

Speech Production and Perception

TTT4185 Machine Learning for Signal Processing

Giampiero Salvi

Department of Electronic Systems
NTNU

HT2021

Why use speech?

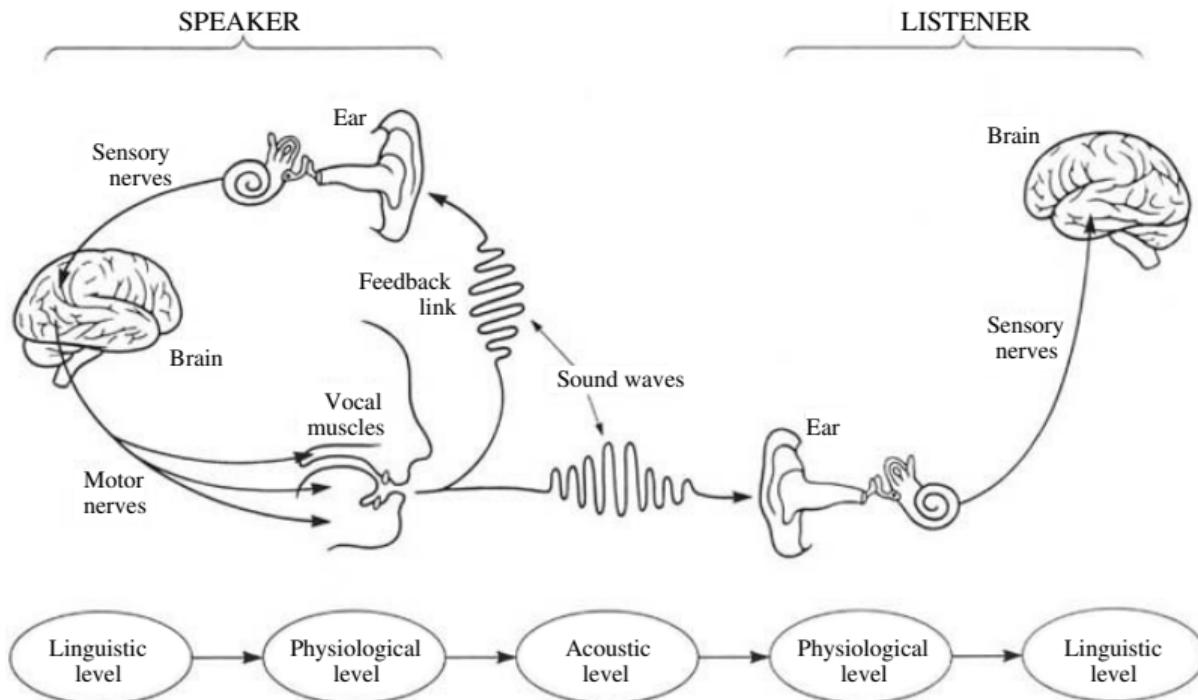
Human-Computer (or -Robot) Interaction

- Natural way of communication (No training needed)
- Leaves hands and eyes free (Good for functionally disabled)
- Effective (Higher data rate than typing)
- Can be transmitted/received inexpensively (phones)

Surveillance/Search

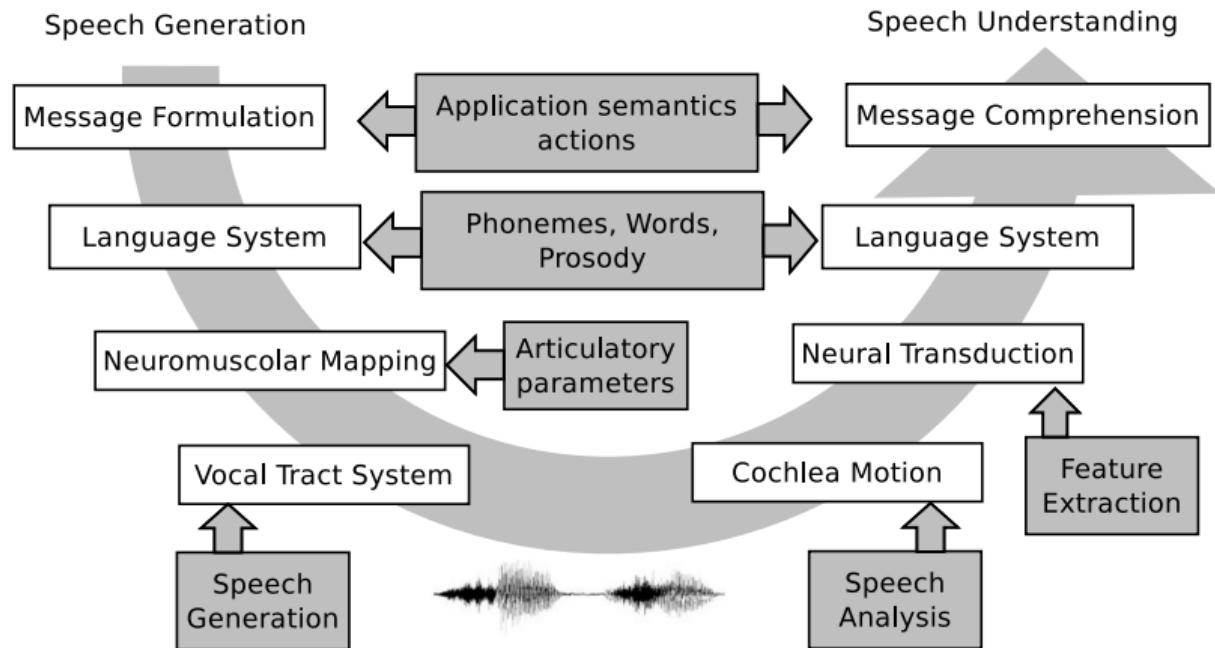
- Transcribe human-human conversations
- produce indexing for broadcast material
- produce subtitles for movies/news

The Speech Chain



Peter Denes, Elliot Pinson, 1963

The Speech Chain revisited



from Huang, Acero and Hon's book

Outline

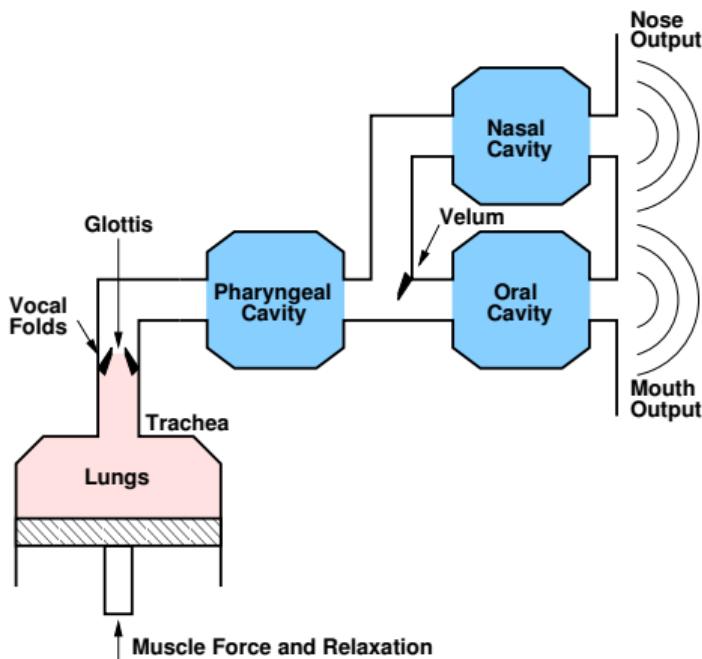
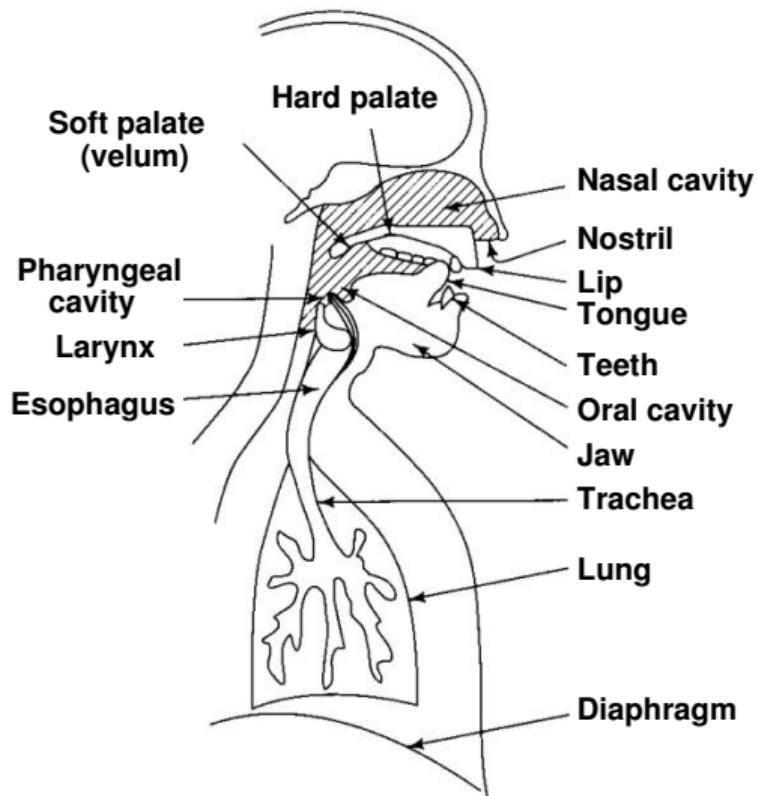
1 Speech Production

