Chapter 4

Hearing, Auditory Models, and Speech Perception

听觉,听觉模型与语音感知

Topics to be Covered

- The Speech Chain (语音链) Production and Human Perception
- Auditory mechanisms (听觉机理)— the human ear and how it converts sound to auditory representations
- Speech perception (语音感知) and what we know about physical and psychophysical measures of sound
- Auditory masking (听觉掩蔽)
- Sound and word perception in noise

Auditory Mechanisms

Speech Perception

- understanding how we hear sounds and how we perceive speech leads to better design and implementation of robust and efficient systems for analyzing and representing speech
- the better we understand signal processing in the human auditory system, the better we can (at least in theory) design practical speech processing systems
 - speech and audio coding (MP3 audio, cellphone speech)
 - speech recognition
- try to understand speech perception by looking at the physiological models of hearing



- The Speech Chain comprises the processes of:
 - speech production,
 - auditory feedback to the speaker,
 - speech transmission (through air or over an electronic communication system) to the listener, and
 - speech perception and understanding by the listener.

The Speech Chain

- The message to be conveyed by speech goes through five levels of representation between the speaker and the listener, namely:
 - the linguistic level (where the basic sounds of the communication are chosen to express some thought of idea)
 - the physiological level (where the vocal tract components produce the sounds associated with the linguistic units of the utterance)
 - the acoustic level (where sound is released from the lips and nostrils and transmitted to both the speaker (sound feedback) and to the listener)
 - the physiological level (where the sound is analyzed by the ear and the auditory nerves), and finally
 - the linguistic level (where the speech is perceived as a sequence of linguistic units and understood in terms of the ideas being communicated)



Auditory System

- the acoustic signal first converted to a neural representation by processing in the ear
 - the conversion takes place in stages at the outer, middle and inner ear
 - these processes can be measured and quantified
- the neural transduction step takes place between the output of the inner ear and the neural pathways to the brain
 - consists of a statistical process of nerve firings at the hair cells of the inner ear, which are transmitted along the auditory nerve to the brain
 - much remains to be learned about this process
- the nerve firing signals along the auditory nerve are processed by the brain to create the perceived sound corresponding to the spoken utterance
 - these processes not yet understood

The McGurk Effect



The Black Box Model of the Auditory System

- researchers have resorted to a "black box" behavioral model of hearing and perception
 - model assumes that an acoustic signal enters the auditory system causing behavior that we record as psychophysical (精神物理学) observations
 - psychophysical methods and sound perception experiments determine how the brain processes signals with different loudness levels, different spectral characteristics, and different temporal properties
 - characteristics of the physical sound are varied in a systematic manner and the psychophysical observations of the human listener are recorded and correlated with the physical attributes of the incoming sound
 - we then determine how various attributes of sound (or speech) are processed by the auditory system



The Black Box Model Examples

Physical Attribute	Psychophysical Observation
Intensity 强度	Loudness 响度
Frequency 频率	Pitch 音高

- Experiments with the "black box" model show:
 - correspondences between sound intensity and loudness, and between frequency and pitch are complicated and far from linear
 - attempts to extrapolate from psychophysical measurements to the processes of speech perception and language understanding are, at best, highly susceptible to misunderstanding of exactly what is going on in the brain

Overview of Auditory Mechanism



 begin by looking at ear models including processing in cochlea (耳蜗)

The Human Ear



- Outer ear (外耳): pinna (耳廓) and external canal
- Middle ear (中耳): tympanic membrane (鼓膜) or eardrum
- Inner ear (内耳): cochlea(耳蜗), neural connections

Human Ear

- Outer ear: funnels (使经过漏斗) sound into ear canal
- Middle ear: sound impinges (撞击) on tympanic membrane; this causes motion
 - middle ear is a mechanical transducer, consisting of the hammer (锤骨), anvil (砧骨) and stirrup (镫骨); it converts acoustical sound wave to mechanical vibrations along the inner ear
- Inner ear: the cochlea is a fluid-filled chamber partitioned by the basilar membrane (基底膜)
 - the auditory nerve is connected to the basilar membrane via inner hair cells
 - mechanical vibrations at the entrance to the cochlea create standing waves (of fluid inside the cochlea) causing basilar membrane to vibrate at frequencies commensurate with the input acoustic wave frequencies (formants) and at a place along the basilar membrane that is associated with these frequencies

The Outer Ear





The Middle Ear



 The Hammer (锤骨), Anvil (砧骨) and Stirrup (镫骨) are the three tiniest bones in the body. Together they form the coupling between the vibration of the eardrum and the forces exerted on the oval window (卵圆窗) of the inner ear.

