

# Spoken Language Structure

Berlin Chen

Department of Computer Science & Information Engineering  
National Taiwan Normal University

## References:

- X. Huang et. al., Spoken Language Processing, Chapter 2
- 王小川教授，語音訊號處理，Chapters 2~3

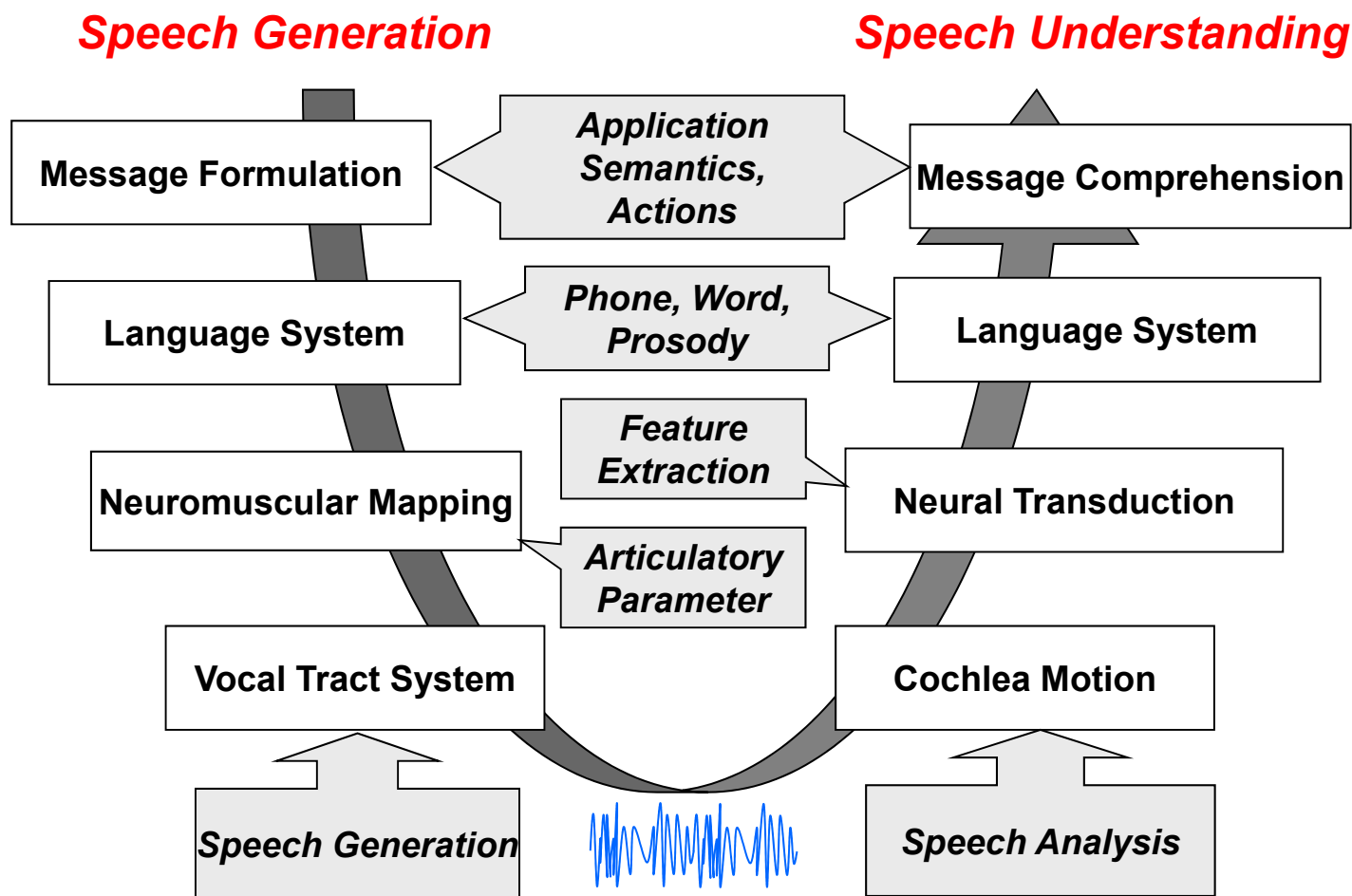
# Introduction

- Take a bottom-up approach to introduce the basic concepts from sound to phonetics (語音學) and phonology (音韻學)
  - Syllables (音節) and words (詞) are followed by syntax (語法) and semantics (語意), which form the structure of spoken language processing
- Topics covered here
  - Speech Production
  - Speech Perception
  - Phonetics and Phonology
  - Structural Features of the Chinese Language

# Determinants of Speech Communication

- Spoken language is used to communicate information from a speaker to a listener. Speech production and perception are both important components of the speech chain
- Speech signals are composed of analog sound patterns that serve as the basis for a discrete, symbolic representation of the spoken language – phonemes, syllables and words
- The production and interpretation of these sounds are governed by the ***syntax*** and ***semantics*** of the language spoken

# Determinants of Speech Communication (cont.)



# Computer Counterpart

- The Speech Production Process
  - **Message formulation:** create the concept (message) to be expressed
  - **Language system:** convert the message into a sequence of words and find the pronunciation of the words (or the phoneme sequence).
    - Apply the prosodic pattern: duration of phoneme, intonation (語調) of the sentence, and the loudness of the sounds
  - **Neuromuscular (神經肌肉) Mapping:** perform articulatory (發聲的) mapping to control the vocal cords, lips, jaw, tongue, etc., to produce the sound sequence

# Computer Counterpart (cont.)

- The Speech Understanding Process
  - **Cochlea (耳蝸) motion:** the signal is passed to the cochlea in the inner ear, which performs the **frequency analysis** as a filter bank
  - **Neural transduction:** converts the spectral signal into activity signals on the auditory nerve, corresponding to a **feature extraction** component

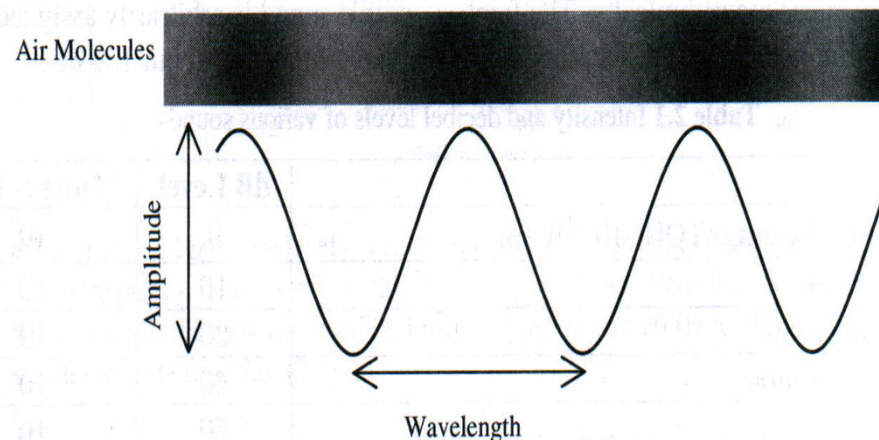
*It's unclear how neural activity is mapped into the language system and how message comprehension (理解) is achieved in the brain*

# Sound

- Sound is a longitudinal (縱向的) pressure wave formed of compressions (壓縮) and rarefactions (稀疏) of air molecules (微粒), in a direction parallel to that of the application of energy
- *Compressions* are zones where air molecules have been forced by the application of energy into a tighter-than-usual configuration
- *Rarefactions* are zones where air molecules are less tightly packed

## Sound (cont.)

- The alternating configurations of compression and rarefaction of air molecules along the path of an energy source are sometimes described by the graph of a **sine wave**
- The use of the sine graph is only a notational convenience for charting local pressure variations over time



**Figure 2.2** Application of sound energy causes alternating compression/rarefaction of air molecules, described by a sine wave. There are two important parameters, amplitude and wavelength, to describe a sine wave. Frequency [cycles/second measured in Hertz (Hz)] is also used to measure of the waveform.



# Measures of Sounds

- **Amplitude** is related to *the degree of displacement of the molecules* from their resting position
  - Measured on a logarithmic scale in *decibels* (dB, 分貝)
  - A decibel is a means for comparing the intensity (強度) of two sounds:  
$$10 \log_{10} (I / I_0)$$
 $I, I_0$  are two intensity levels
  - The intensity is proportional to the **square** of the sound pressure  $P$ . The **Sound Pressure Level** (SPL) is a measure of the absolute sound pressure  $P$  in dB  
$$SPL \text{ (dB)} = 20 \log_{10} (P / P_0)$$
  - The reference 0 dB corresponds to the **threshold of hearing**, which is  $P_0=0.00002 \mu\text{bar}$  for a tone of 1KHz
    - E.g., the speech conversation level at 3 feet is about 60dB SPL; a jackhammer's level is about 120 db SPL

# Measures of Sound (cont.)

- **Absolute threshold of hearing:** is the maximum amount of energy of a pure tone that cannot be detected by a listener in a noise free environment

Table 2.1 Intensity and decibel levels of various sounds.

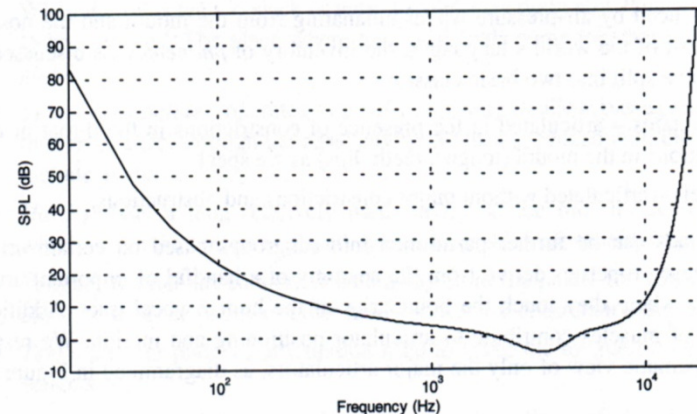
Sound	dB Level	Times > TOH
Threshold of hearing (TOH: $10^{-12} W/m^2$ )	0	$10^0$
Light whisper	10	$10^1$
Quiet living room	20	$10^2$
Quiet conversation	40	$10^4$
Average office	50	$10^5$
Normal conversation	60	$10^6$
Busy city street	70	$10^7$
Acoustic guitar – 1 ft. away	80	$10^8$
Heavy truck traffic	90	$10^9$
Subway from platform	100	$10^{10}$
Power tools	110	$10^{11}$
Pain threshold of ear	120	$10^{12}$
Airport runway	130	$10^{13}$
Sonic boom	140	$10^{14}$
Permanent damage to hearing	150	$10^{15}$
Jet engine, close up	160	$10^{16}$
Rocket engine	180	$10^{18}$
Twelve ft. from artillery cannon muzzle ( $10^{10} W/m^2$ )	220	$10^{22}$

in sound pressure level

The absolute threshold of hearing is the maximum amount of energy of a **pure tone** that cannot be detected by a listener in a noise free environment. The absolute threshold of hearing is a function of frequency that can be approximated by

$$T_q(f) = 3.64(f/1000)^{-0.8} - 6.5e^{-0.6(f/1000-3.3)^2} + 10^{-3}(f/1000)^4 \quad (\text{dB SPL}) \quad (2.3)$$

and is plotted in Figure 2.3.



**Figure 2.3** The sound pressure level (SPL) level in dB of the absolute threshold of hearing as a function of frequency. Sounds below this level are inaudible. Note that below 100 Hz and above 10 kHz this level rises very rapidly. Frequency goes from 20 Hz to 20 kHz and is plotted in a logarithmic scale from Eq. (2.3).