- Source/Filter Model

2 Speech Perception

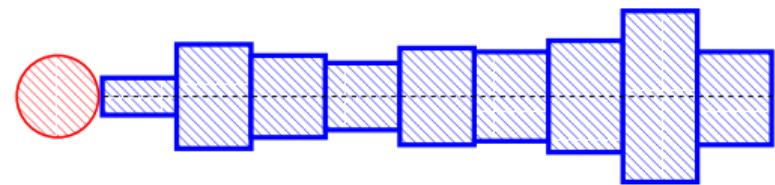
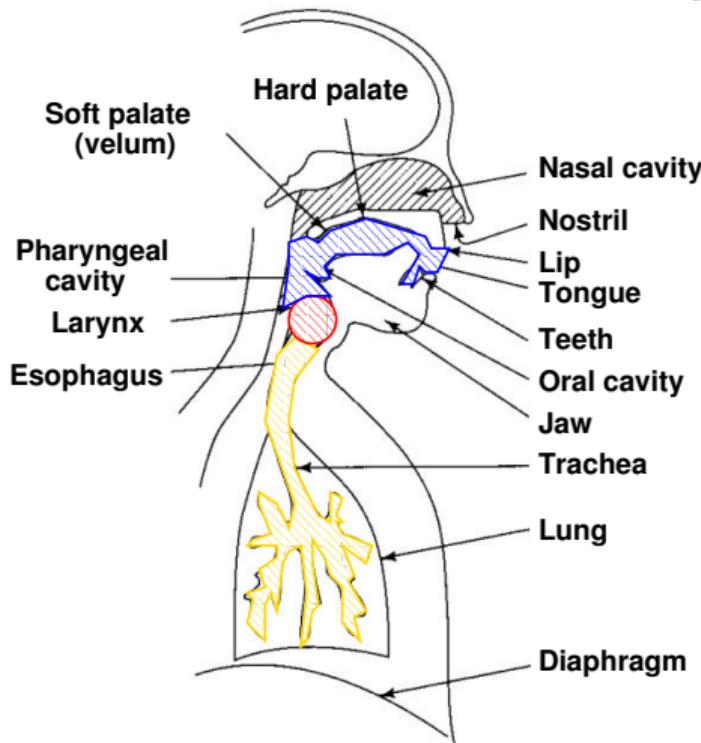
3 Challenges

Physiology



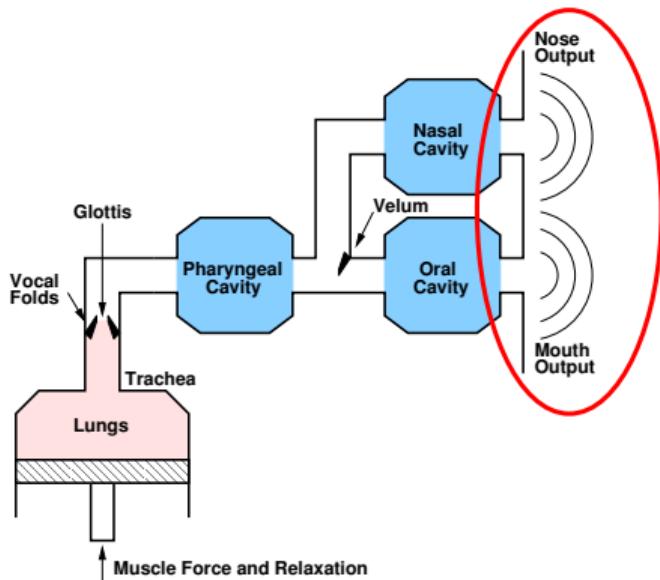
Source/Filter Model, Vowel-like sounds

Vowels



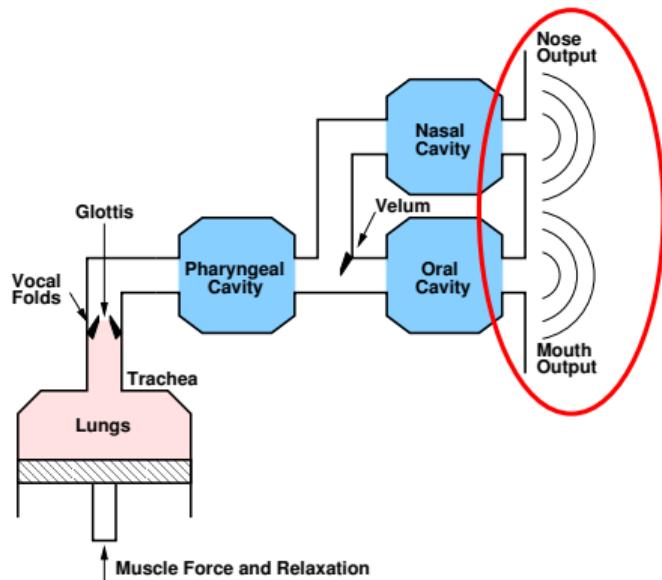
- Source (periodic)
- Front Cavity
- Back Cavity
- Back Cavity (2nd approx.)

Radiation form the Lips/Nose

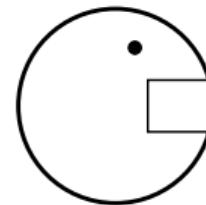


Problem of radiation at the lips plus diffraction about the head too complicated.

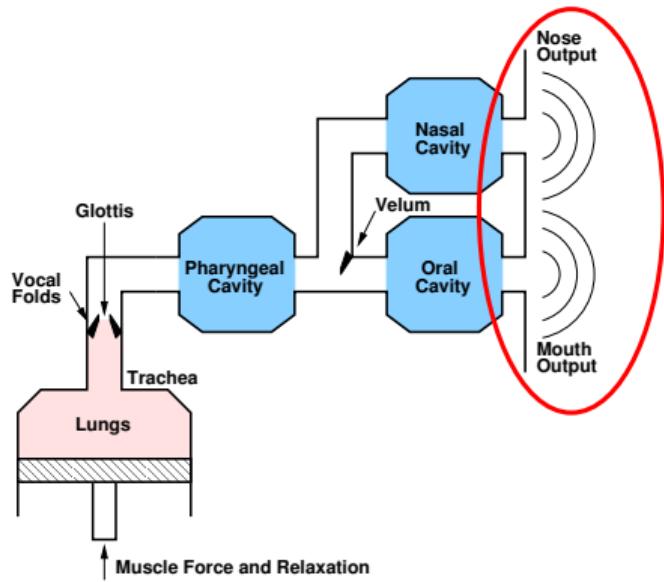
Radiation form the Lips/Nose



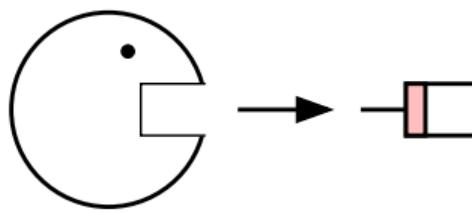
Approx. with a piston in a rigid sphere: solved
but not in closed form



Radiation form the Lips/Nose



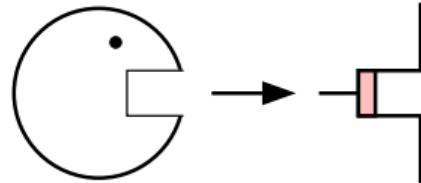
2nd approx: piston in an infinite wall



$$R(z) \approx 1 - \alpha z^{-1}$$

Radiation form the Lips/Nose

2nd approx: piston in an infinite wall



$$R(z) \approx 1 - \alpha z^{-1}$$

Question:

Given $R(z) = 1 - \alpha z^{-1}$ is the transfer function of a linear system, what is the relationship between the system input $x[n]$ and its output $y[n]$?