These bones can be thought of as a compound lever which achieves a multiplication of force—by a factor of about three under optimum conditions. (They also protect the ear against loud sounds by attenuating the sound.)

Transfer Functions at the Periphery



The Inner Ear



- The inner ear can be thought of as two organs, namely
 - the semicircular canals which serve as the body's balance organ and
 - the cochlea which serves as the body's microphone, converting sound pressure signals from the outer ear into electrical impulses which are passed on to the brain via the auditory nerve.



Taking electrical impulses from the cochlea and the semicircular canals, the auditory nerve makes connections with both auditory areas of the brain.

Stretched Cochlea & Basilar Membrane



Basilar Membrane Mechanics

- characterized by a set of frequency responses at different points along the membrane
- mechanical realization of a bank of filters
- filters are roughly constant Q (center frequency/bandwidth) with logarithmically decreasing bandwidth
- distributed along the Basilar Membrane is a set of about 3000 sensors, called Inner Hair Cells (IHC), which act as mechanical motion-to-neural activity converters
- mechanical motion along the BM is sensed by local IHC causing firing activity at nerve fibers that innervate bottom of each IHC
- each IHC connected to about 10 nerve fibers, each of different diameter => thin fibers fire at high motion levels, thick fibers fire at lower motion levels
- 30,000 nerve fibers link IHC to auditory nerve
- electrical pulses run along auditory nerve, ultimately reach higher levels of auditory processing in brain, perceived as sound

Basilar Membrane Mechanics











Speech Perception

The Perception of Sound

- Key questions about sound perception:
 - what is the `resolving power' of the hearing mechanism
 - how good an estimate of the fundamental frequency of a sound do we need so that the perception mechanism basically `can't tell the difference'
 - how good an estimate of the resonances or formants (both center frequency and bandwidth) of a sound do we need so that when we synthesize the sound, the listener can't tell the difference
 - how good an estimate of the intensity of a sound do we need so that when we synthesize it, the level appears to be correct

Sound Intensity

- Intensity (音强) of a sound is a physical quantity that can be measured and quantified
- Acoustic Intensity (I) defined as the average flow of energy (power) through a unit area, measured in watts/square meter
- Range of intensities between 10⁻¹² watts/square meter to 10 watts/square meter; this corresponds to the range from the threshold of hearing to the threshold of pain

Threshold of hearing defined to be:

 $I_0 = 10^{-12} \text{ watts/m}^2$

The intensity level of a sound, IL is defined relative to I_0 as:

$$IL = 10 \log_{10} \left(\frac{I}{I_0} \right)$$
 in dB

For a pure sinusoidal sound wave of amplitude P, the intensity

is proportional to P^2 and the sound pressure level (SPL) is defined as:

$$SPL = 10 \log_{10} \left(\frac{P^2}{P_0^2} \right) = 20 \log_{10} \left(\frac{P}{P_0} \right) dB$$

where $P_0 = 2 \times 10^{-5}$ Newtons/m²

Some Facts About Human Hearing

- the range of human hearing is incredible
 - threshold of hearing thermal limit of Brownian motion of air particles in the inner ear
 - threshold of pain intensities of from 10^12 to 10^16 greater than the threshold of hearing
- human hearing perceives both sound frequency and sound direction
 - can detect weak spectral components in strong broadband noise
- masking is the phenomenon whereby one loud sound
 - makes another softer sound inaudible
 - masking is most effective for frequencies around the masker frequency
 - masking is used to hide quantization noise

Anechoic Chamber (no Echos)



Anechoic Chamber (no Echos)





Sound Pressure Levels (dB)

SPL (dB)—Sound Source 160 Jet Engine — close up 150 Firecracker; Artillery Fire 140 Rock Singer Screaming into Microphone: Jet Takeoff 130 Threshold of Pain; .22 Caliber Rifle 120 Planes on Airport Runway; Rock Concert; Thunder 110 Power Tools; Shouting in Ear 100 Subway Trains; Garbage Truck 90 Heavy Truck Traffic; Lawn Mower 80 Home Stereo — 1 foot; Blow Dryer

