



AUTOMATIC SPEECH RECOGNITION

Lecture 4

Speech Production



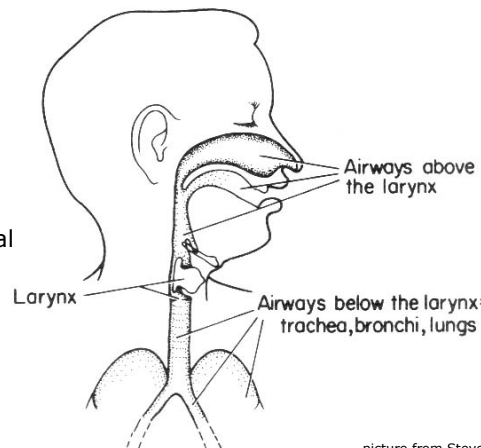
This Lecture

- Physiology of the human speech production mechanism
 - How we generate speech
- A source-filter model of human speech production
 - A link between what happens inside your mouth and the speech that comes out of it



Speech Production System

- 3 parts
 - below the larynx
 - the subglottal system
 - the larynx
 - the vocal folds
 - above the larynx
 - the supraglottal system

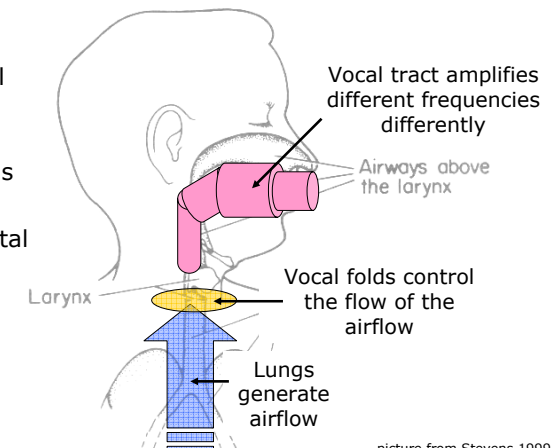


picture from Stevens 1999



Speech Production System

- 3 parts
 - below the larynx
 - the subglottal system
 - the larynx
 - the vocal folds
 - above the larynx
 - the supraglottal system



picture from Stevens 1999

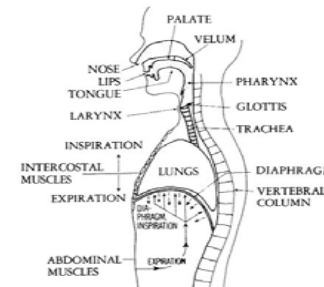


Anatomy of the Subglottal System

- below the larynx
- the trachea
 - 2.5 cm² cross-sectional area
 - 10 to 12 cm length for an adult speaker
- Two bronchi
 - each with one-half the cross-sectional area of the trachea
 - turn into series of smaller airways terminated in the lungs
- Provide airflow for the production of speech



Generating the airflow



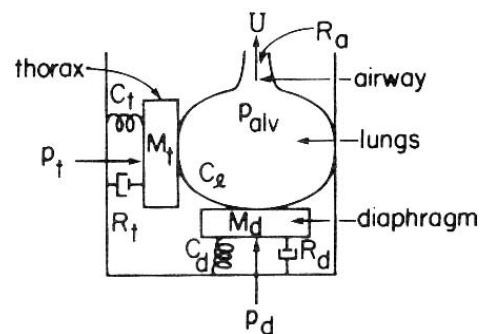
schematic representation of the methods for controlling respiration

picture from Stevens 1999

- mechanical movements of the respiratory system are controlled by:
 - diaphragm
 - abdominal muscles
 - chest wall
- the movement changes the lung's volume
- reflected in the change in the pressure of the air in the lung



Generating the airflow



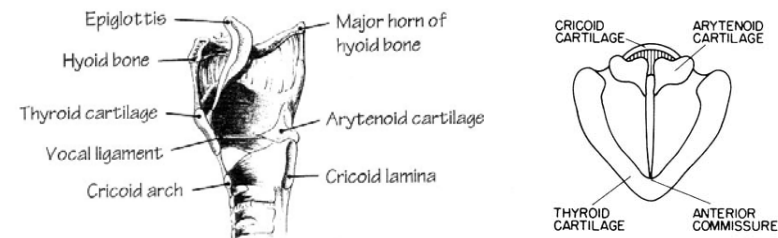
picture from Stevens 1999

Mechanoaerodynamic system for controlling respiratory pressure and flow



The Larynx

- principle structure for speech production are the "vocal folds"
 - two bands of tissue of length 1.0-1.5 cm
 - thickness of 2-3 mm
- The two vocal folds are arranged roughly parallel to each other in an anteroposterior direction.