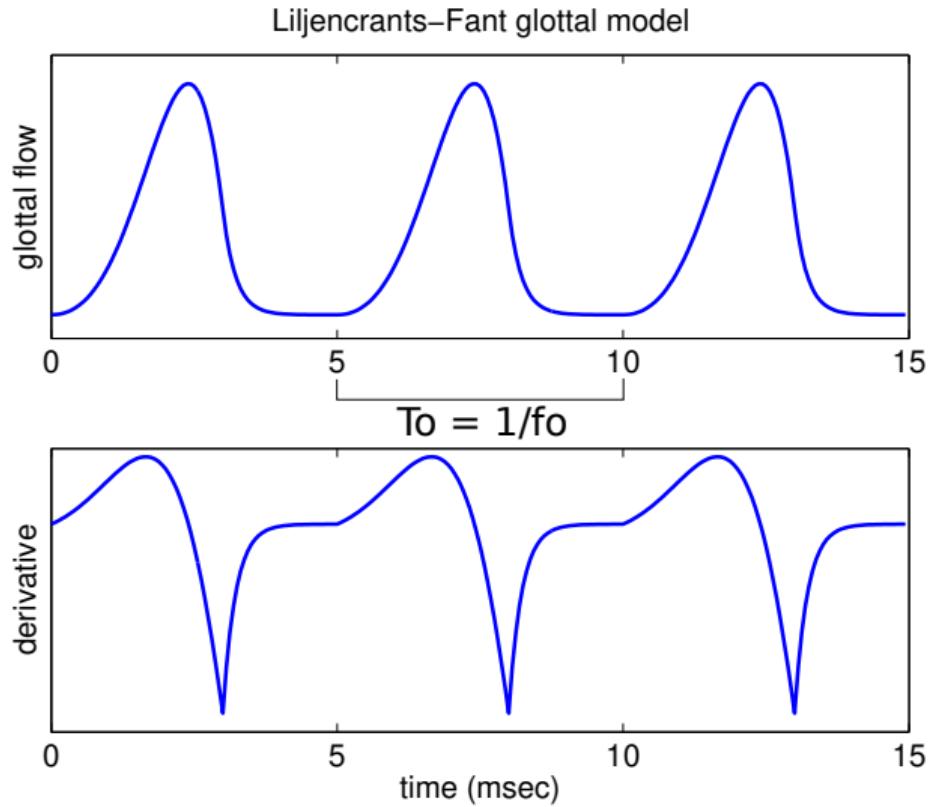
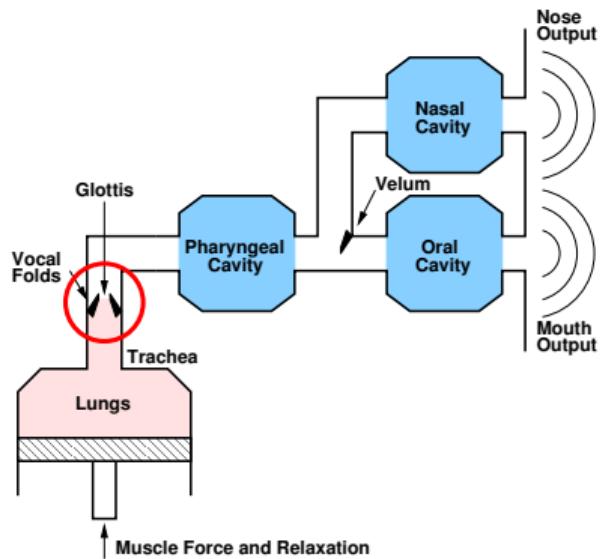
- a) $y[n] = \text{constant}$
- b) $y[n] = 1 - \alpha x[n]$
- c) $y[n] = x[n] - \alpha x[n - 1]$
- d) $y[n] = 1 - \alpha x[n - 1]$

Glottal Flow: Laryngoscopy

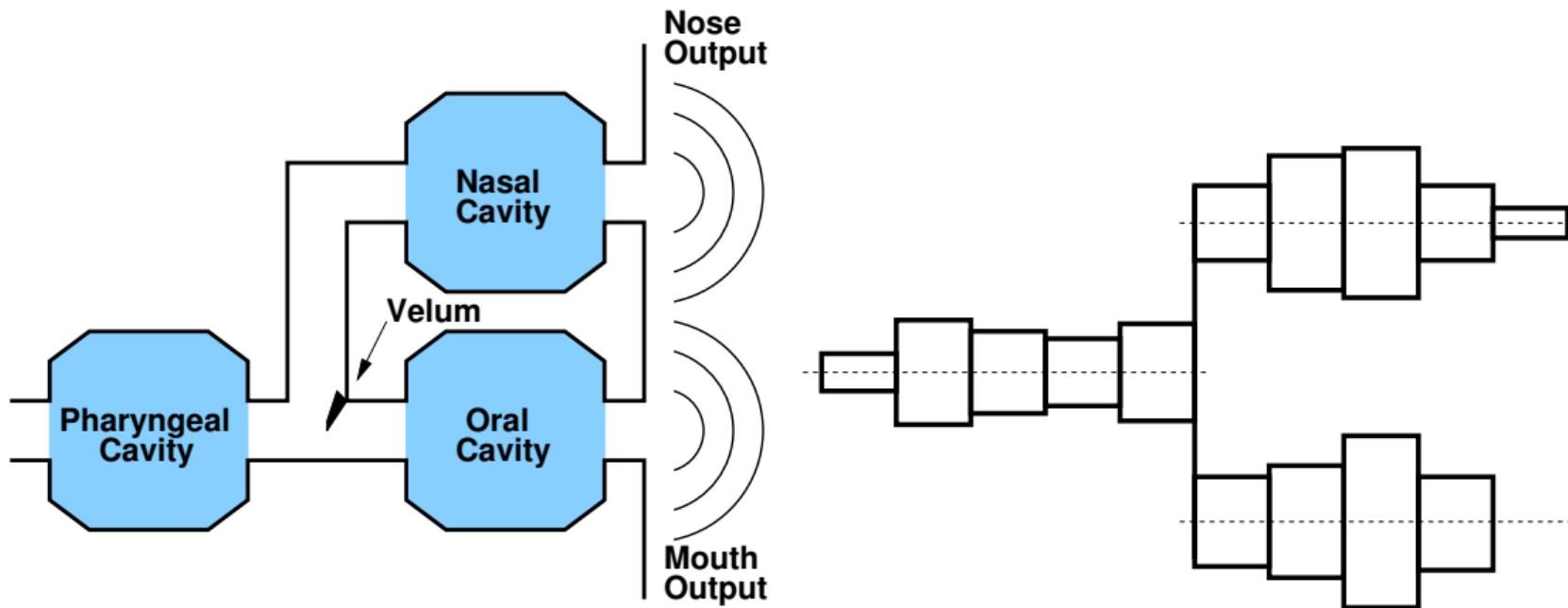


<https://youtu.be/iYpDwhpILkQ>

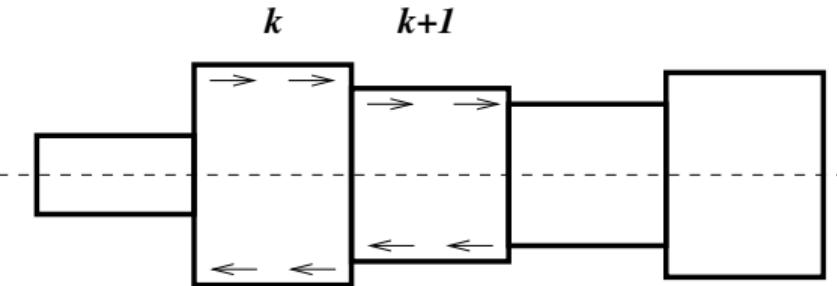
Glottal Flow



Tube Model of the Vocal Tract



Tube Model (cntd.)

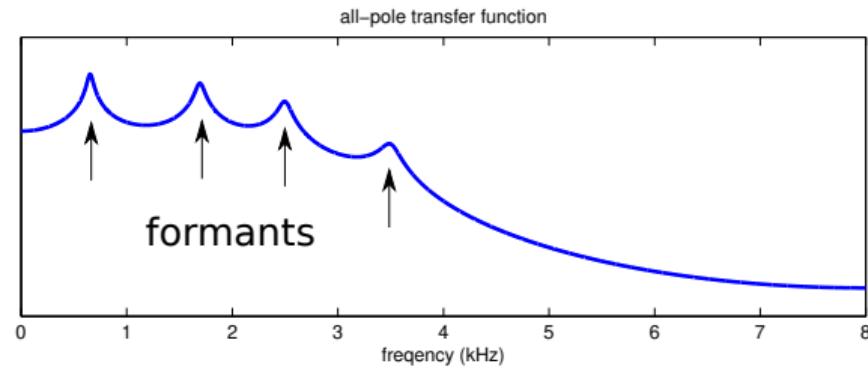
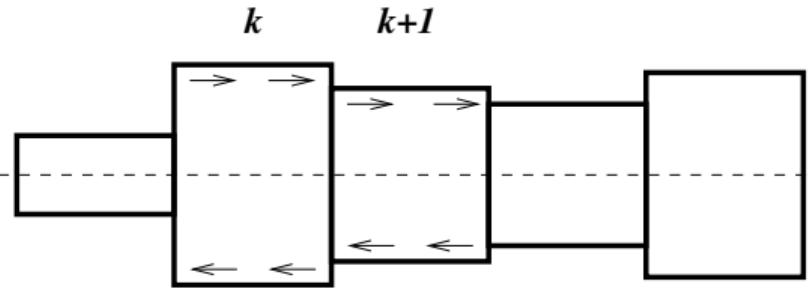


- assume planar wave propagation and lossless tubes
- solve pressure $p(x, t)$ and velocity $u(x, t)$ in each tube according to wave equation
- impose continuity of pressure and velocity at the junctions