# Speech Production

## – *Articulation*

- Speech
  - Produced by air-pressure waves emanating (發出) from the mouth and the nostrils(鼻孔)
  - The inventory of **phonemes** (音素) are the basic units of speech and split into two classes
    - **Consonant** (子音/輔音)
      - Articulated (發音) when constrictions (壓縮) in the throat or obstructions (阻塞) in the mouth
    - **Vowel** (母音/元音)
      - without major constrictions and obstructions

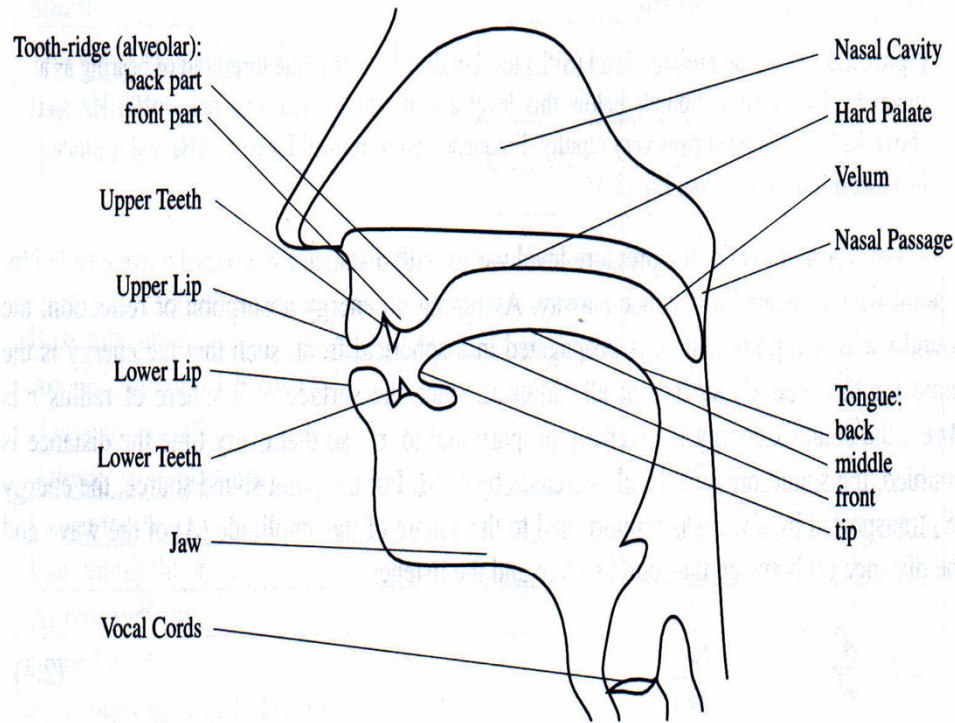
# Speech Production

## – *Articulation* (cont.)

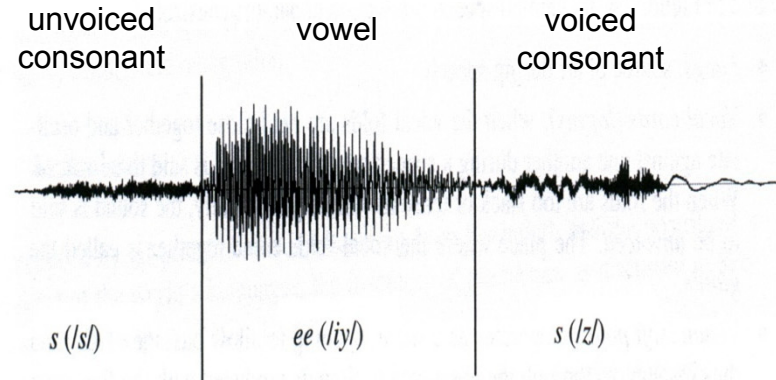
- Human speech production apparatus
  - **Lungs** (肺): source of air during speech
  - **Vocal cords** (larynx, 喉頭): when the vocal folds (聲帶) are held close together and oscillate one another during a speech sound, the speech sound is said to be **voiced** (<=>**unvoiced**)
  - **Soft Palate** (Velum, 軟顎): allow passage of air through the nasal cavity
  - **Hard palate** (硬顎): tongue placed on it to produce certain consonants
  - **Tongue**(舌): flexible articulator, shaped away from palate for vowel, closed to or on the palate or other hard surfaces for consonant
  - **Teeth**: braces (支撐) the tongue for certain consonants
  - **Lips**(嘴唇): round or spread to affect vowel quality, closed completely to stop the oral air flow for certain consonants (*p, b, m*)

# Speech Production

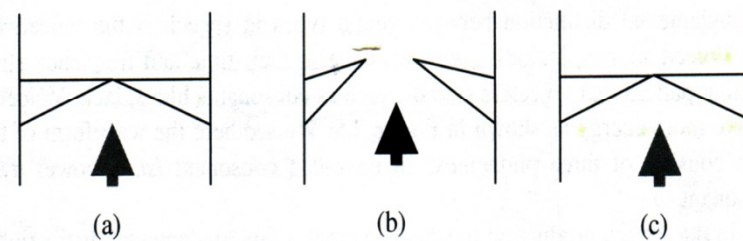
## – Articulation (cont.)



**Figure 2.4** A schematic diagram of the human speech production apparatus.



**Figure 2.5** Waveform of *sees*, showing a voiceless phoneme /s/, followed by a voiced sound, the vowel /iy/. The final sound, /z/, is a type of voiced consonant.



**Figure 2.6** Vocal fold cycling at the larynx. (a) Closed with sub-glottal pressure buildup; (b) trans-glottal pressure differential causing folds to blow apart; (c) pressure equalization and tissue elasticity forcing temporary reclosure of vocal folds, ready to begin next cycle.

# Speech Production

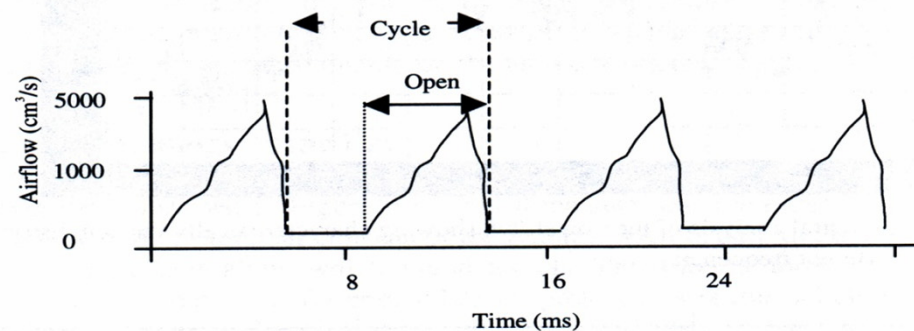
## - *The Voicing Mechanisms*

- Voiced sounds
  - Including vowels, have a roughly **regular pattern** in both time and frequency structures than voiceless sounds
  - Have more energy
  - Vocal folds vibrate during phoneme articulation (otherwise is **unvoiced**)
    - Vocal folds' vibration (60H ~ 300 Hz, cycles in sec.)
    - 男生分佈較低，女生分佈較高
    - The greater mass and length of adult male vocal folds as opposed to female
  - In psychoacoustics, the distinct vowel timbres (of a sound of an instrument, 音質/色) is determined by how the tongue and lips shaping the oral resonance (共鳴/振) cavity

# Speech Production

## - *The Voicing Mechanisms (cont.)*

- Voiced sounds (cont.)
  - The rate of cycling (open and closing) of vocal folds in the larynx during phonation of voiced sounds is called the **fundamental frequency** (基頻)
    - The fundamental frequency contributes more than any other single factor to the perception of **pitch** in speech
    - A prosodic feature for use in recognition of tonal languages (e.g., Chinese) or as a measure of speaker identity or authenticity



**Figure 2.7** Waveform showing air flow during laryngeal cycle.

# Speech Production

## - *Pitch*

The term *pitch* is often used interchangeably with fundamental frequency. However, there is a subtle difference. Psychoacousticians (scientists who study the perception of sound) use the term *pitch* to refer to the *perceived* fundamental frequency of a sound, *whether or not that sound is actually present in the waveform*. Speech transmitted over the commercial phone lines, for example, are usually bandlimited to about 300–3000 Hz. Nevertheless, a person who is phonating at 110 Hz will be *perceived* as phonating at 110 Hz by the listener, even though the fundamental frequency of the received waveform cannot be less than 300 Hz. In this case, the psychoacoustician would say that the pitch of the received speech waveform is 110 Hz, while the lowest frequency in the signal is 330 Hz. This quirk of the human auditory system requires that we be careful with these terms. Nevertheless, with this caution, we will routinely use the word “pitch” to mean “fundamental frequency” in this book, since it is conventional to do so. Since we will not be concerned with perceptual phenomena, this will not cause ambiguities to arise.



# Speech Production

## - Formants

- The resonances (共振/共鳴) of the cavities that are typical of particular articulator configurations (e.g. the different vowel timbres) are called **formants** (共振峰)

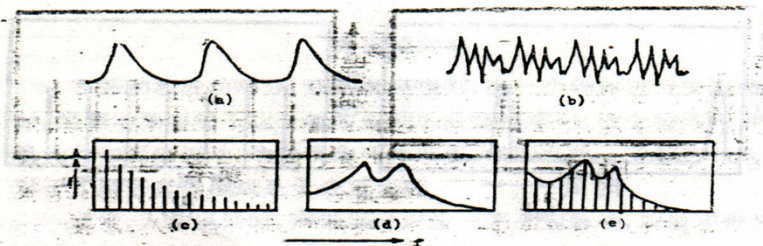
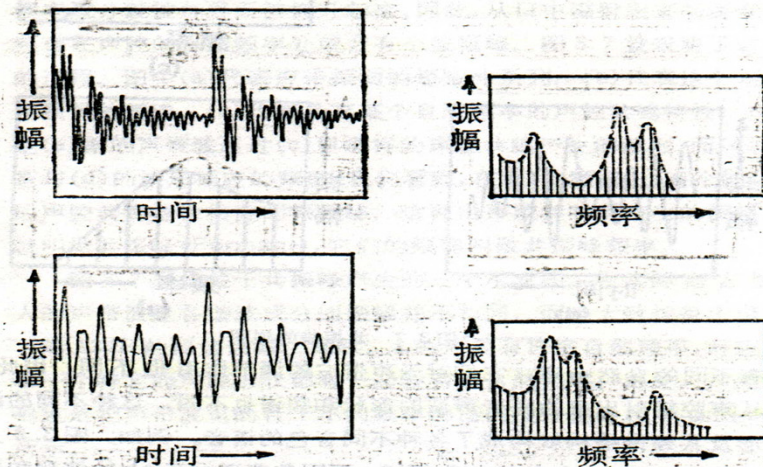
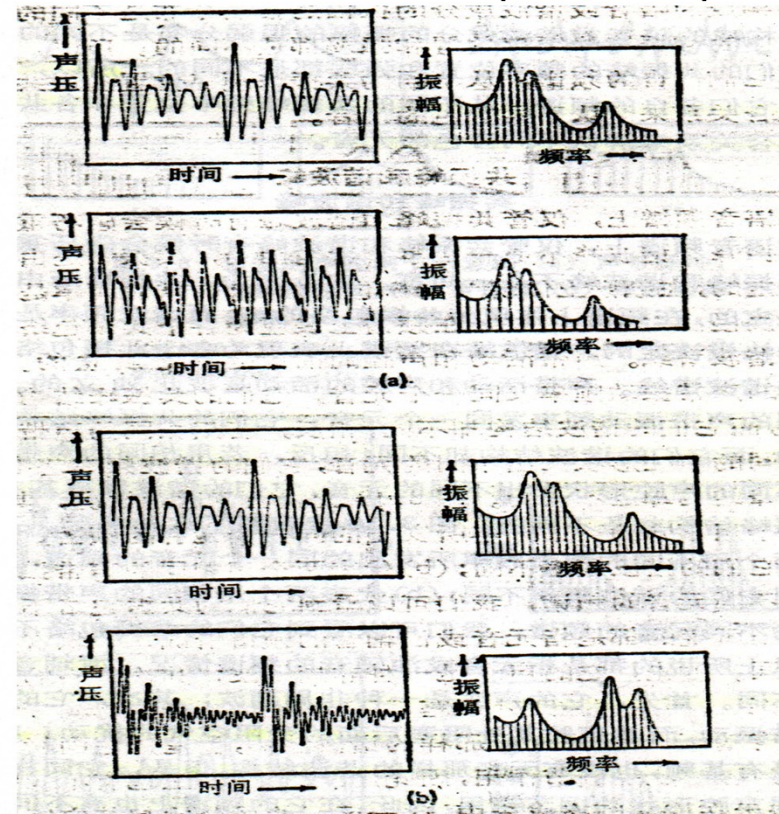