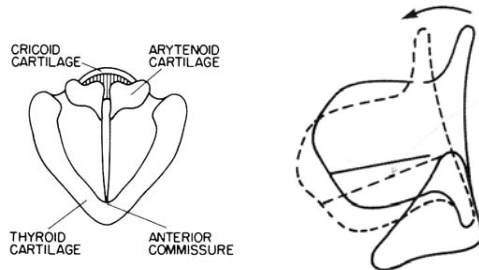
pictures from Stevens 1999



Vocal Folds

- position and stiffness controlled by some surrounding muscles
- "slack" vocal folds → easy to vibrate
- "strict" vocal folds → hard to vibrate

"Feel" Your Vocal Folds
Try 'sssss...' VS. 'zzzzz...'



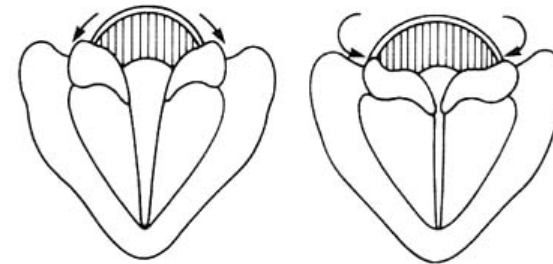
picture from Stevens 1999



Vocal Folds

- abduction and adduction

"Feel" Your Vocal Folds
Try 'haaaaa'

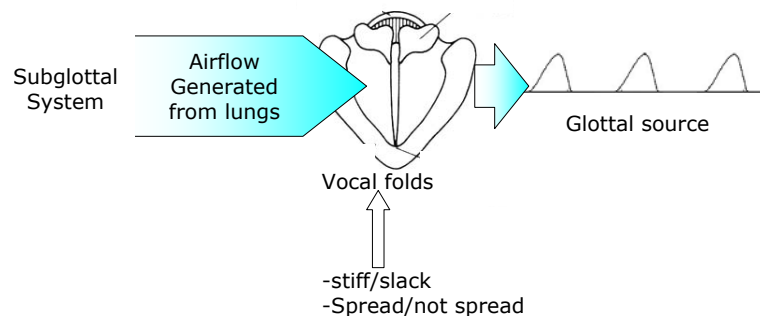


picture from Stevens 1999

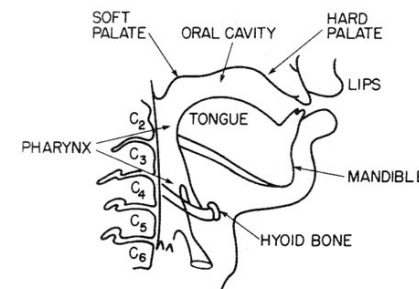


The glottal source

- pressure from the lungs + position/stiffness of the vocal folds



The Supraglottal system



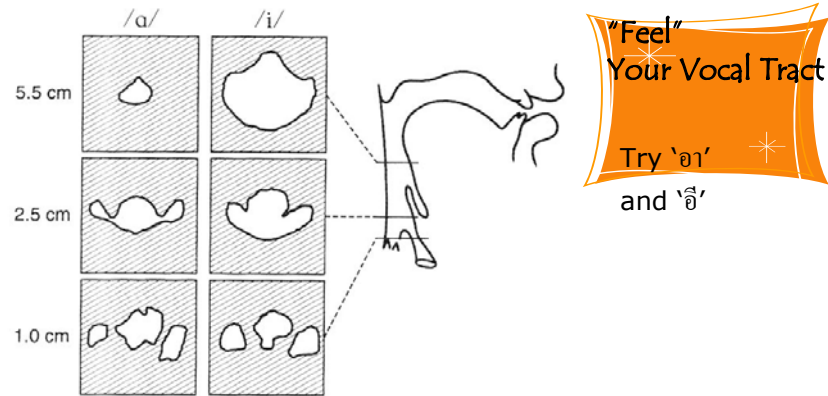
midsagittal section of the vocal tract and surrounding structures

picture from Stevens 1999

- Three pharyngeal constrictor muscles extending from the laryngeal region to the soft palate
- contraction of these muscles produces wide range of shapes and cross-sectional area of the pharyngeal airway
- knowledge of the **cross-sectional area is important** for determining the aerodynamic and acoustic behavior of the vocal tract