⇒ all-pole transfer function ($N = \text{number of tubes}$)

$$V(z) = \frac{Az^{-N/2}}{1 - \sum_{k=1}^N a_k z^{-k}}$$

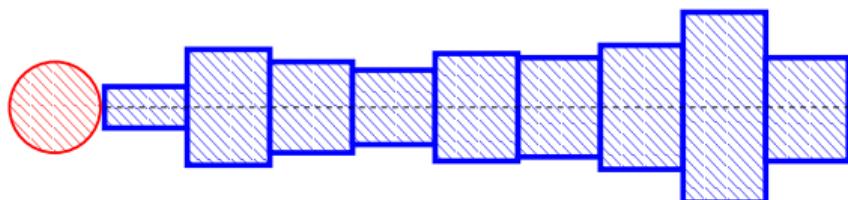
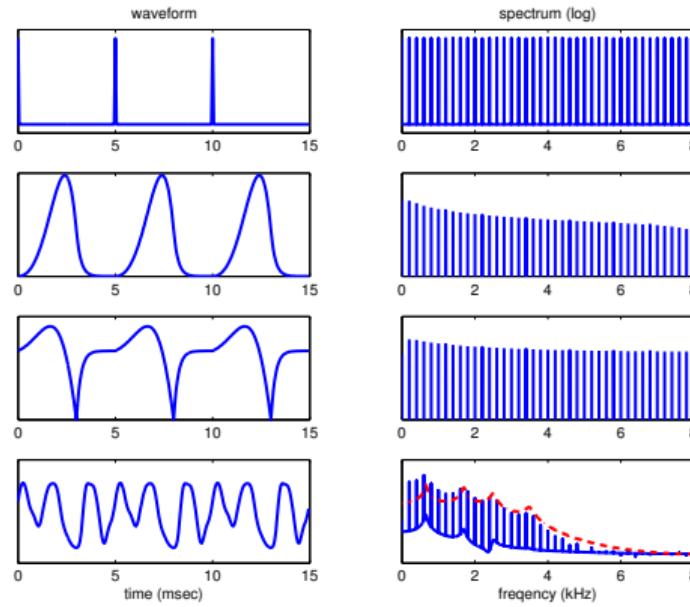
Tube Model (cntd.)



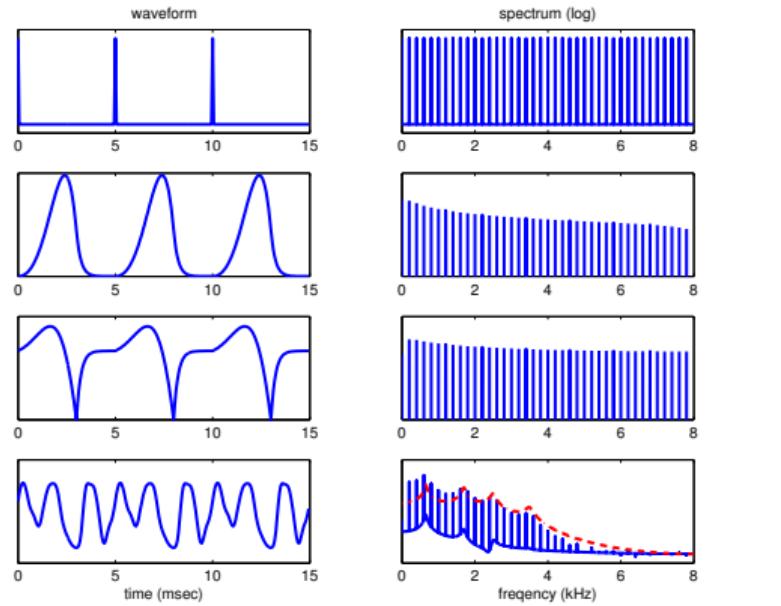
- assume planar wave propagation and lossless tubes
 - solve pressure $p(x, t)$ and velocity $u(x, t)$ in each tube according to wave equation
 - impose continuity of pressure and velocity at the junctions
- ⇒ all-pole transfer function ($N = \text{number of tubes}$)

$$V(z) = \frac{Az^{-N/2}}{1 - \sum_{k=1}^N a_k z^{-k}}$$

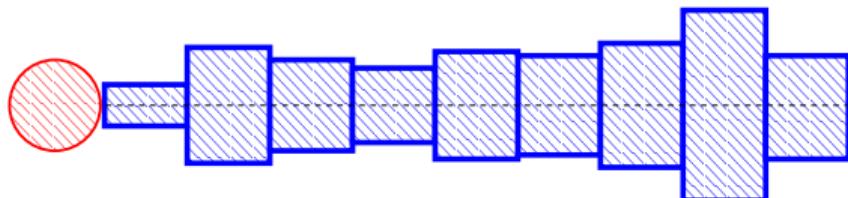
Source/Filter Model: vowel-like sounds



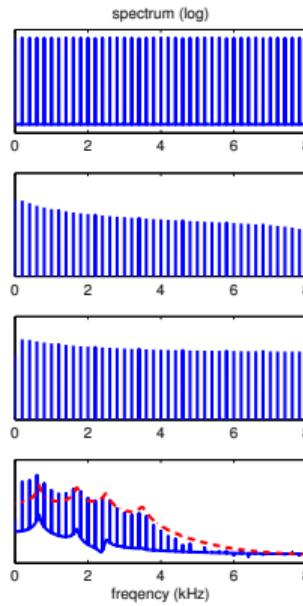
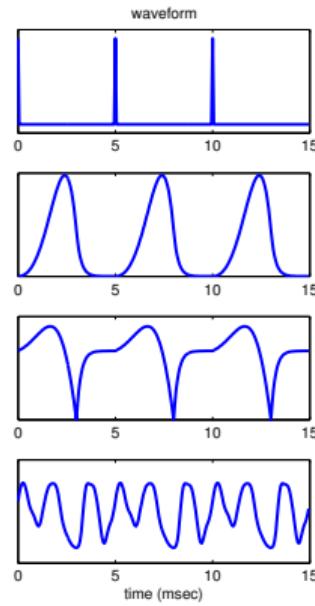
Source/Filter Model: vowel-like sounds



$\leftarrow p[n]$

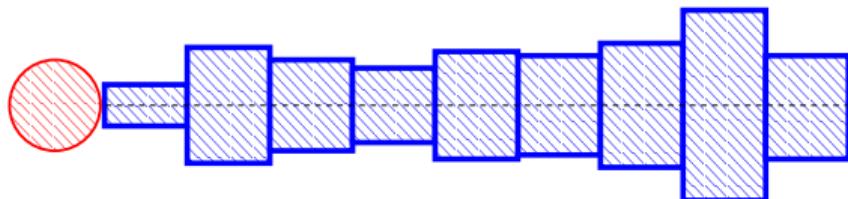


Source/Filter Model: vowel-like sounds

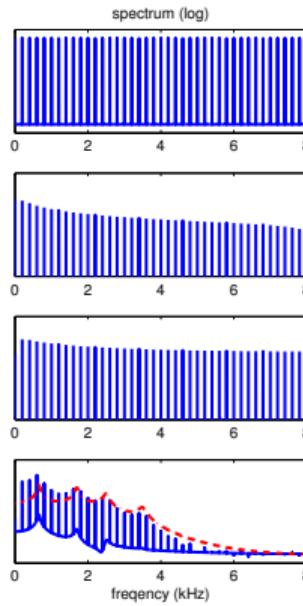
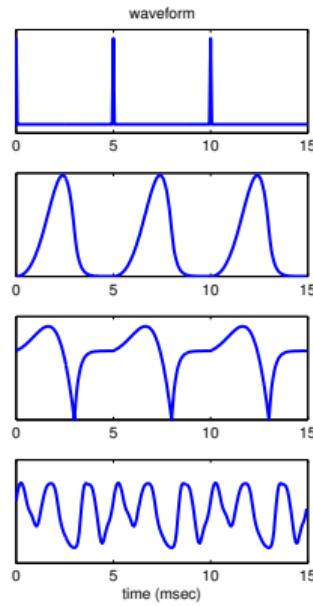


$$\leftarrow p[n]$$

$$\leftarrow p[n] * g[n]$$



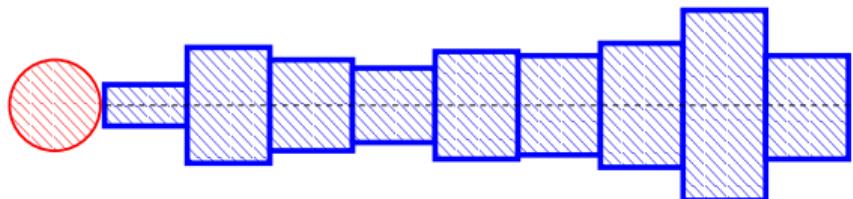
Source/Filter Model: vowel-like sounds



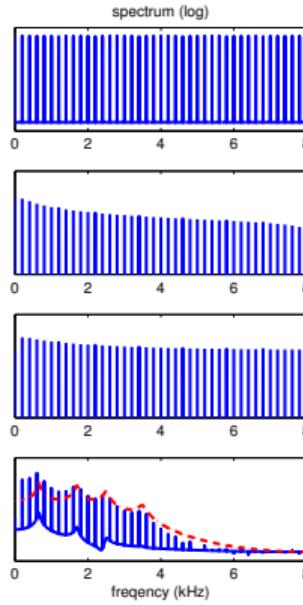
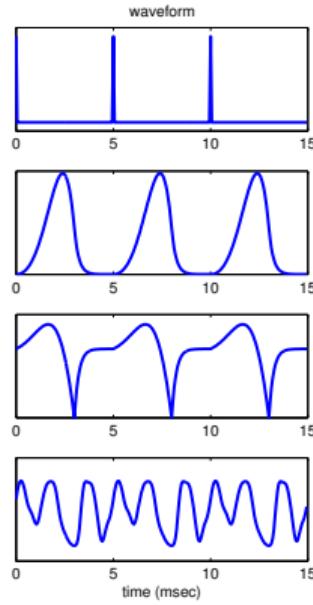
$$\leftarrow p[n]$$