图 3.8 语音的声腔调节示意图: (a) 声带脉冲, (b) 辐射波, (c) 声带谱, (d) 声腔共振包络, (e) 语音频谱  
(此图采自 G. Fant, Acoustic Theory of Speech Production)



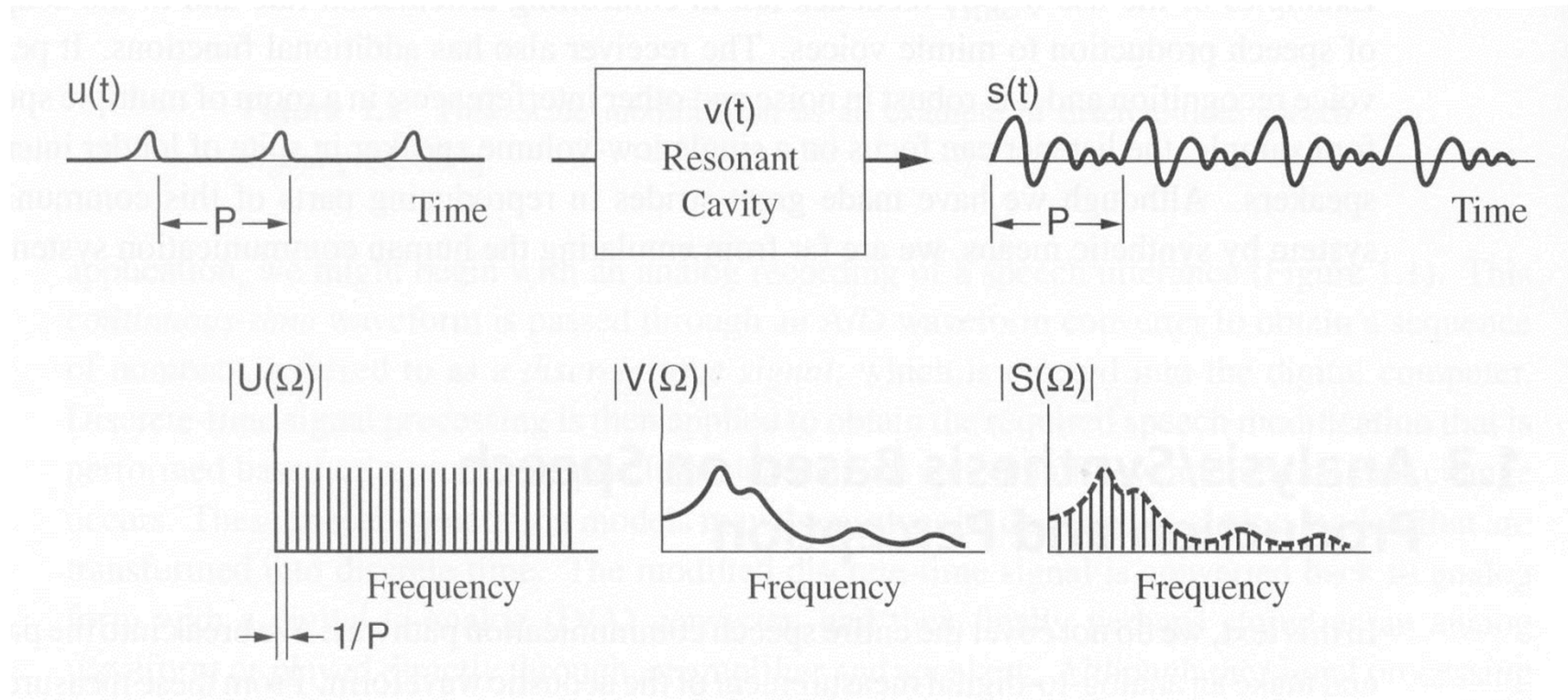
\*图 3.9 英语元音 uh(上图)和 ah(下图) 的波形及频谱



\*图 3.10 (a) 用不同的声带频率所发的同一元音的波形和频谱  
(b) 用相同的声带频率所发的不同元音的波形和频谱

# Speech Production

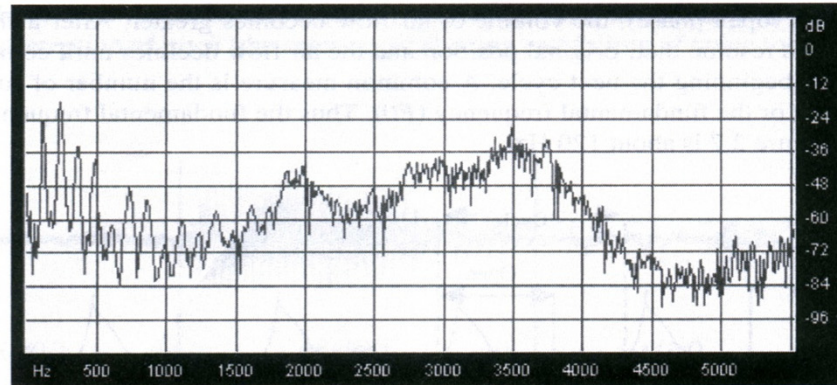
## - Formants (cont.)



**Figure 1.2** Speech production mechanism and model of a steady-state vowel. The acoustic waveform is modeled as the output of a linear time-invariant system with a periodic impulse-like input. In the frequency domain, the vocal tract system function spectrally shapes the harmonic input.

# Speech Production

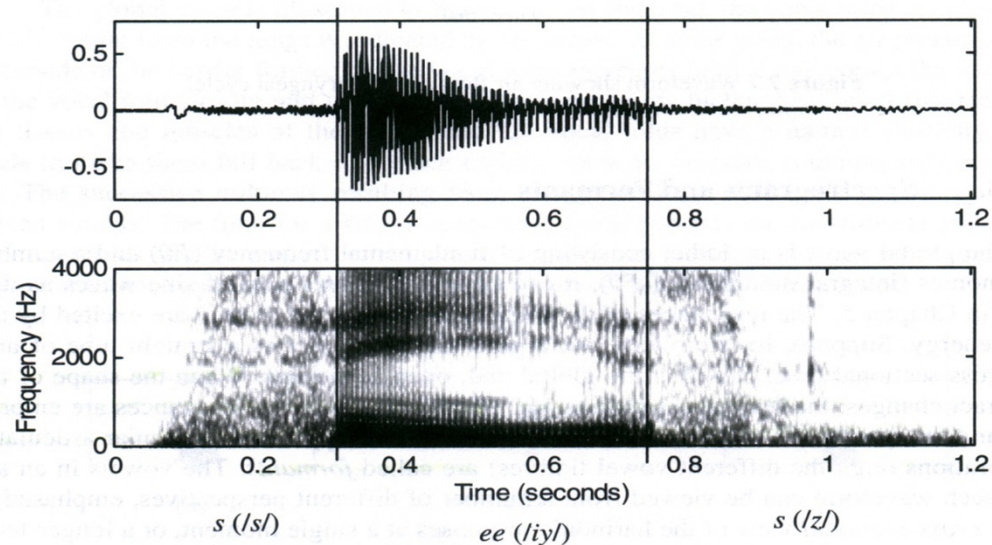
## - Formants (cont.)



Spectrum

頻譜

**Figure 2.8** A spectral analysis of the vowel /iy/, showing characteristically uneven distribution of energy at different frequencies.



Spectrogram

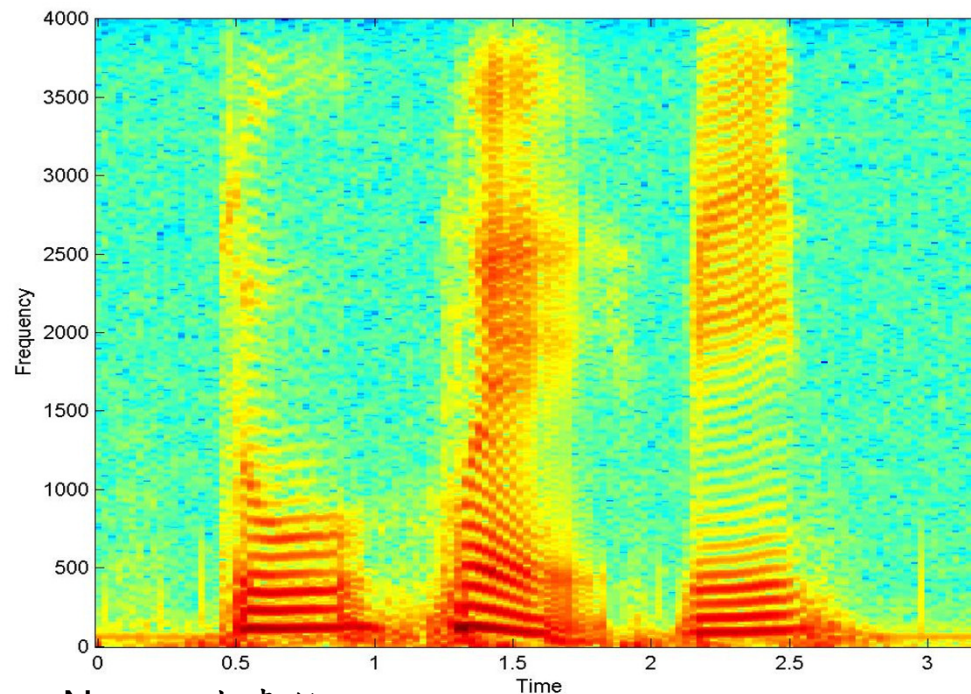
聲譜圖

**Figure 2.9** The spectrogram representation of the speech waveform *sees* (approximate phone boundaries are indicated with heavy vertical lines).

# Speech Production

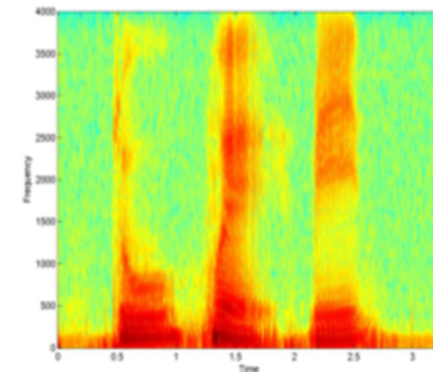
## - Formants (cont.)