Pharynx



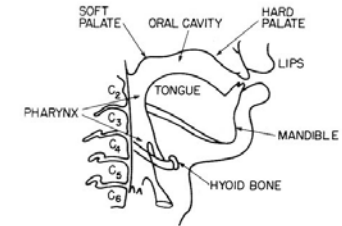
picture from Stevens 1999

cross-sectional shape of the airway at three different positions in the pharyngeal region

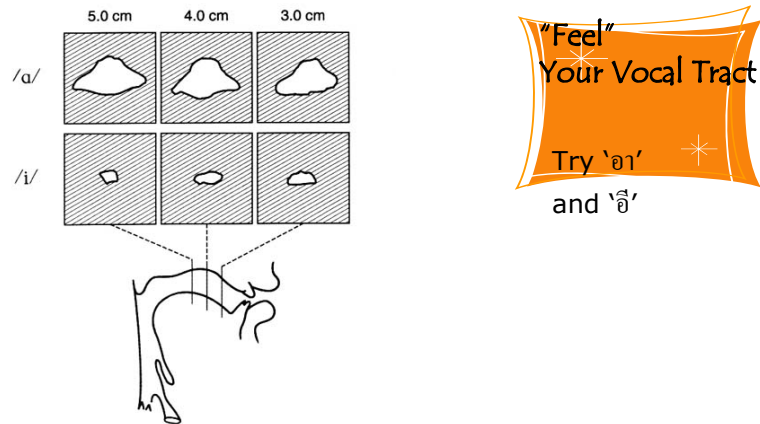


Oral Cavity

- displacement of the tongue body forward and backward changes vocal tract area in the pharyngeal region
- raising and lowering tongue body change vocal tract area in the oral cavity



Oral Cavity



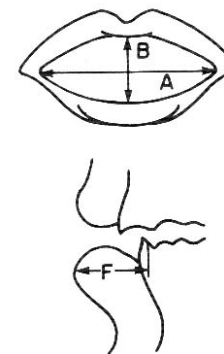
picture from Stevens 1999

cross-sectional shape of the airway



Lips

- The spread and rounding of the lips affect the properties of sound radiated from the mouth opening.



dimensions of the lips

picture from Stevens 1999



Vocal Tract Length

- glottis to the lips opening
- depends on the position of the larynx and the configuration of the lips

	Adult female	Adult male
Vocal Tract Length	14.1 cm	16.9 cm
Pharynx Length	6.3 cm	8.9 cm
Oral Cavity Length	7.8 cm	8.1 cm

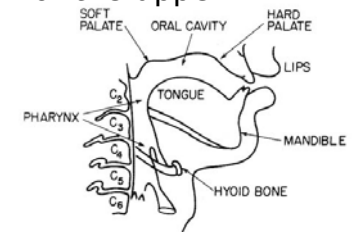
- pharynx length:cavity length ratio is much less in children

Do you think you can make your vocal tract longer or shorter?

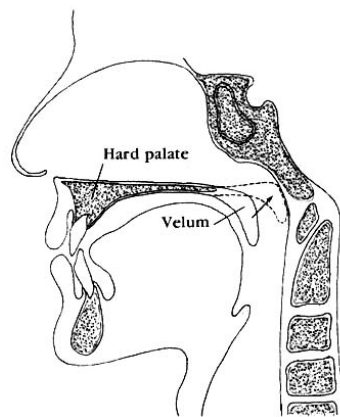


Soft Palate and Nasal Cavity

- upper end of the pharynx extends vertically into the nasal cavity
- the opening into the nasal cavity called the "velopharyngeal opening"
- controlled by the position of the soft palate (velum) and the lateral wall of the upper pharynx



Velopharyngeal Opening



lateral view

- raising and lowering of the soft palate
- cross-sectional area of as large as 1.0 cm²
- 0.2-0.8 cm² when producing sounds that require the participation of the nasal cavity

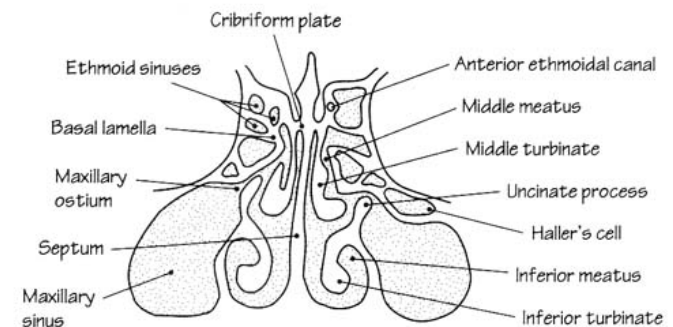
"Feel" Your Vocal Tract
 Try 'ŋ'

vs. 'ŋ'