$$\leftarrow p[n] * g[n]$$

$$\leftarrow p[n] * g[n] * r[n]$$



Source/Filter Model: vowel-like sounds

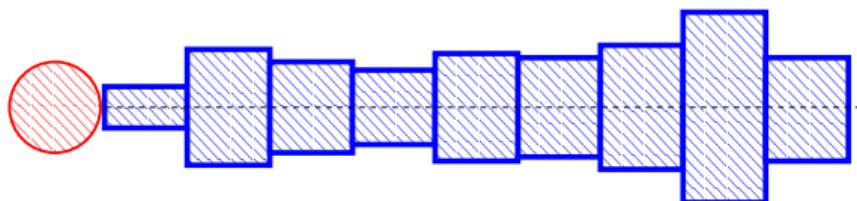


$$\leftarrow p[n]$$

$$\leftarrow p[n] * g[n]$$

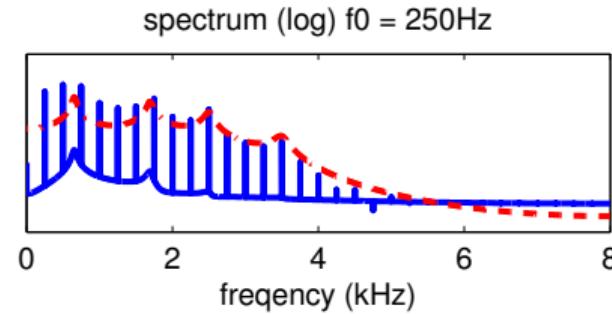
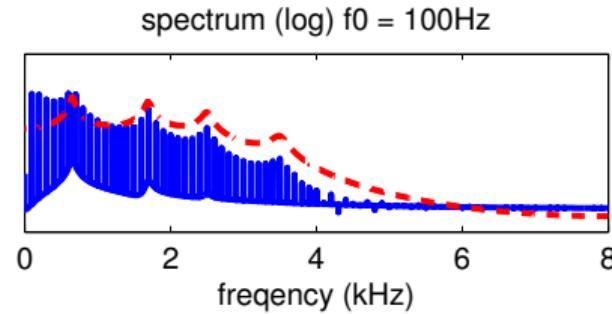
$$\leftarrow p[n] * g[n] * r[n]$$

$$\leftarrow p[n] * g[n] * r[n] * v[n]$$



f_0 and Formants

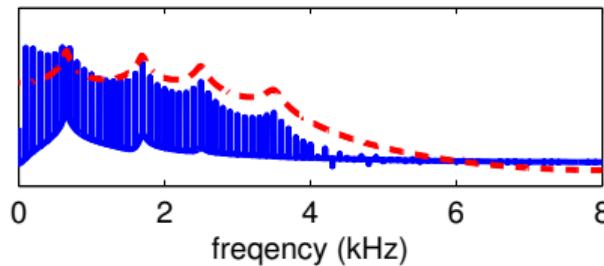
- Varying f_0 (vocal fold oscillation rate)



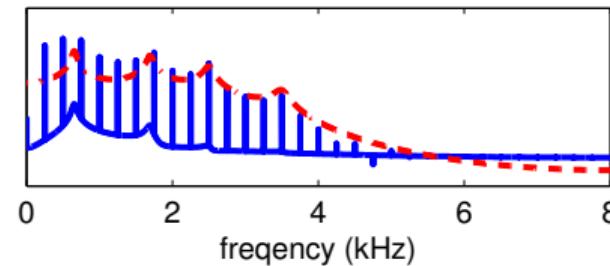
f_0 and Formants

- Varying f_0 (vocal fold oscillation rate)

spectrum (log) $f_0 = 100\text{Hz}$

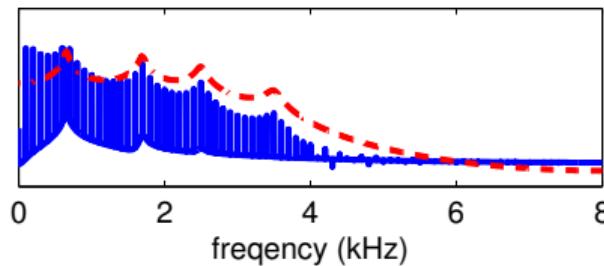


spectrum (log) $f_0 = 250\text{Hz}$

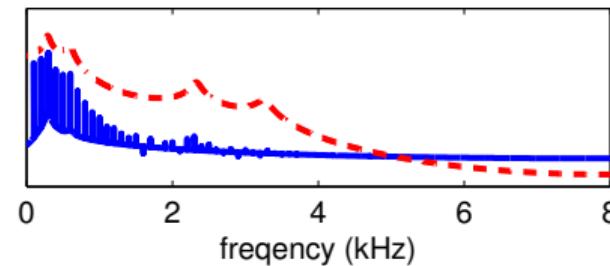


- Varying Formants (vocal tract shape)

spectrum (log) vowel [ɛ]



spectrum (log) vowel [u]



Demonstration of speech sounds

Wavesurfer

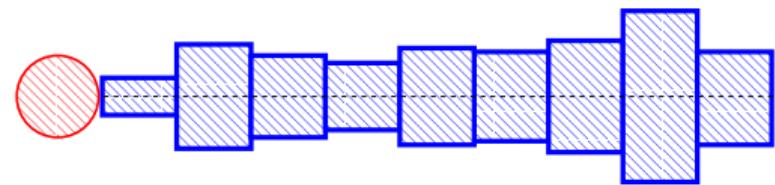
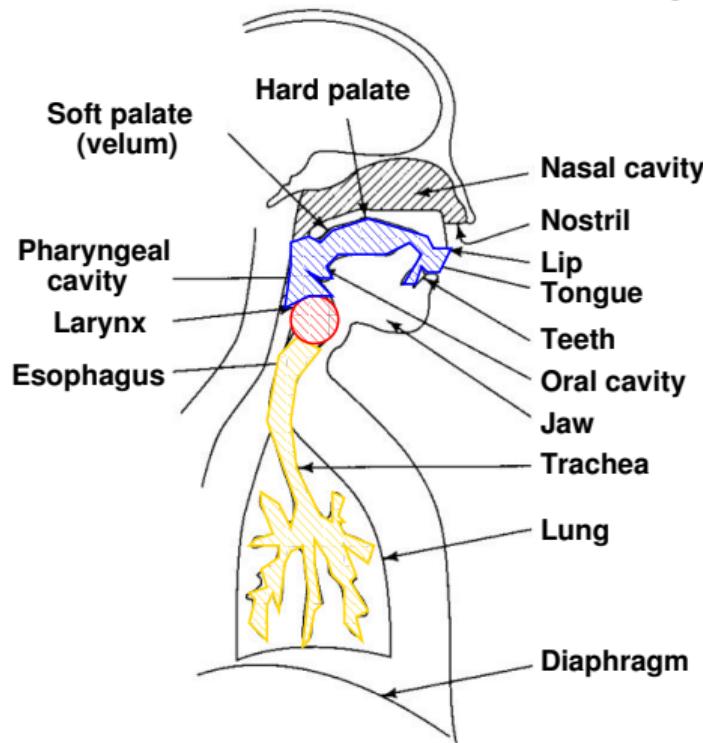
<https://sourceforge.net/projects/wavesurfer/>

Praat

<https://www.fon.hum.uva.nl/praat/>

Source/Filter Model, General Case

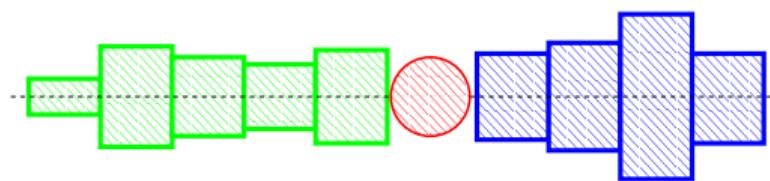
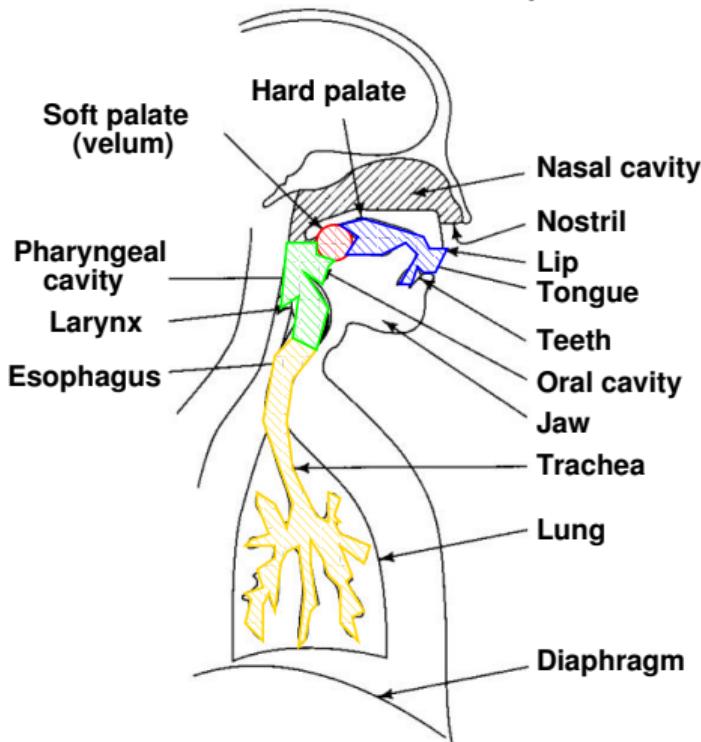
Vowels



- Source (periodic)
- Front Cavity
- Back Cavity
- Back Cavity (2nd approx.)