- Narrowband Spectrogram
  - Both pitch harmonic and formant information can be observed



Name: 朱惠銘

1024-point FFT, 400 ms/frame, 200 ms/frame move



100 ms/frame,  
50 ms/frame move

Wide-band spectrograms : shorter windows (<10ms)

- Have good time resolution

Narrow-band spectrograms : Longer windows (>20ms)

- The harmonics can be clearly seen

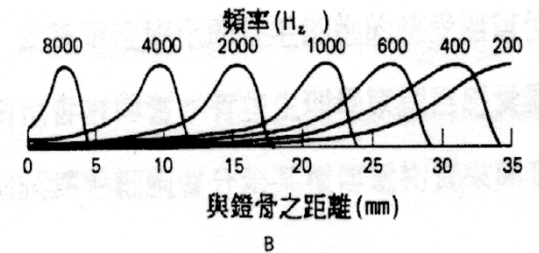
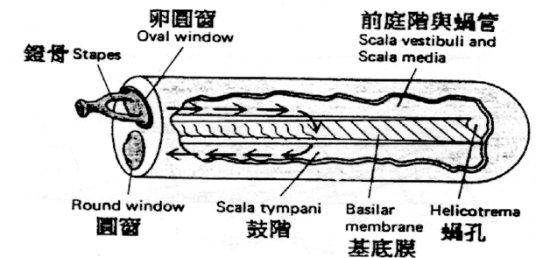
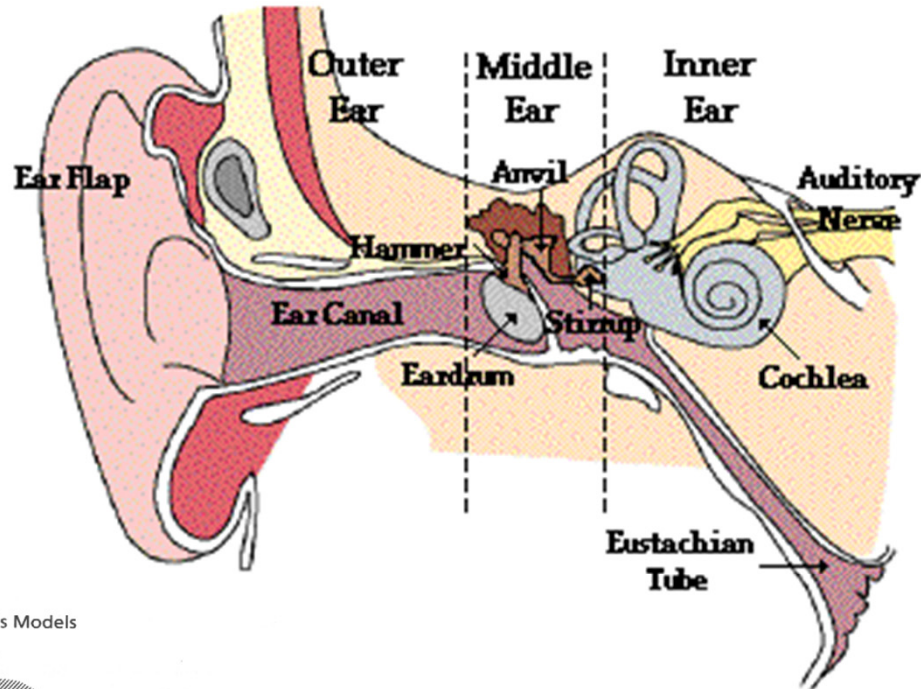
# Speech Perception

## *Physiology of the Ear*

- The ear processes an acoustic pressure signal by
  - First transforming it into a mechanical vibration pattern on the basilar membrane (基底膜)
  - Then representing the pattern by a series of pulses to be transmitted by the auditory nerve
- Physiology of the Ear
  - When air pressure variations reach the eardrum from the outside, it vibrates, and transmits the vibrations to bones adjacent to its opposite side
  - Then the energy is transferred by mechanical action of the stapes into an impression on the membrane stretching over the oval window
  - The cochlea can be roughly regarded as a set of filter banks, whose outputs are ordered by location
    - Frequency-to-place transformation

# Speech Perception

## Physiology of the Ear (cont.)



### Sec. 3.5 Auditory-Based Spectral Analysis Models

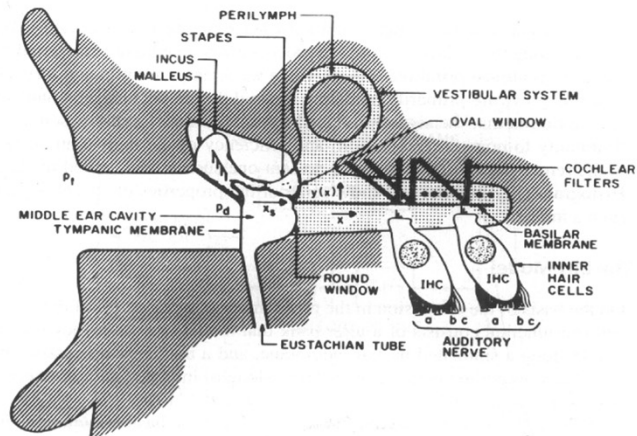
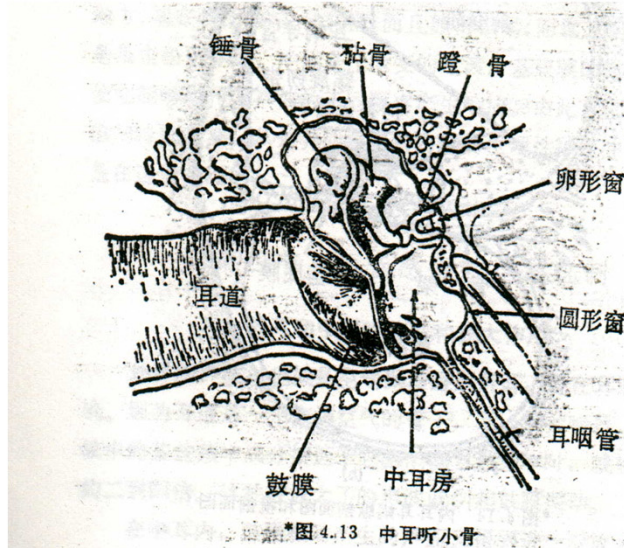


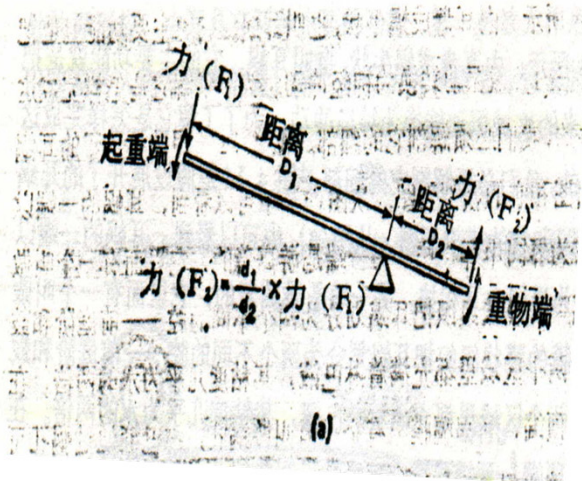
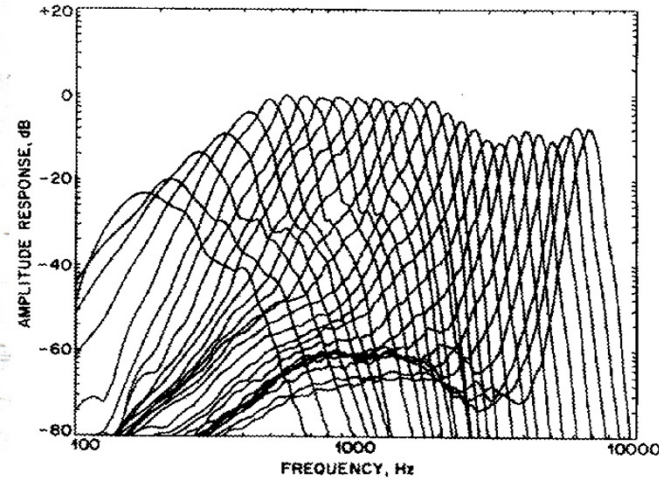
Figure 3.48 Expanded view of the middle and inner ear mechanics.

# Speech Perception

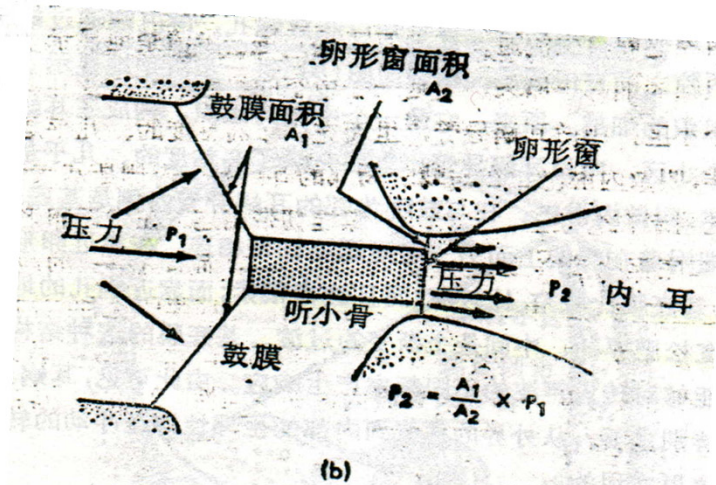
## Physiology of the Ear (cont.)



\*图 4.13 中耳听小骨



(a)



(b)

\*图 4.15 中耳放大作用示意图

# Speech Perception

## *Physical vs. Perceptual Attributes*

- **Non-uniform equal loudness perception** of tones of varying frequencies
  - Tones of different pitch have different perceived loudness
  - **Sensitivity** of the ear varies with the frequency and the quality of sound
  - Hear sensitivity reaches a maximum around 4000 Hz

**Table 2.2** Relation between perceptual and physical attributes of sound.

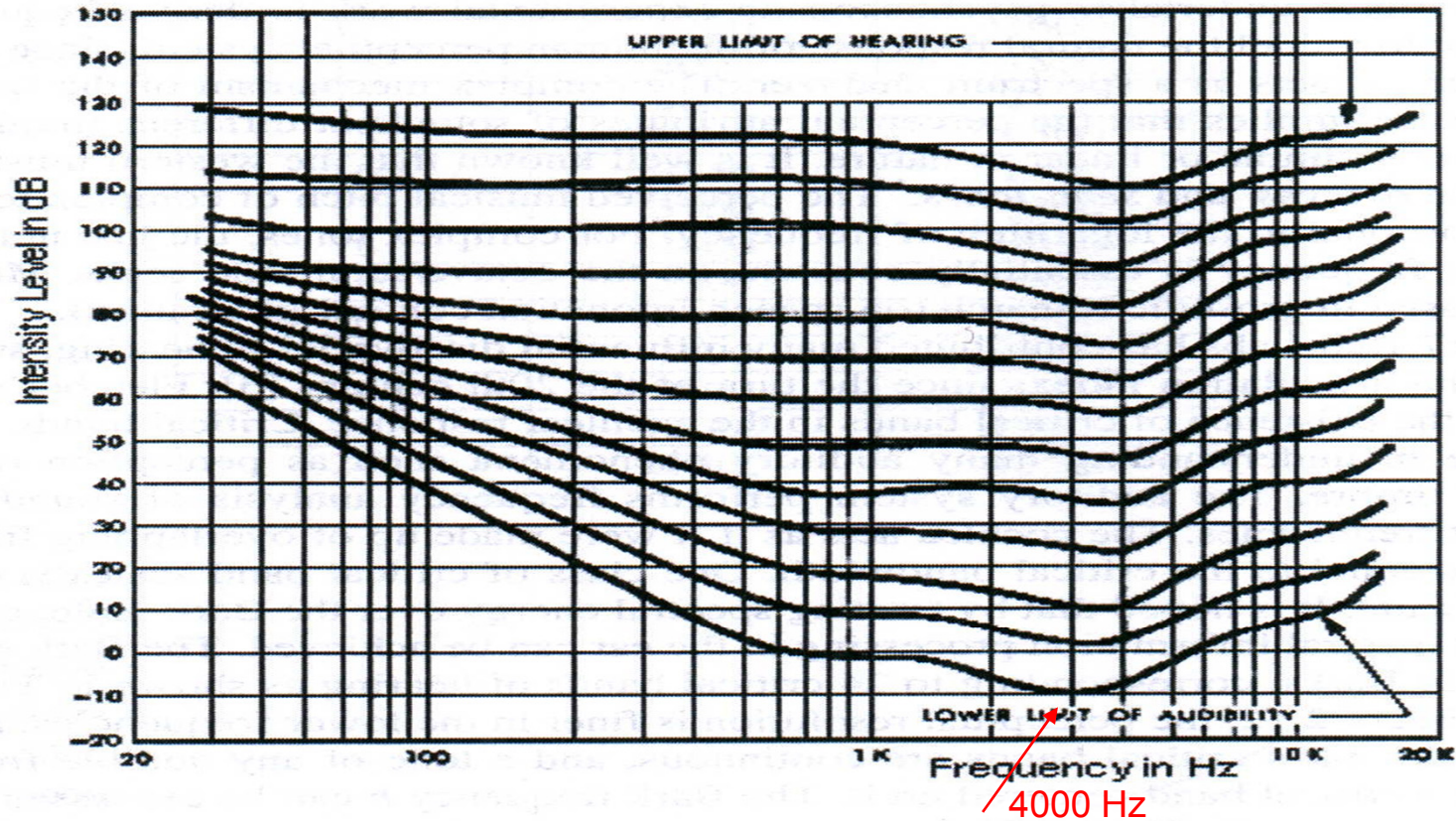
<b>Physical Quantity</b>	<b>Perceptual Quality</b>
Intensity	Loudness
Fundamental frequency	Pitch
Spectral shape	Timbre
Onset/offset time	Timing
Phase difference in binaural hearing	Location