picture from Stevens 1999



Nasal Cavity



coronal view of the nasal cavity

picture from Stevens 1999

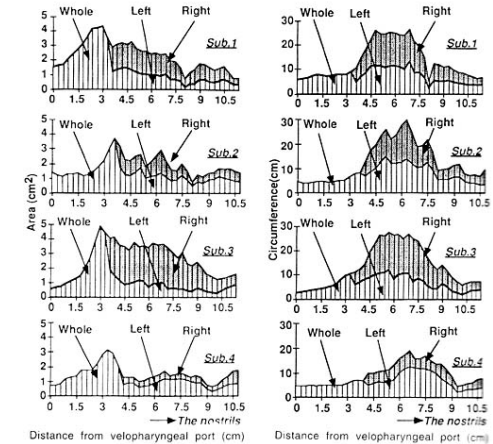


Nasal Cavity

- Over much of its length, the nasal cavity is divided into two passages. (often asymmetrical)
- total length of 11 cm
- total volume of 23 cm³
- narrowest portion at the nostrils (1-2cm²)
- the circumference is 3-5 times greater than that of a circular shape.
- large surface area caused acoustic losses in the nasal cavity. (characteristic of nasal sounds)



Nasal Cavity



picture from Stevens 1999

measurements of the nasal passages for four subjects

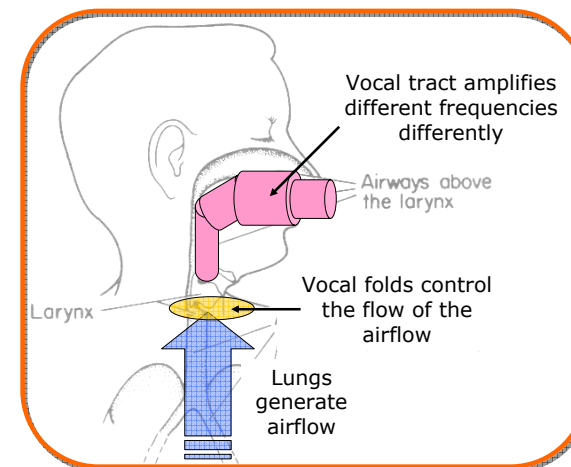


Nasal Cavity

- Cannot be controlled
- Very speaker dependent
 - good for speaker identification
 - bad for speech recognition



The next step...

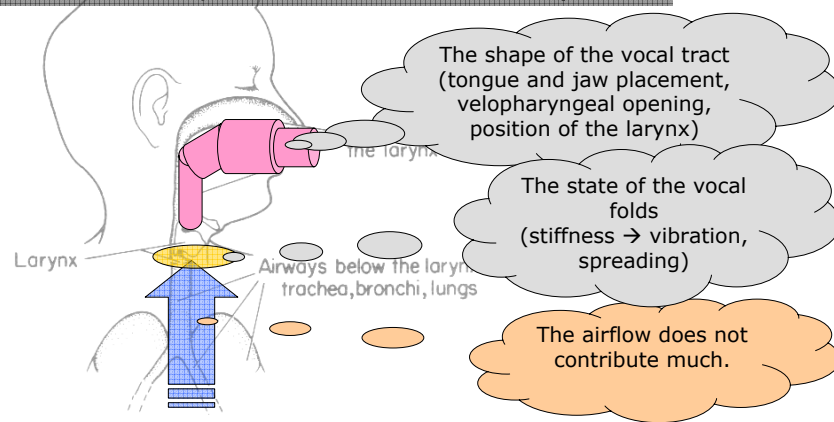


We want to build a model of this system.



Related Factors

How different speech units sound differently to us:



Simple experiment # 1

- Try to feel:
 1. the position of your tongue
 2. the vibration of the vocal folds
 3. the position of the larynx
- while saying these words: อ่า อี้ อู แอ



Simple experiment # 1

- Can you feel the difference in:
 - the position of your tongue?
 - Yes No
 - the vibration of the vocal folds?
 - Yes No
 - the position of the larynx?
 - Yes No



Simple experiment # 2

- Try to feel:
 1. the position of your tongue
 2. the vibration of the vocal folds
 3. the position of the larynx
- while saying these words:

มอ หม่อ หมื่อ ม้อ หมอ



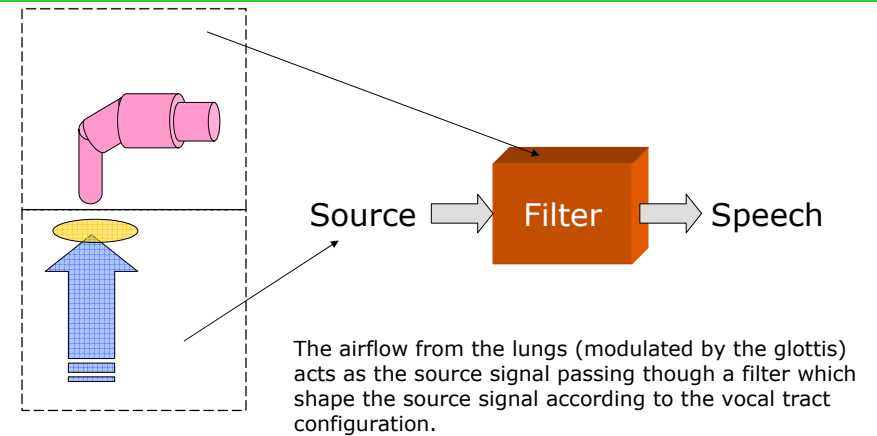


Simple experiment # 2

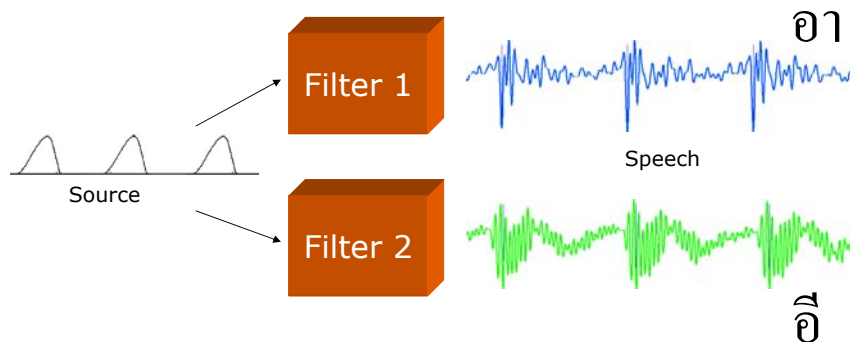
- Can you feel the difference in:
 - the position of your tongue?
 - Yes No
 - the vibration of the vocal folds?
 - Yes No
 - the position of the larynx?
 - Yes No



Source-filter model



Example



Filter 1 is associated with the vocal tract configuration of $oɪ$

Filter 2 is associated with the vocal tract configuration of $əɪ$



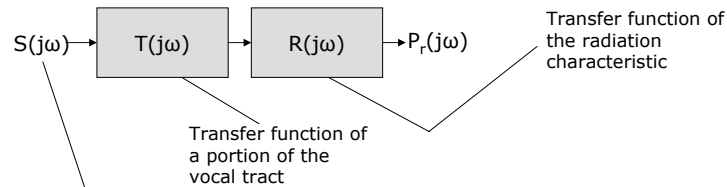
Key Points for Modeling Speech Production

- Sources:
 - Different types of sources used in producing sounds in languages.
- Filters:
 - Amplify different sound frequencies.
 - Link those frequencies with the vocal tract configuration.
 - Other factors.



Source-filter Model of Speech Production

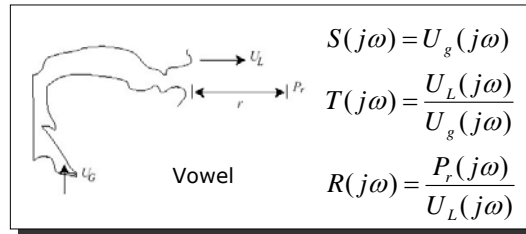
$$P_r(j\omega) = S(j\omega)T(j\omega)R(j\omega)$$



Transfer function of the source signal

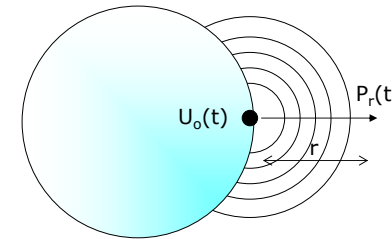
The source can be:

- Glottal Source
- Noise Source



Radiation Characteristic

- The sound pressure at distance r from the lips, $P_r(t)$, is linearly related to the volume velocity $U_o(t)$ at the lips.
- For up to 4000Hz, the mouth opening can be regarded as a simple source of strength $U_o(t)$