Source/Filter Model, General Case

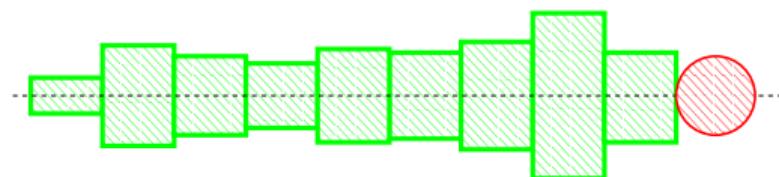
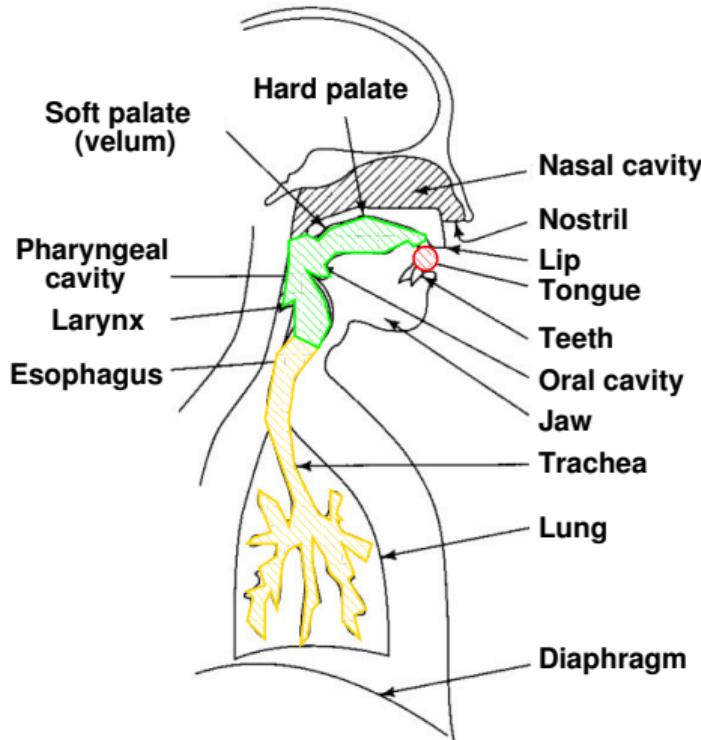
Fricatives (e.g. sh) or Plosive (e.g. k)



- Source (noise or impulsive)
- Front Cavity
- Back Cavity
- Back Cavity (2nd approx.)

Source/Filter Model, General Case

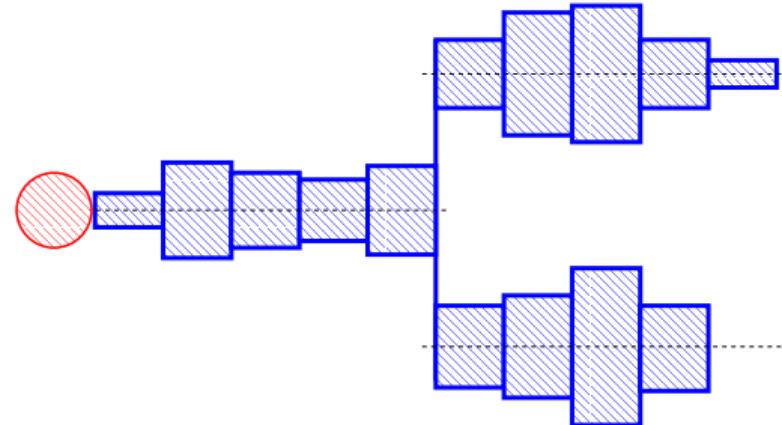
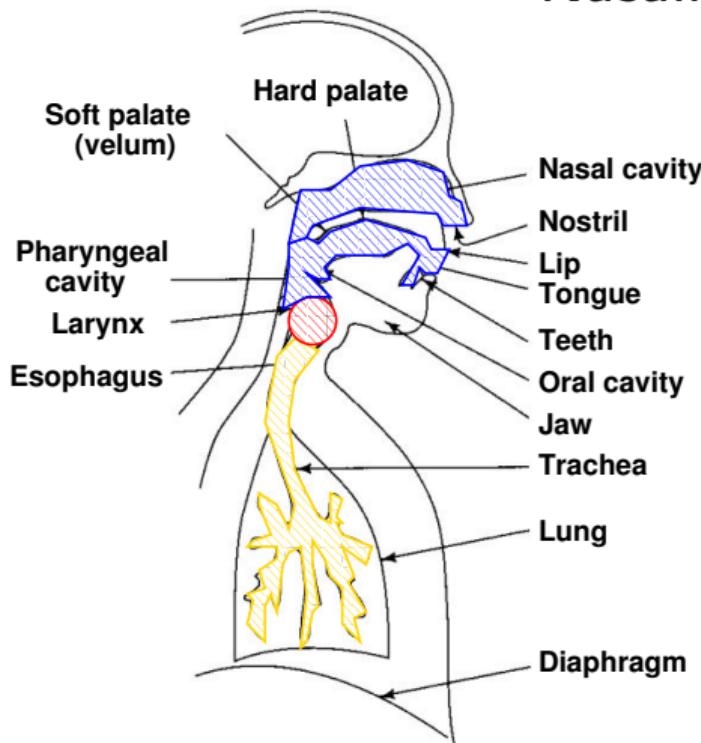
Fricatives (e.g. s) or Plosive (e.g. t)



- Source (noise or impulsive)
- Front Cavity
- Back Cavity
- Back Cavity (2nd approx.)

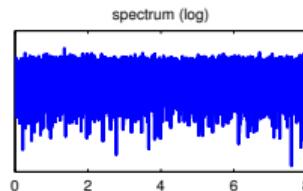
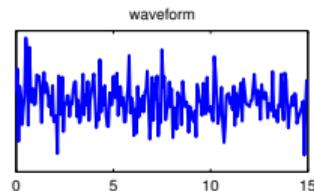
Source/Filter Model, General Case

Nasalised Vowels

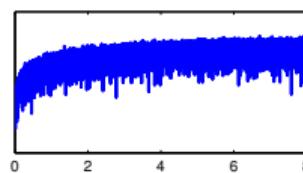
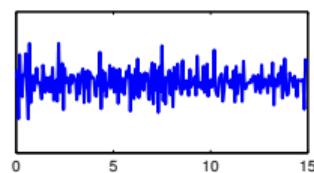


- Source (periodic)
- Front Cavity
- Back Cavity
- Back Cavity (2nd approx.)

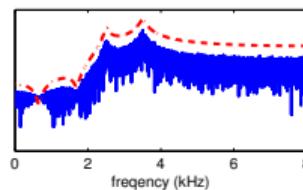
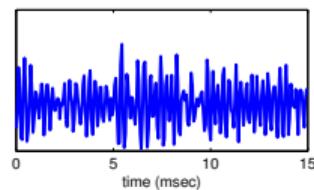
Source/Filter Model: fricative sounds