# Speech Perception

## *Physical vs. Perceptual Attributes*

- Non-uniform equal loudness perception

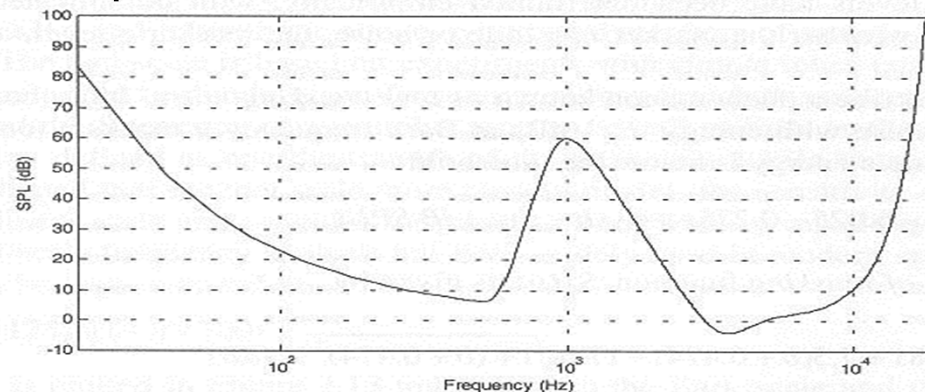


**Figure 2.11** Equal-loudness curves indicate that the response of the human hearing mechanism is a function of frequency and loudness levels. This relationship again illustrates the difference between physical dimensions and psychological experience (after ISO 226).

# Speech Perception

## *Physical vs. Perceptual Attributes (cont.)*

- Masking: when the ear is exposed to two or more different tones, it's a common experience that one tone may *mask* others
  - An **upward shift** in the hearing threshold of the weaker tone by the louder tone
  - A pure tone masks of higher frequency more effectively than those of lower frequency
  - The greater the intensity of the masking tone, the broader the range of frequencies it can mask



**Figure 2.15** Absolute threshold of hearing and spread of masking threshold for a 1 kHz sine-wave masker with a 69 dB SPL. The overall masked threshold is approximately the largest of the two thresholds.

# Speech Perception

## *Physical vs. Perceptual Attributes (cont.)*

- The **sense of localization attention** (Lateralization)
  - Binaural listening greatly enhances our ability to sense the direction of the sound source
  - Time and intensity cues have different impacts for low frequency and high frequency, respectively
    - **Low-frequency** sounds are lateralized mainly on the basis of interaural **time** differences
    - **High-frequency** sounds are lateralized mainly on the basis of interaural **intensity** differences
- The question of distinct voice quality
  - Speech from different people sounds different, e.g., different fundamental frequencies, different vocal-tract length
  - **The concept of timbre** (音質) is defined as that the attribute of auditory sensation by which a subject can judge that two sounds similarly presented and having the same loudness and pitch are dissimilar

# Speech Perception

## *Frequency Analysis*

- Researchers undertook psychoacoustic (心理聲學) experimental work to derive frequency scales that attempt to model the natural response of the human perceptual system (*the cochlea acts as a spectrum analyzer*)
  - The perceptual attributes of sounds at different frequencies may not be entirely simple or linear in natural
- **Bark Scale:** Fletcher's work (1940) pointed to the existence of critical bands in the cochlear response
  - The **cochlea** acts as if it were made up of **overlapping filters** having bandwidth equal to the critical bandwidth
  - One class of critical band scales is called **Bark frequency scale** (24 critical bands)

# Speech Perception

## *Frequency Analysis (cont.)*

- **Bark Scale:** (cont.)
  - Treat spectral energy over the Bark scale, a more natural fit with spectral information processing in the ear can be achieved
  - The perceptual **resolution** (解析度) is finer in the lower frequencies
  - The critical bands are continuous such that a tone of any audible frequency always finds a critical band centered on it

$$b(f) = 13 \arctan(0.00076 f) + 3.5 \arctan\left[\left(\frac{f}{7500}\right)^2\right]$$

# Speech Perception

## Frequency Analysis (cont.)

- Bark Scale: (cont.)**

Table 2.3 The Bark frequency scale.

Bark Band #	Edge (Hz)	Center (Hz)
1	100	50
2	200	150
3	300	250
4	400	350
5	510	450
6	630	570
7	770	700
8	920	840
9	1080	1000
10	1270	1170
11	1480	1370
12	1720	1600
13	2000	1850
14	2320	2150
15	2700	2500
16	3150	2900
17	3700	3400
18	4400	4000
19	5300	4800
20	6400	5800
21	7700	7000
22	9500	8500
23	12000	10500
24	15500	13500

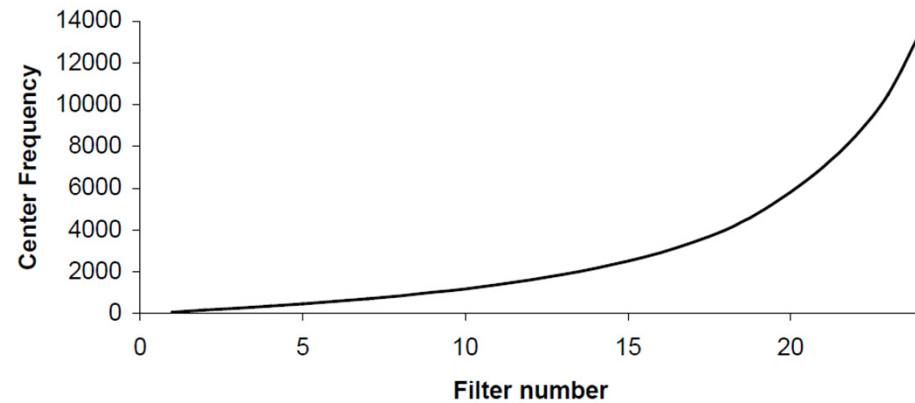
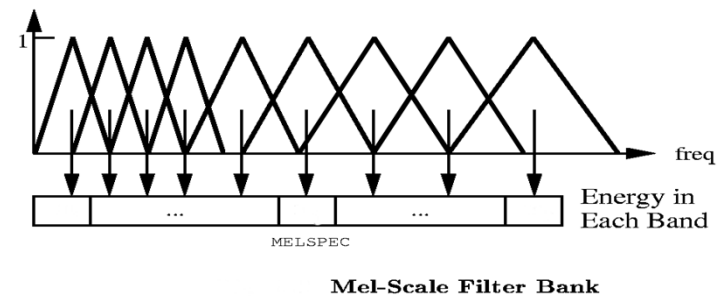


Figure 2.12 The center frequency of 24 Bark frequency filters as illustrated in Table 2.3.



# Speech Perception

## *Frequency Analysis (cont.)*

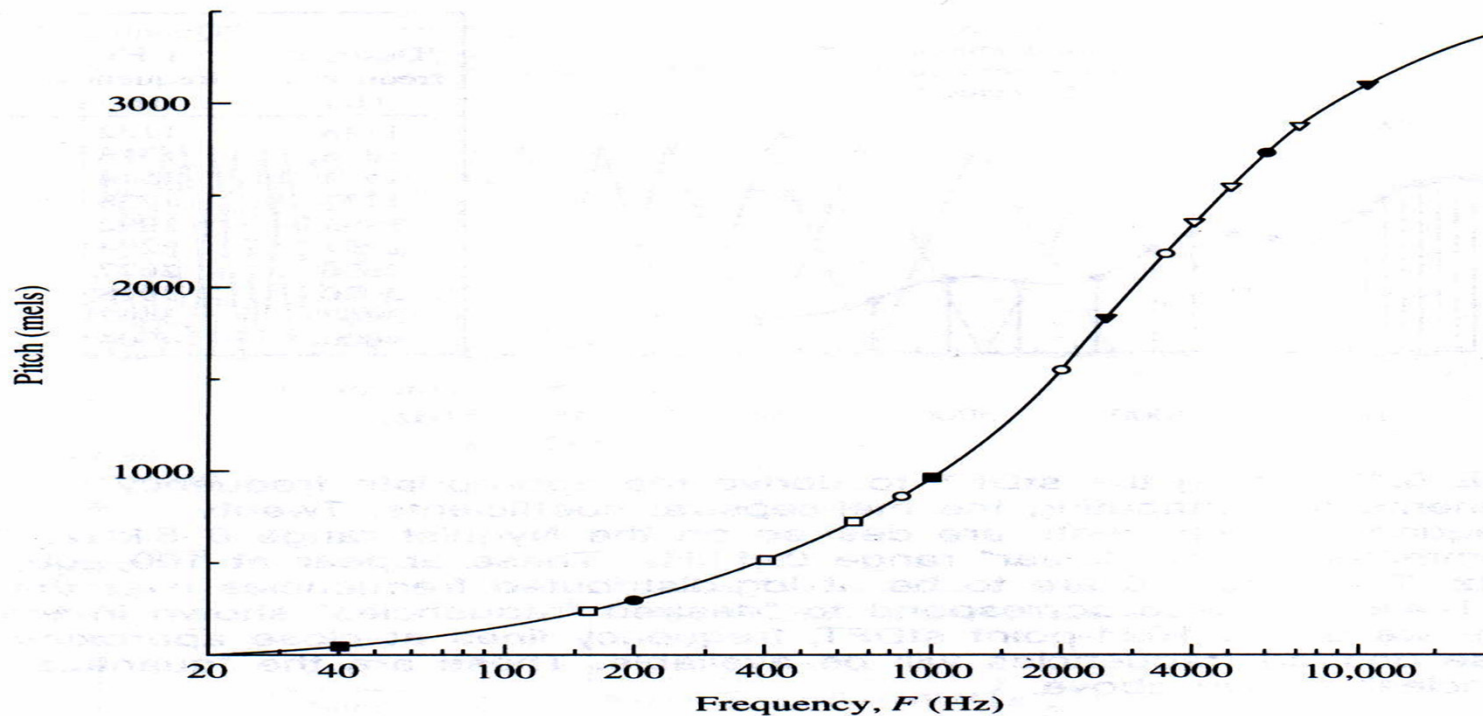
- **Mel Frequency Scale (Mel): linear below 1 KHz and logarithmic above**
  - Model the sensitivity of the human ear
  - **Mel**: a unit of measure of perceived pitch or frequency of a tone
- **Steven and Volkman (1940)**
  - Arbitrarily chose the frequency 1,000 Hz as “1,000 mels”.
  - Listeners were then asked to change the physical frequency until the pitch they perceived was twice the reference, then 10 times, and so on; and then half the reference, 1/10, and so on
    - These pitches were labeled 2,000, 10,000 mels and so on; and 500 and 100 mels, and so on
  - Determine a mapping between the real frequency scale (Hz) and the perceptual frequency (Mel)
  - Have been widely used in modern speech recognition system