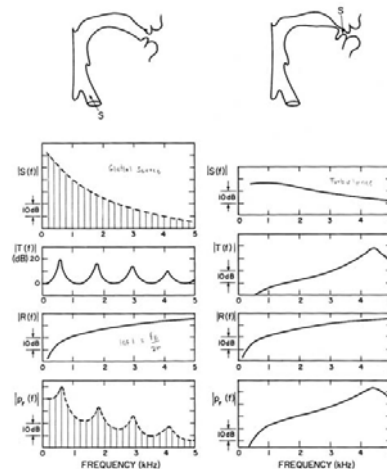
$$R(j\omega) = \frac{P_r(j\omega)}{U_o(j\omega)}$$

$$= \frac{j\omega}{4\pi r} e^{-j\frac{\omega}{c}r}$$

$$|R(f)| = \frac{f\rho}{2r}$$



Examples

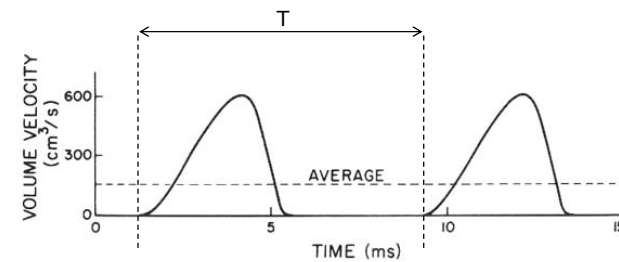


picture from Stevens 1999



Glottal source

- Volume Velocity Source at Glottis U_g
- generated from the vibration of the vocal folds
- periodic

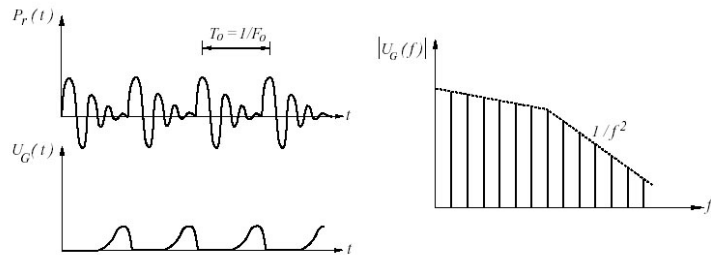


typical shape of the glottal waveform

picture MIT OCW



Glottal source



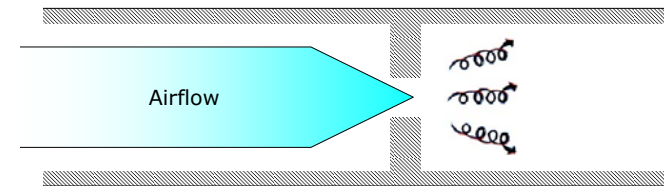
	F ₀ ave (Hz)	F ₀ min (Hz)	F ₀ max (Hz)
Men	125	80	200
Women	225	150	350
Children	300	200	500

picture from MIT OCW

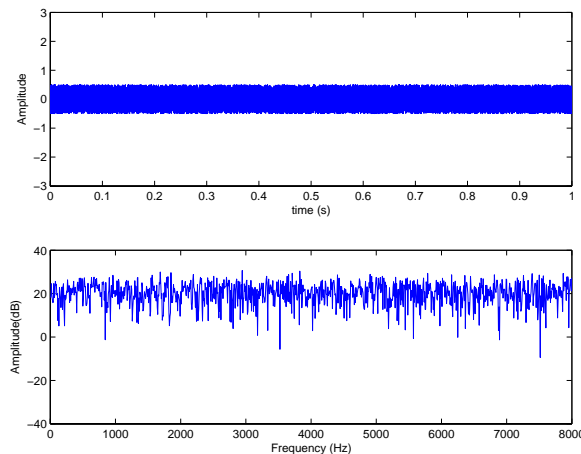


Noise Source

- Noise source
 - generated from air turbulence at some constrictions in the vocal tract
 - white noise



Noise Source

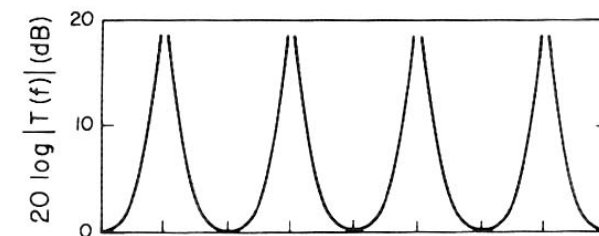


Waveform and spectrum of the white noise



Vocal tract transfer function

- $T(j\omega) = U_l(j\omega)/U_s(j\omega)$
- usually characterized by several peaks corresponding to resonances of the cavity of the vocal tract



picture from Stevens 1999



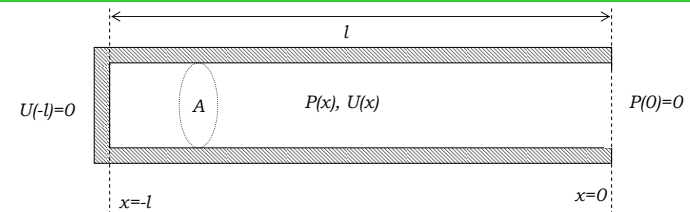
A Little Summary

- So, now we know that there are two kinds of source.
- Whichever the source, it will pass through a filter which is a frequency-dependent amplifier.
- The filter will amplify frequency components of the source. This results in the speech waveform that comes out of the mouth (and reach a microphone, if there is one.)