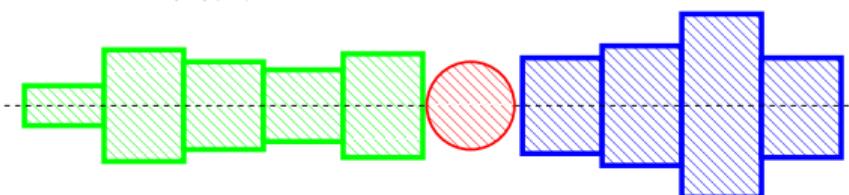
$\leftarrow p[n]$



$\leftarrow p[n] * r[n]$

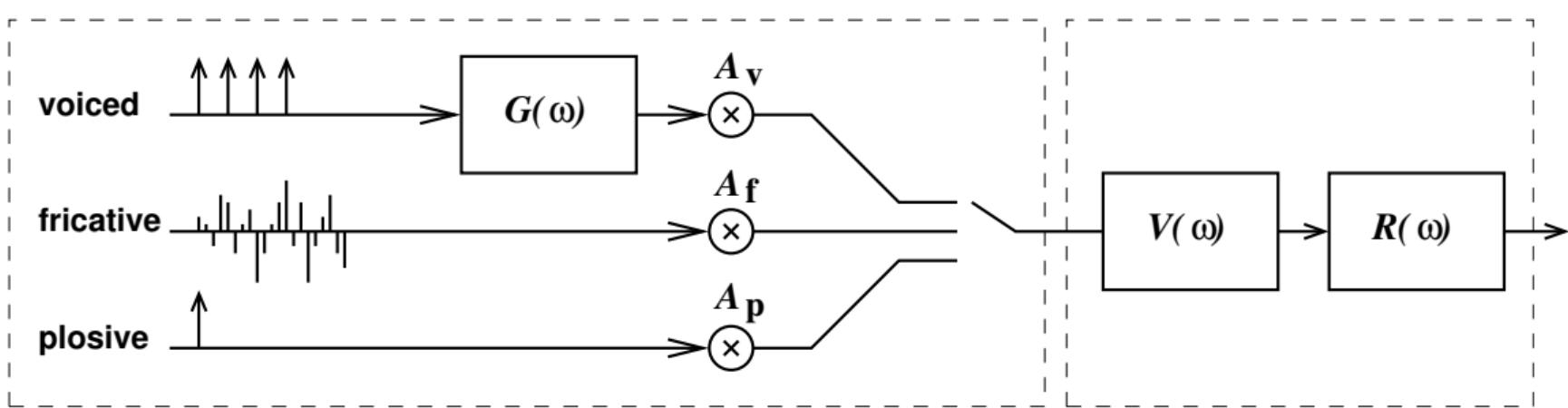


$\leftarrow p[n] * r[n] * v[n]$



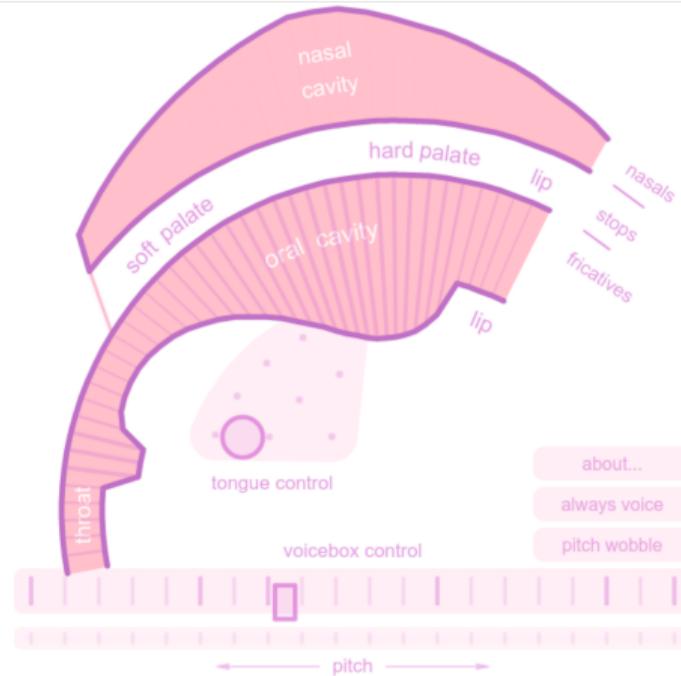
Complete Source/Filter Model

Source



Pink Trombone!

<http://dood.al/pinktrombone/>



IPA Chart: Consonants

THE INTERNATIONAL PHONETIC ALPHABET (2005)

CONSONANTS (PULMONIC)

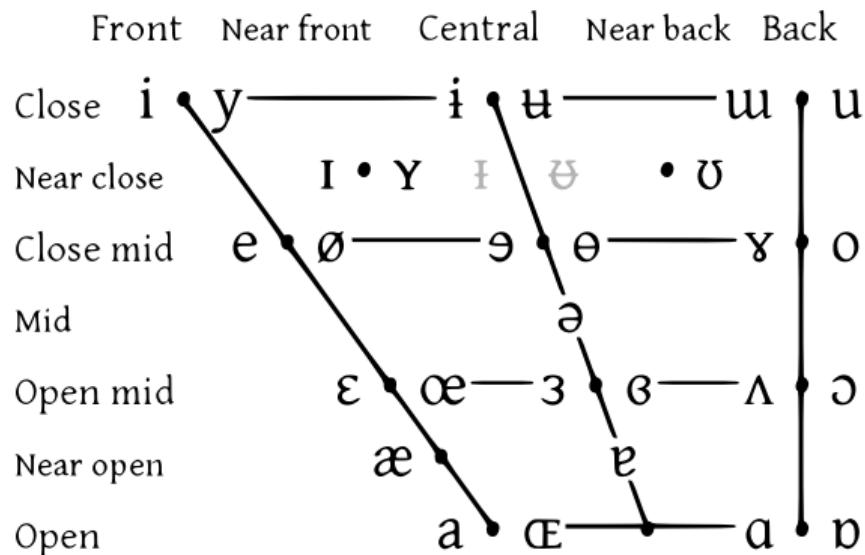
	LABIAL		CORONAL				DORSAL			RADICAL		LARYNGEAL
	Bilabial	Labio-dental	Dental	Alveolar	Palato-alveolar	Retroflex	Palatal	Velar	Uvular	Pharyngeal	Epi-glottal	Glottal
Nasal	m	n̪		n		ɳ	ɲ	ŋ	ɳ			
Plosive	p b	ɸ β		t d		t̪ d̪	c ɟ	k ɡ	q ɢ		ʔ	ʔ
Fricative	ɸ β	f v	θ ð	s z	ʃ ʒ	ʂ ʐ	ç ɟ	x ɣ	χ ʁ	h	h	h
Approximant		v		ɹ		ɬ	j	w		l̪	ɫ	h̪
Trill	B			r						R		Я
Tap, Flap		v		t̪		t̪						
Lateral fricative			ɬ	ɺ	ɭ	ɭ	ʎ	ɻ				
Lateral approximant			l̪	l̪	ɭ	ɭ	ʎ	ɻ				
Lateral flap			ɬ	ɬ	ɭ	ɭ						

Where symbols appear in pairs, the one to the right represents a modally voiced consonant, except for murmured h̪.
 Shaded areas denote articulations judged to be impossible. Light grey letters are unofficial extensions of the IPA.

IPA Chart: Vowels

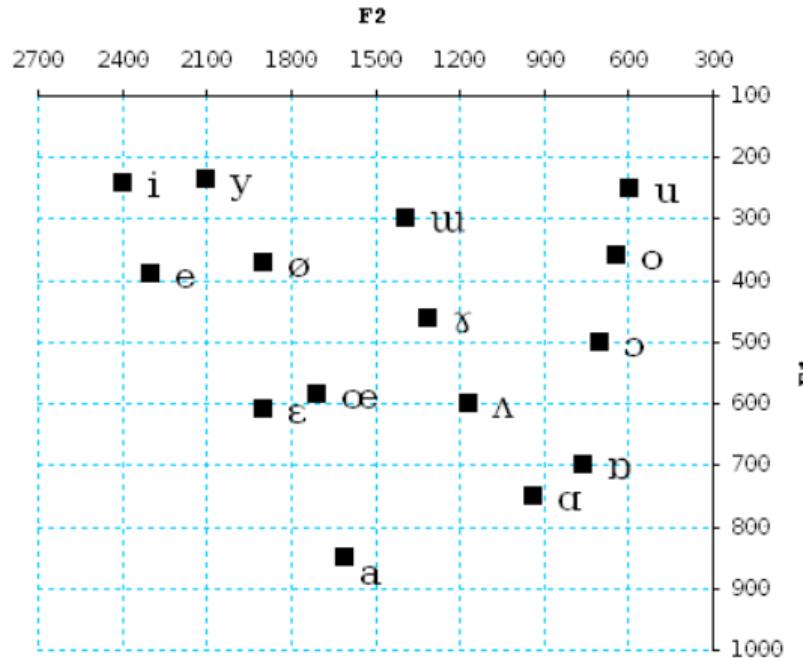
THE INTERNATIONAL PHONETIC ALPHABET (2005)

VOWELS



Vowels at right & left of bullets are rounded & unrounded.