SPL (dB)—Sound Source

70	Busy Street; Noisy Restaurant
60	Conversational Speech — 1 foot
50	Average Office Noise; Light Traffic; Rainfall
40	Quiet Conversation; Refrigerator; Library
30	Quiet Office; Whisper
20	Quiet Living Room; Rustling Leaves
10	Quiet Recording Studio; Breathing
0	Threshold of Hearing

Range of Human Hearing



Hearing Thresholds 听阈

- Threshold of Audibility is the acoustic intensity level of a pure tone that can barely be heard at a particular frequency
 - threshold of audibility \approx 0 dB at 1000 Hz
- Thresholds vary with frequency and from person-to-person
- Maximum sensitivity is at about 3000 Hz

Loudness Level

• Loudness Level (响度级 LL) is equal to the IL of a 1000 Hz tone that is judged by the average observer to be equally loud as the tone



Loudness

 Loudness (L) (in sones 未) is a scale that doubles whenever the perceived loudness doubles



$$L = 2^{(LL-40)/10}$$

Log L = 0.03(LL-40)
= 0.03LL-1.2

• for a frequency of 1000Hz, the loudness level, *LL*, in phons is, by definition, numerically equal to the intensity level *IL* in decibels, so that the equation may be rewritten as

 $LL = 10\log(I/I_0)$ Or since $I_0 = 10^{-12}$ watts/m² $LL = 10\log/+120$ Substitution of this value of LL in the equation gives $\log L = 0.03(10\log/+120)-1.2$ $= 0.3 \log/+2.4$ Which reduces to $L = 251I^{0.3}$

Pitch

- Pitch(音高) and fundamental frequency(基频) are not the same thing
- we are quite sensitive to changes in pitch
 - F < 500 Hz, $\Delta F \approx 3 Hz$
 - *F* > 500 Hz, Δ*F*/*F* ≈ 0.003
- relationship between pitch and fundamental frequency is not simple, even for pure tones
 - the tone that has a pitch half as great as the pitch of a 200
 Hz tone has a frequency of about 100 Hz
 - the tone that has a pitch half as great as the pitch of a 5000 Hz tone has a frequency of less than 2000 Hz
- the pitch of complex sounds is an even more complex and interesting phenomenon



FIGURE 4.16

Chart of the subjective pitch (in mels) versus the actual frequency (in Hz) of a pure tone.

Pitch (*mels*) =1127log_e(1+f /700)

Perception of Frequency

• Pure tone

- Pitch is a perceived quantity while frequency is a physical one (cycle per second or Hertz)
- Mel is a scale that doubles whenever the perceived pitch doubles; start with 1000 Hz = 1000 mels, increase frequency of tone until listener perceives twice the pitch (or decrease until half the pitch) and so on to find mel-Hz relationship
- The relationship between pitch and frequency is nonlinear
- Complex sound such as speech
 - Pitch is related to fundamental frequency but not the same as fundamental frequency; the relationship is more complex than pure tones

Auditory Masking

Pure Tone Masking

- Masking is the effect whereby some sounds are made less distinct or even inaudible by the presence of other sounds
- Make threshold measurements in presence of masking tone; plots below show shift of threshold over non-masking thresholds as a function of the level of the tone masker



Auditory Masking



Masking and Critical Bandwidth

 Critical Bandwidth (临界带宽) is the bandwidth of masking noise beyond which further increase in bandwidth has little or no effect on the amount of masking of a pure tone at the center of the band



The noise spectrum used is essentially rectangular, thus the notion of equivalent rectangular bandwidth (ERB)