# Speech Perception

## *Frequency Analysis (cont.)*

- **Mel Frequency Scale (cont.)**

$$\text{Mel}(f) = 1125 \ln\left(1 + \frac{f}{700}\right)$$

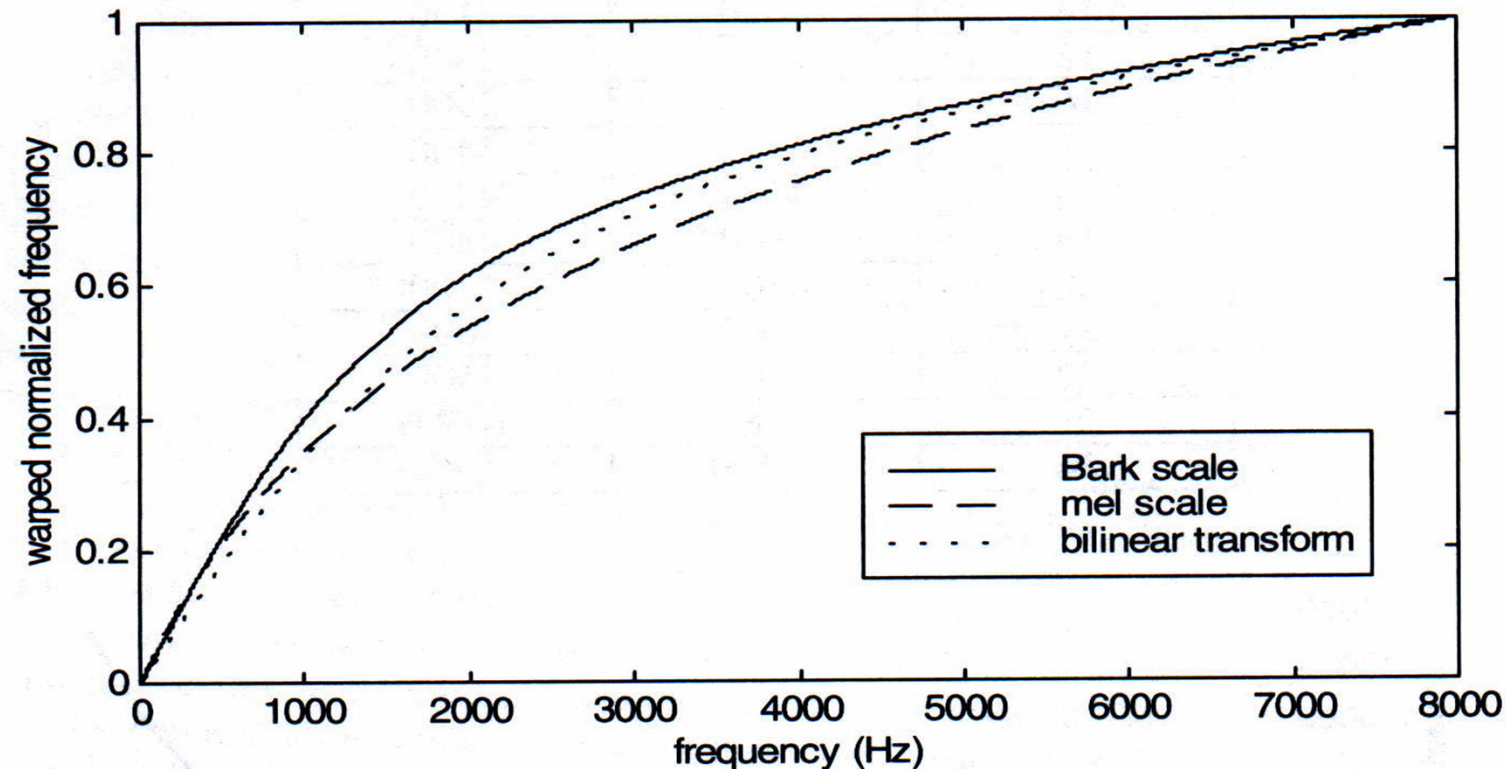


**FIGURE 6.12.** The mel scale. After Stevens and Volkman (1940).



# Speech Perception

## *Frequency Analysis (cont.)*



**Figure 2.13** Frequency warping according to the Bark scale, ERB scale, mel-scale, and bilinear transform for  $\alpha = 0.6$ : linear frequency in the  $x$ -axis and normalized frequency in the  $y$ -axis.

# Phonetics and Phonology

- Phonetics (語音學): The study of speech sounds and their production, classification, and transcription
- Phonology (音韻學): The study of the distribution and patterning of speech sounds in a language and of the tacit rules governing the speech pronunciation

# Phoneme and Phone

- **Phoneme and Phone**

- In speech science, the term *phoneme* (音素/音位) is used to denote any of the *minimal units of speech sound* in a language that can serve to distinguish one word from another
  - E.g., m*ea*n /iy/ and ma*n* /ae/
- The term *phone* is used to denote a phoneme's acoustic realization
  - E.g., phoneme /t/ has two very different acoustic realizations in the word *sat* and *meter*. We had better treat them as **two different phones** when building a spoken language system
  - E.g., phoneme // : *like* and *sail*

# Phoneme and Phone

- Phoneme and phone interchangeably used to refer to the speaker-independent and context-independent units of meaningful sound contrast
  - The set of phonemes will differ in realization across individual speakers

**Table 2.4** English phonemes used for typical spoken language systems.

Phonemes	Word Examples	Description
<i>iy</i>	<i>feel, eve, me</i>	front close unrounded
<i>ih</i>	<i>fill, hit, lid</i>	front close unrounded (lax)
<i>ae</i>	<i>at, carry, gas</i>	front open unrounded (tense)
<i>aa</i>	<i>father, ah, car</i>	back open unrounded
<i>ah</i>	<i>cut, bud, up</i>	open-mid back unrounded
<i>ao</i>	<i>dog, lawn, caught</i>	open-mid back round
<i>ay</i>	<i>tie, ice, bite</i>	diphthong with quality: aa + ih
<i>ax</i>	<i>ago, comply</i>	central close mid (schwa)
<i>ey</i>	<i>ate, day, tape</i>	front close-mid unrounded (tense)
<i>eh</i>	<i>pet, berry, ten</i>	front open-mid unrounded
<i>er</i>	<i>turn, fur, meter</i>	central open-mid unrounded rhotic
<i>ow</i>	<i>go, own, tone</i>	back close-mid rounded
<i>aw</i>	<i>foul, how, our</i>	diphthong with quality: aa + uh
<i>oy</i>	<i>toy, coin, oil</i>	diphthong with quality: ao + ih
<i>uh</i>	<i>book, pull, good</i>	back close-mid unrounded (lax)
<i>uw</i>	<i>tool, crew, moo</i>	back close round
<i>b</i>	<i>big, able, tab</i>	voiced bilabial plosive
<i>p</i>	<i>put, open, tap</i>	voiceless bilabial plosive
<i>d</i>	<i>dig, idea, wad</i>	voiced alveolar plosive
<i>t</i>	<i>talk, sat</i>	voiceless alveolar plosive & alveolar flap
<i>g</i>	<i>meter</i>	voiced velar plosive
<i>k</i>	<i>gut, angle, tag</i>	voiceless velar plosive
<i>f</i>	<i>cut, ken, take</i>	voiceless labiodental fricative
<i>v</i>	<i>fork, after, if</i>	voiced labiodental fricative
<i>s</i>	<i>vat, over, have</i>	voiceless alveolar fricative
<i>z</i>	<i>sit, cast, toss</i>	voiced alveolar fricative
<i>th</i>	<i>zap, lazy, haze</i>	voiceless dental fricative
<i>dh</i>	<i>thin, nothing, truth</i>	voiced dental fricative
<i>sh</i>	<i>then, father, scythe</i>	voiceless postalveolar fricative
<i>zh</i>	<i>she, cushion, wash</i>	voiced postalveolar fricative
<i>l</i>	<i>genre, azure</i>	alveolar lateral approximant
<i>l</i>	<i>lid</i>	velar lateral approximant
<i>r</i>	<i>elbow, sail</i>	retroflex approximant
<i>y</i>	<i>red, part, far</i>	palatal sonorant glide
<i>w</i>	<i>yacht, yard</i>	labiovelar sonorant glide
<i>hh</i>	<i>with, away</i>	voiceless glottal fricative
<i>m</i>	<i>help, ahead, hotel</i>	bilabial nasal
<i>n</i>	<i>mat, amid, aim</i>	alveolar nasal
<i>ng</i>	<i>no, end, pan</i>	velar nasal
<i>ch</i>	<i>sing, anger</i>	voiceless alveolar affricate: t + sh
<i>jh</i>	<i>chin, archer, march</i>	voiced alveolar affricate: d + zh
	<i>joy, agile, edge</i>	

# Vowels

- The tongue shape and positioning on the oral cavity do not form a major constriction (壓縮) of air flow during vowel articulation
  - Variations of tongue placement give each vowel its distinct character by changing the resonances (the positions of formants)
    - Just as different sizes and shapes of bottles give rise to different acoustic effects when struck
  - The linguistically important dimensions of the tongue movements are generally the ranges [front <-> back] and [high <-> low]
- F1 and F2
  - The primary energy entering the pharyngeal (咽) and oral (口腔) cavities in vowel production vibrates at the fundamental frequency. The major resonances of the oral and pharyngeal cavities for vowels are called F1 and F2

# Vowels (cont.)

- F1 and F2 (cont.)
  - The major resonances of these two cavities for vowels are called **F1** and **F2**, the first and second formants
    - **Determined by** the tongue placement and oral tract shape in vowels
    - **Determine** the characteristic timbre or quality of the vowel
  - English vowels can be described by the relationship of F1 and F2 to one another
  - F2 is determined by **the size of the and shape of the oral portion**, forward of the major tongue extrusion(擠壓)
  - F1 corresponds to **the back or pharyngeal portion** of the cavity (the cavity from the glottis (聲門) to the tongue extrusion), which is longer than the forward part. Its resonance would be lower
  - **Rounding the lips** has the effect of extending the front-of-tongue cavity, thus lowering F2

# Vowels (cont.)

- The characteristic F1 and F2 values are ideal locations for perception

**Table 2.5** Phoneme labels and typical formant values for vowels of English.

<b>Vowel Labels</b>	<b>Mean F1 (Hz)</b>	<b>Mean F2 (Hz)</b>
<i>iy (feel)</i>	300	2300
<i>ih (fill)</i>	360	2100
<i>ae (gas)</i>	750	1750
<i>aa (father)</i>	680	1100
<i>ah (cut)</i>	720	1240
<i>ao (dog)</i>	600	900
<i>ax (comply)</i>	720	1240
<i>eh (pet)</i>	570	1970
<i>er (turn)</i>	580	1380
<i>ow (tone)</i>	600	900
<i>uh (good)</i>	380	950
<i>uw (tool)</i>	300	940

嘴唇愈成圓形或愈開

## Vowels (cont.)

- The **tongue** hump (彎曲、隆起) is the major actor in vowel articulation. The most important secondary vowel mechanism for English and many other language is **lip** rounding
- E.g. /iy/ (*see*) and /uw/ (*blue*)
  - When you say /iy/, your tongue will be in the high/front position and your lips will be flat, slightly open, and somewhat spread
    - Lower F1 and Higher F2
  - When you say /uw/, your tongue will be in the high/back position and your lips begin to round out, ending in a more puckered (縮攏的) position
    - Higher F1 and Lower F2