NOW: We would like to know that, given a vocal tract configuration, how the frequency-dependent amplifier is going to be like.



Natural Frequencies of a Uniform Tube



Acoustic equations

$$\frac{\partial P(x,t)}{\partial x} = -\frac{\rho}{A} \frac{\partial U(x,t)}{\partial t} \quad \text{①}$$

$$\frac{\partial U(x,t)}{\partial x} = -\frac{A}{\gamma P_0} \frac{\partial P(x,t)}{\partial t} \quad \text{②}$$

ρ : density of air (0.00114 gm/cm³)
 γ : ratio of specific heat at constant pressure to specific heat at constant volume (1.4 for air)
 P_0 : ambience pressure level (cm H₂O)
 c : velocity of sound (35,400 cm/s)



Natural Frequencies of a Uniform Tube

Assume:

$$\left. \begin{aligned} P(x,t) &= P(x)e^{j\omega t} \\ U(x,t) &= U(x)e^{j\omega t} \end{aligned} \right\} \text{Air particles at position } x \text{ oscillates with angle frequency } \omega$$

From ① and ② \Rightarrow $\frac{\partial^2 P(x)}{\partial x^2} + k^2 P(x) = 0 \quad \text{③} \quad k = \frac{\omega}{c}, c = \sqrt{\frac{\gamma P_0}{\rho}}$

- The solution of ③ is either in the form of *sin* or *cos*
- The requirements for natural frequencies are $p(0)=0$ and $U(-l)=0$

$\Rightarrow P(0)=0 \Rightarrow P(x) = P_m \sin(kx)$



Natural Frequencies of a Uniform Tube

- From ① and $P(x,t) = P(x)e^{j\omega t} \Rightarrow U(x) \propto \frac{\partial P(x)}{\partial x}$

$$\frac{\partial P(x)}{\partial x} = kP_m \cos(kx)$$

$$\therefore U(-l) = C \cos(-kl) = C \cos(kl) = 0$$

$$\cos(kl) = 0 \Leftrightarrow kl = \frac{\pi}{2}, \frac{3\pi}{2}, \frac{5\pi}{2}, \dots$$

So, formant frequencies are

$$f = \frac{c}{4l}, \frac{3c}{4l}, \frac{5c}{4l}, \dots$$



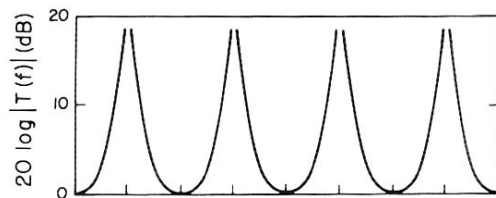
Transfer Function

$$U(-l, t) = C \cos(kl)e^{j\omega t}$$

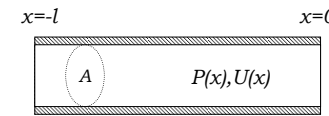
$$U_{-l}(j\omega) = C \cos(kl)F(e^{j\omega t})$$