z	F_1, F_u	F _c	z	ΔF_G	z	F_{l},F_{u}	Fc	z	AFG
Bark	Hz	Hz	Bark	Hz	Bark	Hz	Hz	Bark	Hz
0	0		111-111		12	1720			
		50	0.5	100		1120	1850	12.5	280
1	100			arel dans	13	2000			
2	200	150	1.5	100			2150	13.5	320
2	200	250			14	2320			
3	200	250	2.5	100			2500	14.5	380
5	300	350	25	100	15	2700			
4	400	330	3.5	100	16	21.50	2900	15.5	450
	400	450	4.5	110	16	3150	2400	16.5	550
5	510	450	4.5	110	17	3700	3400	16.5	550
		570	5.5	120	17	3700	4000	17.5	700
6	630			120	18	4400	4000	17.5	700
		700	6.5	140			4800	18.5	900
7	770				19	5300	10000		
		840	7.5	150			5800	19.5	1100
8	920				20	6400			
0	1000	1000	8.5	160			7000	20.5	1300
9	1080	1170	0.5		21	7700			
10	1270	1170	9.5	190	Coldina -	a margarita	8500	21.5	1800
10	1270	1370	10.5	210	22	9500			116
11	1480	1370	10.5	210	22	12 000	10,500	22.5	2500
	1100	1600	11.5	240	23	12,000	12 500	22.5	250
12	1720	1000	11.0	240	24	15 500	13,500	23.3	3500
		1850	12.5	280	24	15,500			

TABLE 4.2 Critical band rate z and lower (F_l) and upper (F_u) frequencies of critical bandwidths ΔF_G centered at F_c . (After Zwicker and Fastl [428].)

Temporal Masking



Exploiting Masking in Coding



Parameter Discrimination

JND – Just Noticeable Difference (最小可觉差) Similar names: differential limen (DL), ...

Parameter	JND/DL
Fundamental Frequency	0.3-0.5%
Formant Frequency	3-5%
Formant bandwidth	20-40%
Overall Intensity	1.5 dB

Auditory Models

Auditory Models

- Auditory models
 - To predict auditory phenomena for speech applications
- Perceptual effects included in most auditory models:
 - spectral analysis on a non-linear frequency scale (usually mel or Bark scale)
 - spectral amplitude compression (dynamic range compression)
 - loudness compression via some logarithmic process
 - decreased sensitivity at lower (and higher) frequencies based on results from equal loudness contours
 - utilization of temporal features based on long spectral integration intervals (syllabic rate processing)
 - auditory masking by tones or noise within a critical frequency band of the tone (or noise)

Perceptual Linear Prediction



Perceptual Linear Prediction

- Included perceptual effects in PLP
 - critical band spectral analysis using a Bark frequency scale with variable bandwidth trapezoidal shaped filters
 - asymmetric auditory filters with a 25 dB/Bark slope at the high frequency cutoff and a 10 dB/Bark slope at the low frequency cutoff
 - use of the equal loudness contour to approximate unequal sensitivity of human hearing to different frequency components of the signal
 - use of the non-linear relationship between sound intensity and perceived loudness using a cubic root compression method on the spectral levels
 - a method of broader than critical band integration of frequency bands based on an autoregressive, all-pole model utilizing a fifth order analysis

Human Speech Perception Experiments

Sound Perception in Noise

	р	t	k	f	θ	\$	ſ	b	d	g	ν	δ	z	3	m	n
p t k	240 1 18	252 3	41 1 219	2 1	1					1						
<i>f</i> ម ទ្	9		1	225 69	24 185	232	236	5 3			2	1				
b d g					1 1			242	213 33	22 203	24	12 3	1 1			
ν δ 23					1			6 1	2	3 4	171 22 1	30 208 7	4 238	244	1	1
m n												1			274	1 252

FIGURE 17.4 Confusion matrix for S/N = +12 dB and a frequency response of 200–6500 Hz. From [13].

Confusions as to sound PLACE, not MANNER

Sound Perception in Noise

	р	t	k	f	θ	s	ſ	Ь	d	g	τ	δ	z	3	m	n
p t k	80 71 66	43 84 76	64 55 107	17 5 12	14 9 8	630	284 4	1 1	1			1 1 1	2			
f ፀ §]	18 19 8 1	12 17 56	9 16 4 3	175 104 23 4	48 65 39 6	11 32 107 29	1 7 45 195	7 5 4	2423	1 5 3	2 6 1	2 4 1	5 3	2		
b d g	1			5	4 2	4	8	136 5 3	10 80 63	9 45 66	47 11 3	16 20 19	6 20 37	26 56		
τδ v3				2	6 1	2 1	1	48 31 7 1	5 20 26	5 17 27 18	145 86 16 3	45 58 28	12 21 94 45	5 44 129		
m n	1				4			4 1	5	2	4	$\frac{1}{7}$	3 1	6	17 4	

FIGURE 17.5 Confusion matrix for S/N = -6 dB and a frequency response of 200–6500 Hz. From [13].

