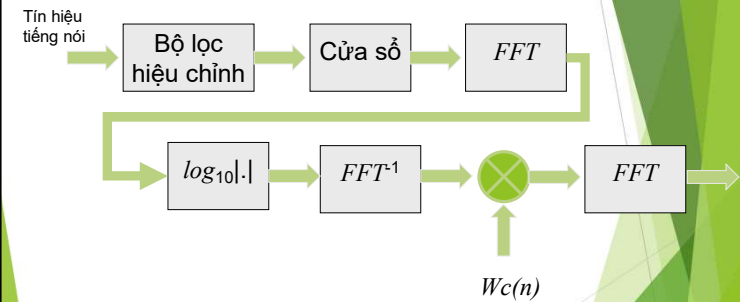
56

Xác định formant

- ▶ Tham số cần xác định
 - ▶ Formant F_k
 - ▶ Dải thông B_k
- ▶ Phương pháp
 - ▶ Xử lý đồng hình
 - ▶ LPC

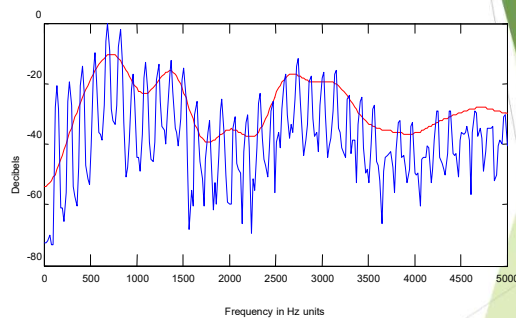
57

Xử lý đồng hình



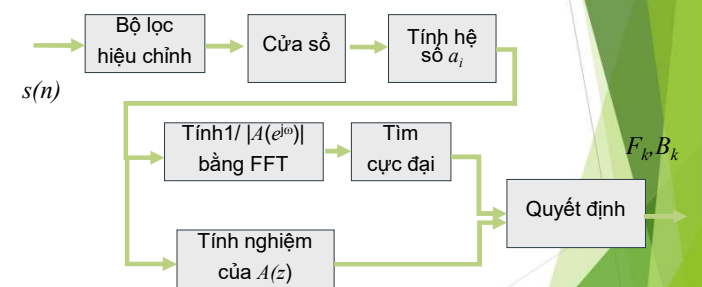
58

Xử lý đồng hình



59

Phương pháp LPC



60

57

58

59

60