$$U_0(j\omega) = C \cos(0)F(e^{j\omega t})$$

$$T(j\omega) = \frac{U_0(j\omega)}{U_{-l}(j\omega)} = \frac{1}{\cos(kl)}$$

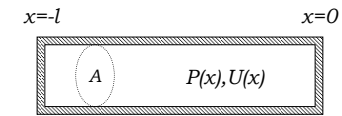


Uniform Tube open/closed at both ends



$$P(0) = 0, U(0) = U_m$$

$$P(-l) = 0, U(-l) = U_m$$



$$P(0) = P_m, U(0) = 0$$

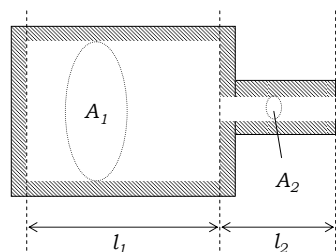
$$P(-l) = P_m, U(-l) = 0$$

Formant frequencies are

$$f = 0, \frac{c}{2l}, \frac{c}{l}, \frac{3c}{2l}, \dots$$



Helmholtz Resonator



$$f_1 = \frac{1}{2\pi\sqrt{M_A C_A}}$$

$$M_A = \frac{\rho l_2}{A_2}, C_A = \frac{A_1 l_1}{\rho c^2}$$

$$f_1 = \frac{c}{2\pi} \sqrt{\frac{A_2}{A_1 l_1 l_2}}$$

M_A : Acoustic Mass
 C_A : Acoustic Compliance

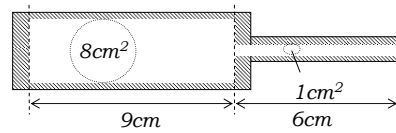


Vocal Tract Modeling by Concatenated Tubes

- A number of vocal tract shapes, both vowels and consonants, can be approximated by two or more resonators or uniform tubes of different cross-sectional areas connected together.
- neglect coupling between components
- combination of natural frequencies for all of the components (+Helmholtz structure)



Concatenated Tubes



$$f = \frac{354}{2\pi} \sqrt{\frac{1 \times 10^{-4}}{(8 \times 10^{-4})(9 \times 10^{-2})(6 \times 10^{-2})}} \quad f_1 = \frac{c}{2l} = \frac{354}{2 \times 9 \times 10^{-2}} = 1967 \text{ Hz} \quad f_1 = \frac{c}{2l} = \frac{354}{2 \times 6 \times 10^{-2}} = 2950 \text{ Hz}$$

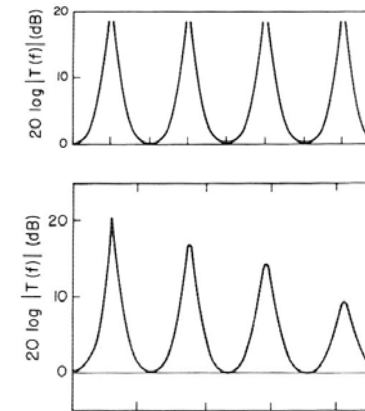
$$f = 271 \text{ Hz} \quad f_2 = \frac{c}{l} = 1967 \times 2 = 3933 \text{ Hz} \quad f_2 = \frac{c}{l} = 2950 \times 2 = 5180 \text{ Hz}$$

$$f_3 = \dots \quad f_3 = \dots$$

formant frequencies = 271Hz, 1967Hz, 2950Hz



Loss in the Vocal Tract

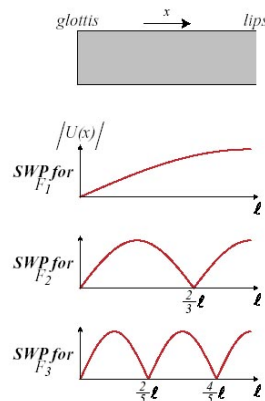


picture from Stevens 1999



Standing Wave Pattern

A uniform tube open at one end and closed at the other is often referred to as a "quarter wavelength resonator"



picture from MIT OCW

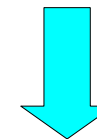


Perturbation Theory

- Narrowing the tube where $U(x)$ is maximum in the standing wave pattern for a given formant
- Narrowing the tube where $P(x)$ is maximum in the standing wave pattern for a given formant



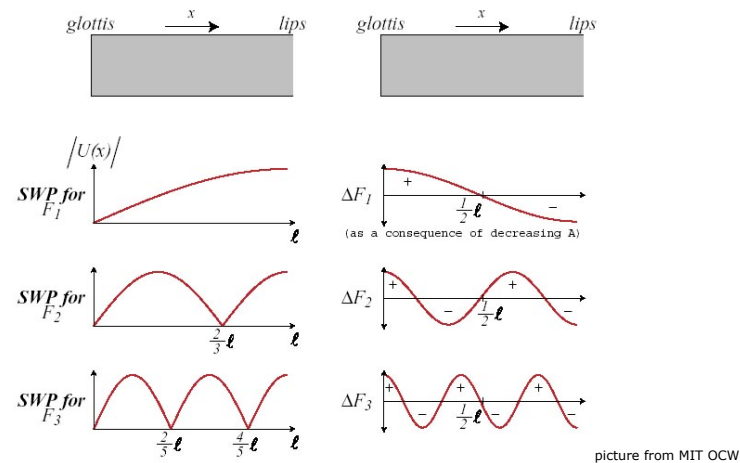
decrease the value of that formant



increase the value of that formant



Summary of Perturbation Theory



References

- Stevens, Kenneth. *Acoustic Phonetics*. Cambridge, MA: MIT Press, 1999. ISBN: 0-262-19404-X.