# Vowels (cont.)

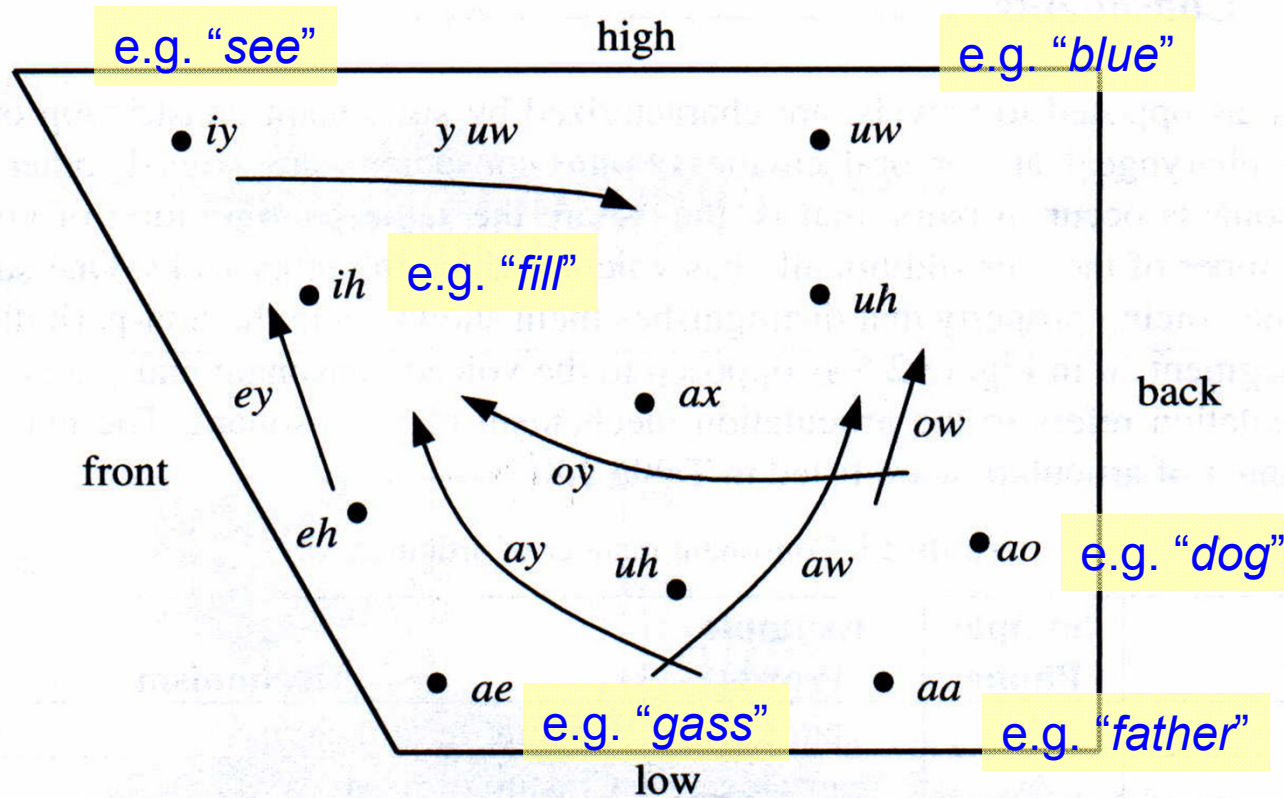


Figure 2.18 Relative tongue positions of English vowels [24].

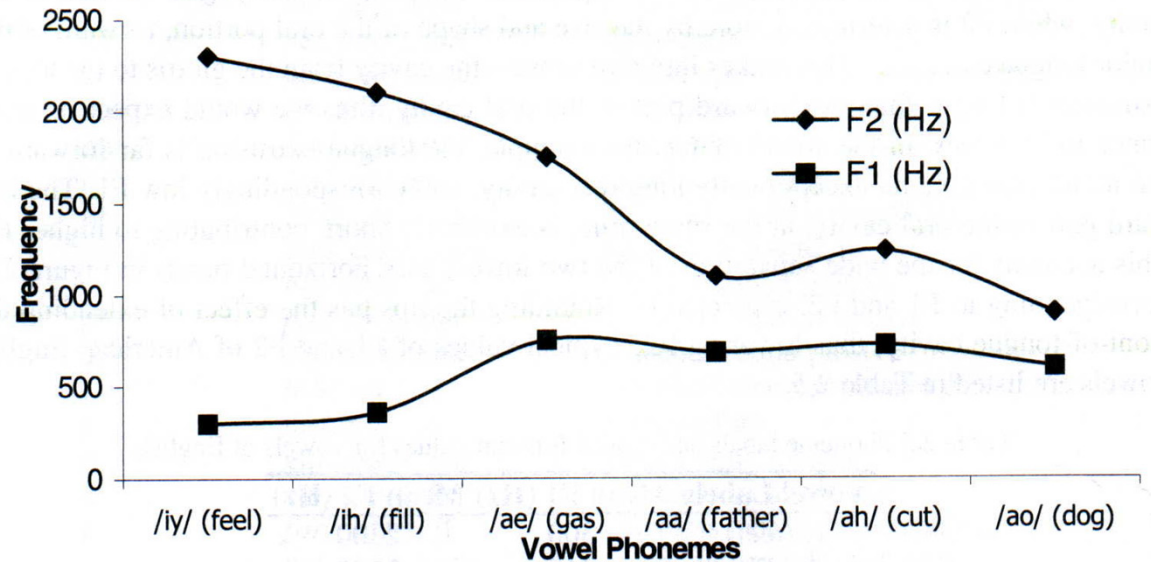
## Vowels (cont.)

- Diphthongs(雙母音)
  - A special class of vowels that combine two distinct sets of F1/F2 values

**Table 2.6** The diphthongs of English.

<b>Diphthong Labels</b>	<b>Components</b>
ay (tie)	/aa/ → /iy/
ey (ate)	/eh/ → /iy/
oy (coin)	/ao/ → /iy/
aw (foul)	/aa/ → /uw/

## Vowels (cont.)



**Figure 2.17** F1 and F2 values for articulations of some English vowels.

The major articulator for English vowels is the middle to rear portion of the tongue.

- Note: not only tongue hump (彎曲、隆起) but also lip rounding is the two major actor in vowel articulation for most languages

# Consonants

- Characterized by significant constriction (壓縮) or obstruction (阻塞) in the pharyngeal and/or oral cavities
  - Some consonants are voiced; others are not
  - Many consonants occur in pairs, i.e., sharing the same configuration of articulators and one member of the pair additionally has voicing while the other lacks (e.g. /z, s/)

Table 2.8 Consonant manner of articulation.

Manner	Sample Phone	Example Words	Mechanism
Plosive	/p/	tat, tap	Closure in oral cavity
Nasal	/m/	team, meet	Closure of nasal cavity
Fricative	/s/	sick, kiss	Turbulent airstream noise
Retroflex liquid	/r/	rat, tar	Vowel-like, tongue high and curled back
Lateral liquid	/l/	lean, kneel	Vowel-like, tongue central, side airstream
Glide	/y/, /w/	yes, well	Vowel-like

破裂音  
鼻音  
摩擦音  
捲舌音  
舌邊音  
滑音

# Consonants (cont.)

- Plosives (破裂音)
  - E.g., /b, p/, /d, t/, /g, k/
  - Consonant that involve complete blockage of oral cavity
- Fricatives (摩擦音)
  - E.g., /z, s/
  - Consonants that involve nearly complete blockage of oral cavity
- Nasals (鼻音)
  - E.g., /m, n, ŋ/
  - Consonants that let the oral cavity significantly constricted, velar (軟顎) open, voicing and air pass through the nasal cavity
- Retroflex liquids (捲舌音)
  - E.g., /r/
  - The tip of the tongue is curled back slightly

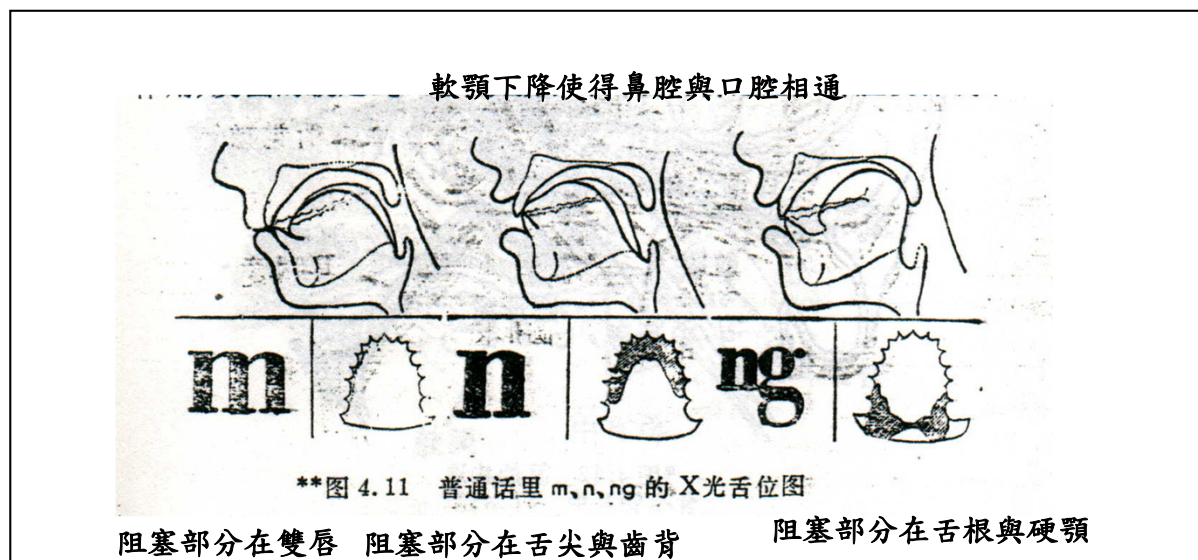
# Consonants (cont.)

- Lateral liquids (舌邊音)
  - E.g., /l/
  - Air stream flows around the sides of the tongue
- Glides (滑音)
  - E.g. /y, w/
  - Be a little shorted and lack the ability to be stressed, usually at the initial position within a syllable (e.g., **y**es, **w**ell)

# Consonants (cont.)

- Semi-vowels
  - Have voicing without complete constriction or obstruction of the vocal tract
  - Include the liquid group /r, l/ and glide group /y, w/
  - {vowels, semi-vowels}: sonorant (響音)
- Non-sonorant consonants
  - Maintain some voicing before or during the obstruction until the pressure differential across the glottis (聲門) to disappear, due to the closure 帶聲的子音
  - E.g., /b, d, g, z, zh, v/ (voicing) and their counterparts  
/p, t, k, s, sh, f/ 不帶聲的子音

# Consonants (cont.)





# Phonetic Typology (語音的類型)

- **Length:** Japanese vowels have a characteristic distinction of the length that can be hard for non-natives to perceive and use when learning the language
  - The word *kado* (corner) and *kaado* (card) are spectrally identical, differing in their durations
  - Length is phonemically distinctive for Japanese
- **Pitch:**
  - The primary dimension lacks in English
  - Many Asia and Africa language are tonal
    - E.g. Chinese
  - For tonal language, they have lexical meaning contrasts cued by pitch
    - E.g. Mandarin Chinese has four primary tones

# Phonetic Typology (cont.)

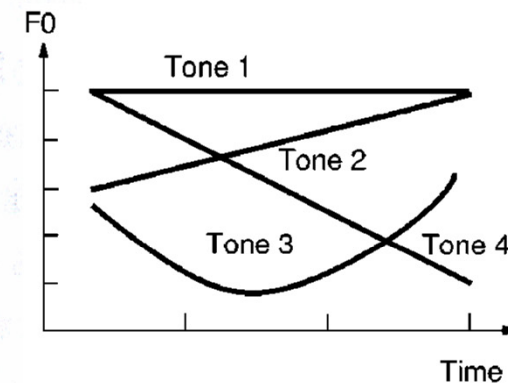
- **Pitch:** (cont.)

- Though English don't make systematic use of pitch in its inventory of word contrasts, we always see with any possible phonetic effect:

- Pitch is systematically viewed in English to signal a speaker's emotions, intentions and attitudes
- Pitch has some linguistic function in signaling grammatical structure as well

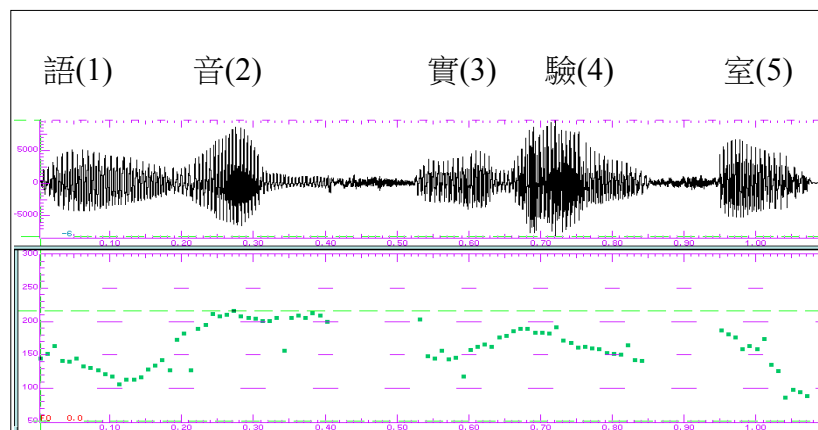
**Table 2.11** The contrastive tones of Mandarin Chinese.