Confusions in both sound PLACE and MANNER

Speech Perception

50% S/N level for correct responses:

- -14 db for digits
- -4 db for major words
- +3 db for nonsense syllables

- Speech Perception depends on multiple factors including the perception of individual sounds (based on distinctive features) and the predictability of the message (think of the message that comes to mind when you hear the preamble 'To be or not to be ...', or 'Four score and seven years ago ...')
- the importance of linguistic and contextual structure cannot be overestimated (e.g., the Shannon Game where you try to predict the next word in a sentence i.e., 'he went to the refrigerator and took out a ...' where words like plum, potato etc are far more likely than words like book, painting etc.)

Word Intelligibility

Fig. 7.22. Effects of vocabulary size upon the intelligibility of monosyllabic words. (After MILLER, HEISE and LICHTEN)

Intelligibility - Diagnostic Rhyme Test (诊断押韵测试)

Voicing	Nasality	Sustenation	Sibilation	Graveness	Compactness	
veal feel bean peen gin chin dint tint zoo sue dune tune vole foal goat coat zed said dense tense vast fast gaff calf vault fault daunt taunt jock chock bond pond	meat beat need deed mitt bit nip dip moot boot news dues moan bone note dote mend bend neck deck mad bad nab dab moss boss gnaw daw mom bomb knock dock	veebeesheetcheatvillbillthicktickfoopoohshoeschoosethosedozethoughdoughthendenfencepencethandanshadchadthongtongshawchawvonbonvoxbox	zeetheecheepkeepjiltgiltsingthingjuicegoosechewcoojoegosoletholejestguestchaircarejabgabsankthankjawsgauzesawthawjotgotchopcop	weedreedpeakteakbiddidfinthinmoonnoonpooltoolbowldoleforethormetnetpenttentbankdankfadthadfoughtthoughtbongdongwadrodpottot	yield wield key tea hit fit gill dill coop poop you rue ghost boast show so keg peg yen wren gat bat shag sag yawl wall caught thought hop fop got dot	

תת	$T_{100} R_d - W_d$	_					
DR	$T = 100 \times \frac{1}{T_d}$	Coder	Rate (kb/s)	Male	Female	All	MOS
R	= right	FS1016	4.8	94.4	89.0	91.7	3.3
W	= wrong	IS54	7.95	95.2	91.4	93.3	3.6
Т	= total	GSM	13	94.7	90.7	92.7	3.6
d	= one of the six	G.728	16	95.1	90.9	93.0	3.9

Quantification of Subjective Quality

Absolute category rating (ACR) – MOS, mean opinion score

Quality description	Rating
Excellent	5
Good	4
Fair	3
Poor	2
Bad	1

Degradation category rating (DCR) – D(egradation)MOS; need to play reference

Quality description	Rating
Degradation not perceived	5
perceived but not annoying	4
slightly annoying	3
annoying	2
very annoying	1

Description	Rating
Much better	3
Better	2
Slightly better	1
About the same	0
Slightly worse	-1
Worse	-2
Much worse	-3

Comparison category rating (CCR) – randomized (A,B) test

MOS (Mean Opinion Scores 平均意见分)

- Why MOS:
 - SNR is just not good enough as a subjective measure for most coders (especially model-based coders where waveform is not preserved inherently)
 - noise is not simple white (uncorrelated) noise
 - error is signal correlated
 - clicks/transients
 - frequency dependent spectrum—not white
 - includes components due to reverberation and echo
 - noise comes from at least two sources, namely quantization and background noise
 - delay due to transmission, block coding, processing
 - transmission bit errors—can use Unequal Protection Methods
 - tandem encodings

MOS for Range of Speech Coders

Lecture Summary

- the ear acts as a sound canal, transducer, spectrum analyzer
- the cochlea acts like a multi-channel, logarithmically spaced, constant Q filter bank
- frequency and place along the basilar membrane are represented by inner hair cell transduction to events (ensemble intervals) that are processed by the brain
 - this makes sound highly robust to noise and echo
- hearing has an enormous range from threshold of audibility to threshold of pain
 - perceptual attributes scale differently from physical attributes—e.g., loudness, pitch
- masking enables tones or noise to hide tones or noise => this is the basis for perceptual coding (MP3)
- perception and intelligibility are tough concepts to quantify—but they are key to understanding performance of speech processing systems