Tone	Shape	Example	Chinese	Meaning
1	High level	<i>ma</i>	妈	mother
2	High rising	<i>ma</i>	麻	numb
3	Low rising	<i>ma</i>	马	horse
4	High falling	<i>ma</i>	骂	to scold



**Figure 1:** Pitch patterns of four lexical tones.

# Phonetic Typology (cont.)

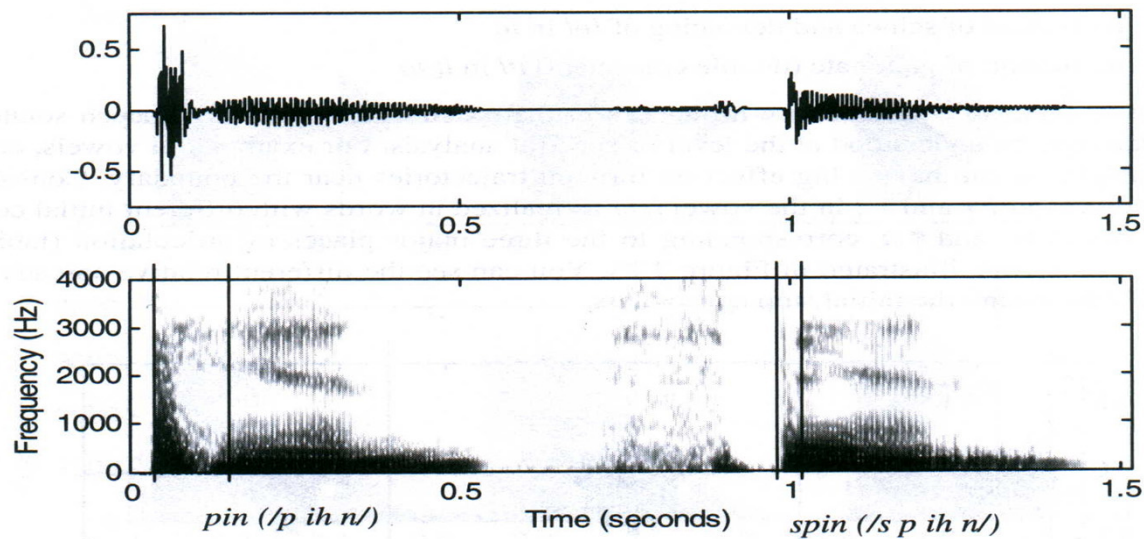


	Tone 1	Tone 2	Tone 3	Tone 4	neutral tone
number of models	4	6	6	4	3
typical tone	1	2	3	4	5
concatenation	1-(2)	2-(2)	3-(1)	4-(1)	(1)-5
combinations	(3)-1	(1)-2	(1)-3	(3)-4	(3)-5
	(3)-1-(2)	(1)-2-(2)	(1)-3-(1)	(3)-4-(1)	
		(3)-2	(3)-3		
		(3)-2-(2)	(3)-3-(1)		

# The Allophone: Sound and Context

- Phonetic units should be correlated with potential meaning distinctions
  - *mean* /m iy n/ and *men* /m eh n/
- However, the fundamental meaning-distinguishing sound is often modified in some systematic way by its phonetic neighbors
  - **Coarticulation**: the process by which the neighbor sounds influence one another
  - **Allophone**: when the variations resulting from coarticulatory processes **can be consciously perceived, the modified phonemes** are called allophones
  - E.g. :
    - *p* in (*pin*, /p ih n/) produces a notice puff (噴出) of air, called aspiration (送氣), but loses its aspiration in (*spin*, /s p ih n/)
    - A vowel before a voicing consonant, .e.g., *bad* /d/, seems typically longer than the same vowel before the unvoiced counterpart, in this case *bat* /t/

# The Allophone: Sound and Context (cont.)



**Figure 2.22** Spectrogram: bursts of *pin* and *spin*. The relative duration of a *p*-burst in different phonetic contexts is shown by the differing width of the area between the vertical lines.

# Structural Features of Chinese Language

- Not Alphabetic (字母的)
- At Least 10,000 Commonly Used Characters (字)
  - Almost all morphemes (詞素) with their own meaning
  - All monosyllabic
- Unlimited Number of Words (詞), at Least 100,000 Commonly Used, Each Composed of One to Several Characters (字)
  - The meaning of the word can be directly or partly related, or even completely irrelevant to the meaning of the component characters  
書店, 大學, 和尚, 光棍
- Chinese is a Tonal Language
  - 4 lexical tones, 1 neutral tone (the number is for Mandarin)

# Structural Features of Chinese Language (cont.)

- About 1,335 Syllables Only (the number is for Mandarin)
  - About 408 base-syllables if differences in tone disregarded (the number is for Mandarin)
- Large Number of Homonym Characters (同音字) Sharing the Same Syllable
- Monosyllabic Structure of Chinese Language
  - Each syllable stands for many characters with different meaning
  - Combination of syllables (characters) gives unlimited number of words
  - Small number of syllables carries plurality (多重性) of linguistic information
- Almost Each Character with Its Own Meaning, thus Playing Some Linguistic Role Independently

# Structural Features of Chinese Language (cont.)

- No Natural Word Boundaries in a Chinese Sentence

電腦科技的進步改變了人類的生活和工作方式

=====

- Word segmentation not unique
- Words not well defined
- Commonly accepted lexicon not existing
- Open Vocabulary Nature with Flexible Wording Structure

- New words easily created everyday

電 (electricity) + 腦 (brain) → 電腦 (computer)

- Long word arbitrarily abbreviated

臺灣大學 (Taiwan University) → 臺大

- Name/title

李登輝總統 (President T.H. Lee) → 李總統登輝

- Unlimited number of compound words

高 (high) + 速 (speed) + 公路 (highway) → 高速公路 (freeway)



# Structural Features of Chinese Language (cont.)

- Difficult for Word-based Approaches Popularly Used in Alphabetic Languages
  - Serious out of vocabulary (OOV) problem
- Considering Phonetic Structure of Mandarin Syllables
  - INITIAL / FINAL's
  - Phone-like-units / phonemes
- Different Degrees of Context Dependency
  - Intra-syllable only
  - Intra-syllable plus inter-syllable
  - Right context dependent only
  - Both right and left context dependent

# Structural Features of Chinese Language (cont.)

- Examples
  - 22 INITIAL's extended to 113 right-context-dependent INITIAL's
  - 33 phone-like-units extended to 145 intra-syllable right-context-dependent phone-like-units, or 481 with both intra/inter-syllable context dependency
  - 4,606 triphones with intra/inter-syllable context dependency

Syllables (1,345)				Tones (4+1)
Base-syllables (408)				
<b>INITIAL's</b> (21)	<b>FINAL's</b> (37)			
	Medials (3)	Nucleus (9)	Ending (2)	
Consonants (21)	Vowels plus Nasals (12)			
Phones (31)				





# Explanations

- 首先要整理自己的思想，決定要說的訊息內容
- 把它們變為適當的語言形式，選擇適當的詞彙，按照某種語言的法則，組成詞句，以表達想說的訊息內容 (遣詞造句)
- 以生理神經式衝動的形式，言運動神經傳播到聲帶、舌唇等器官的肌肉，驅動這些肌肉運動
- 空氣發生壓力變化，經過聲腔的調節，從而產生出通常的語言聲波

# Explanations for Speech Production

## 人的發音器官可分三大部分

- **動力器官：肺和氣管等呼吸器官**
  - 我們大約每五秒呼吸一次，說話是在呼氣的過程中進行
  - 利用肺部呼出的氣流作為動力來激勵聲帶振動
- **發聲器官：聲帶、喉頭及一些軟骨組織等**
  - 來自肺部的穩定氣流由於喉頭的開關節制動作，因此被改變，成為聽得見的、像蜂鳴一樣的聲音。
  - 喉頭的節制動作主要依賴聲帶來完成的。聲帶是發聲體本身，為語音提供主要的聲源。聲帶振動產生的一系列的脈衝(impulses)，是一種週期波，其頻譜含有大量的諧波(harmonics)成分，它們的頻率是基頻 (fundamental frequency) 的整數倍

# Explanations for Speech Production (cont.)

## 人的發音器官可分三大部分 (cont.)

- **共鳴(共振)調節器官:口腔、鼻腔、咽腔 (統稱”聲腔” , vocal tract)**
  - 聲腔是充滿氣體的管腔，具有一定的自然頻率。當來自聲帶的脈衝之某一諧波與聲腔的某一自然頻率相同或相近時，就發生共鳴(resonance)現象，此一脈衝諧波頻率成分被加強而提起。因此，從口中輻射出的語音的頻譜在聲腔的自然頻率處就有共振峰(Formants)，它們的頻率叫做共振峰頻率
  - 發音(articulation)機制、調音機制:指聲腔對於聲帶產生聲音的共鳴和調節作用，它與語音的音色關係極為密切
  - 聲腔變化主要是由舌的高低前後所造成的，像語音學(phonetics)常用的母音舌位圖
  - 雙唇與牙齒是唯一從外部看得見的發音器官，可以額外地為人提供許多語言交際的信息

# Explanations for Speech Production (cont.)

- 聲腔在發母音(vowel)與發子音(consonant)時的表現
  - 發母音時聲腔裡沒有阻塞，但發子音時，聲腔的某兩個部位必定構成阻塞、阻礙，然後突然釋放被阻空氣，氣流通過從狹縫洩出或突然衝出，從而形成噪音
  - 子音的音色跟聲腔阻塞部分的不同和解除的方式的不同有直接相關



# Explanations for Speech Perception

- 聽力形成：
  1. 聲音由耳翼(pinna)接收，並傳至外耳道再傳至耳膜(eardrum)
  2. 耳膜接收聲音的能量，並將它轉變成機械能量，所以第一個能量的轉換是從耳膜開始
  3. 耳膜再把機械能量，傳送到聽小骨鏈
  4. 鐙骨(stapes)的踏板接在卵圓窗上面，它將機械能再轉成液能，這裏是第二個能量轉換處
  5. 前庭階的能量會傳遞到中階，中階液體的移動，會造成柯氏器上面毛髮細胞的移動
  6. 中階再將液能轉為電能量，此為第三個能量轉換處。
  7. 毛髮細胞會刺激在柯氏器基部的神經細胞，再將這些神經訊號經由聽神經傳到腦部
  8. 能源轉換結論：外耳(聲能) → 中耳(機械能) → 內耳(液能及電能)

# Speech Perception

## *Physiology of the Ear (cont.)*

- 外耳：
  - 耳道：是一個充滿氣體的管子，是一種共鳴器，當傳入聲波的某些頻率接近它的一套自然頻率時，就被放大的約二至四倍
- 中耳：
  - 三小聽骨：錘骨、鈹骨、蹬骨。錘骨與鼓膜相連，蹬骨與覆蓋著卵圓窗 (oval window)
  - 兩種主要功能：
    - 放大作用，以提高傳入內耳的聲音能量(槓桿原理)
    - 保會內耳免受特強音的損害
- 內耳：
  - 耳蝸：充滿淋巴液，黏度幾乎為水的兩倍，耳蝸隔膜分隔兩區，淋巴液由蝸孔自由流通兩區。耳蝸隔膜內有耳蝸導管，充滿內淋巴液。
    - 基底膜在靠近卵圓窗處，較窄、薄，繃的緊；而靠近蝸孔部分最為寬鬆肥大
    - 基底膜的這種特性，讓其能最傳入聲波不同的頻率產生響應
  - 主要功能：
    - 把外界機械動能轉換成神經衝動

## Consonants (cont.)

- 最後再看嘴唇、舌頭跟口腔的一些關係
  - 閉唇 (labial): /p/, /b/, /m/, /w/
  - 舌被齒或齒與唇夾 (dental or labio-dental consonants): /f/, /v/, /th/, /dh/
  - 舌頭前端碰齒槽 (alveolar consonants): /t/, /d/, /n/, /s/, /z/, /r/, /l/
  - 舌頭前端碰上顎 (palatal consonants): /sh, zh, y/
  - 舌頭後端碰軟顎 (velar consonants): /k/, /g/, /ng/