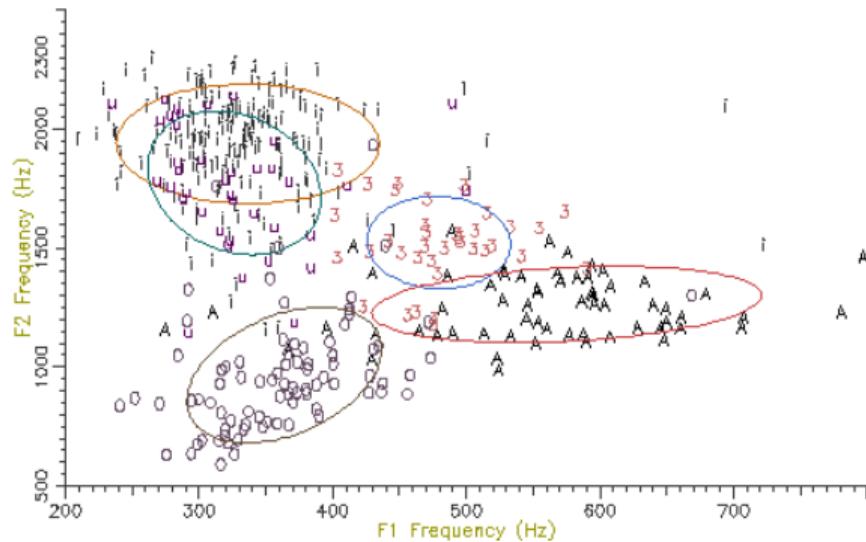
Average formant frequencies (vowels)



Source: <https://en.wikipedia.org/wiki/Formant>

Formant frequencies, single speaker

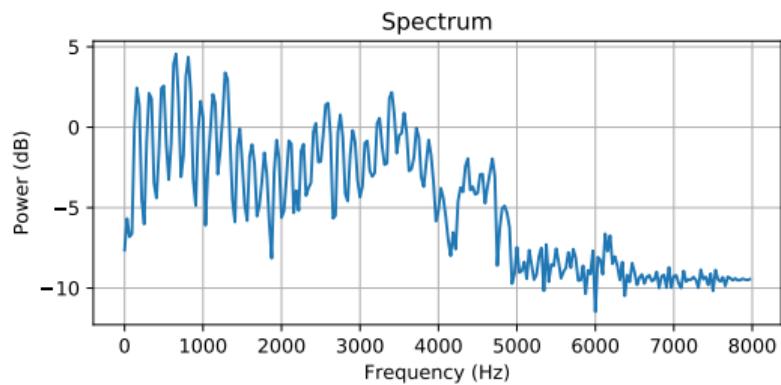
Vowel variability
single speaker, all contexts



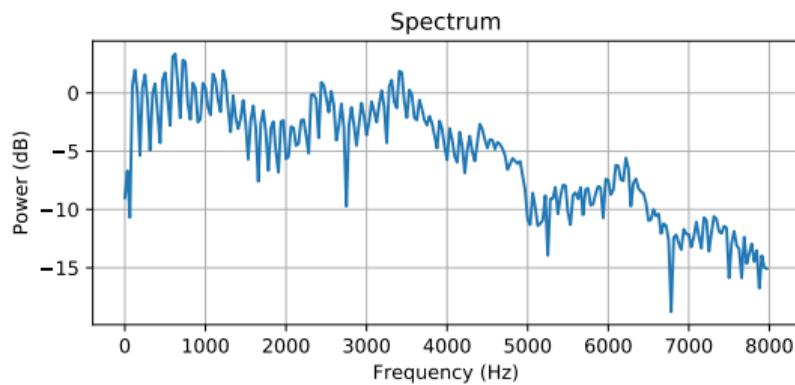
Source: <https://www.phon.ucl.ac.uk/resource/sfs/howto/formant.php>

Question: which is the speech sound with higher f_0 ?

a)

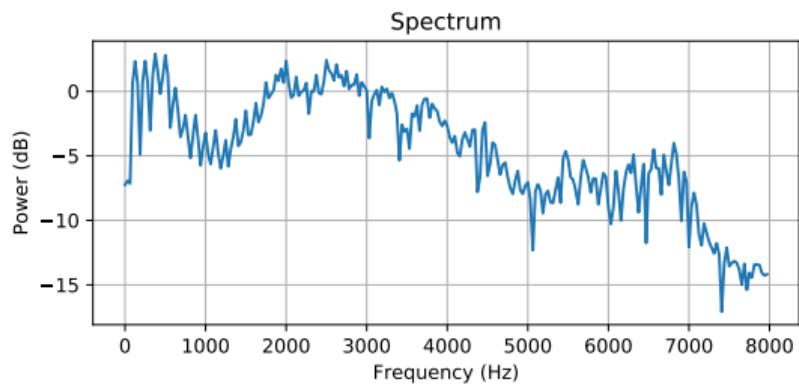


b)

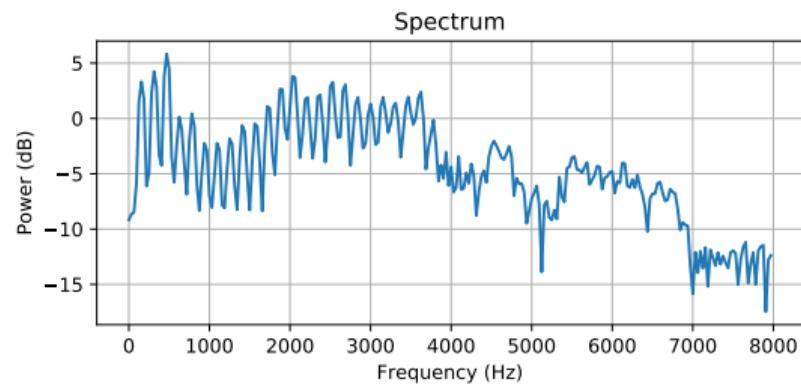


Question: which is the speech sound with higher f_0 ?

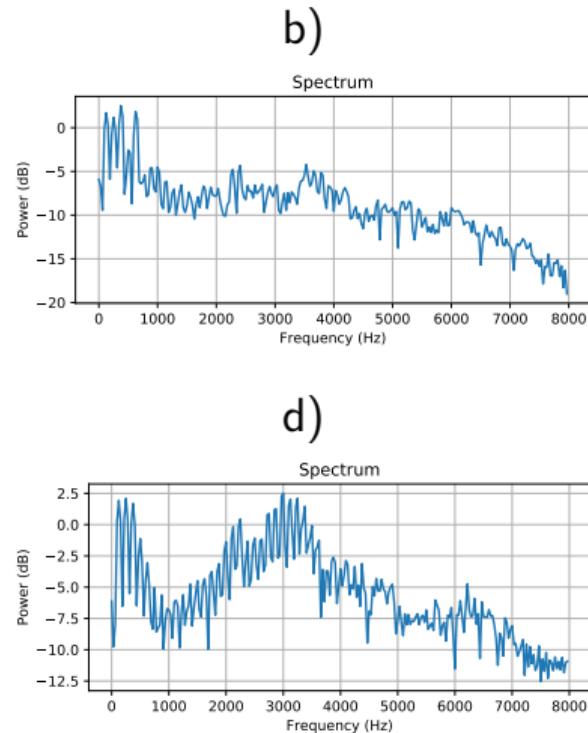
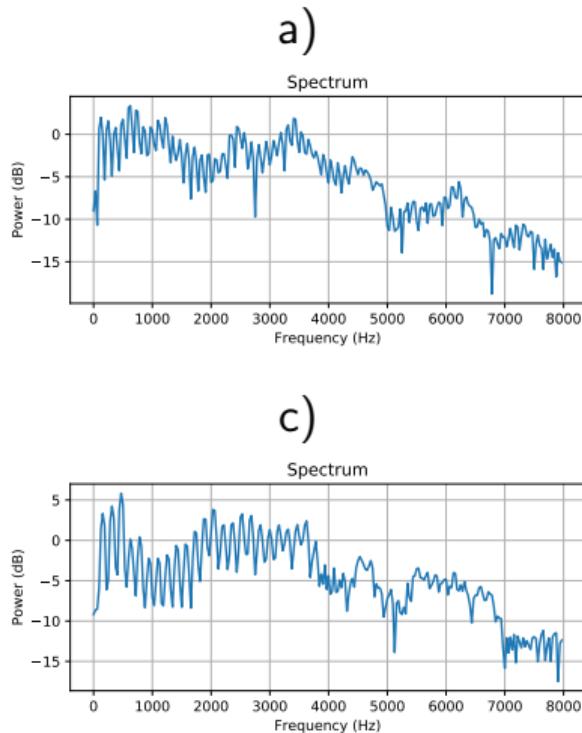
a)



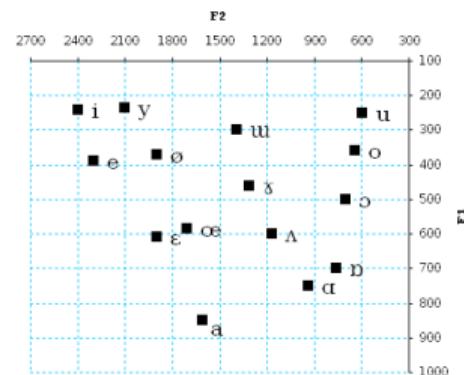
b)



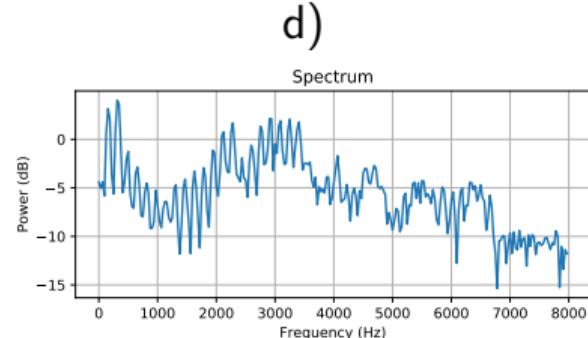
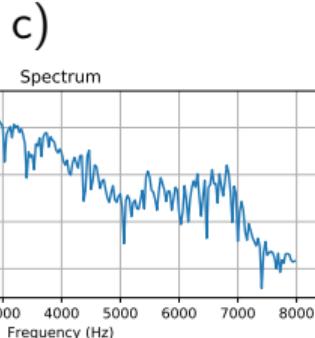
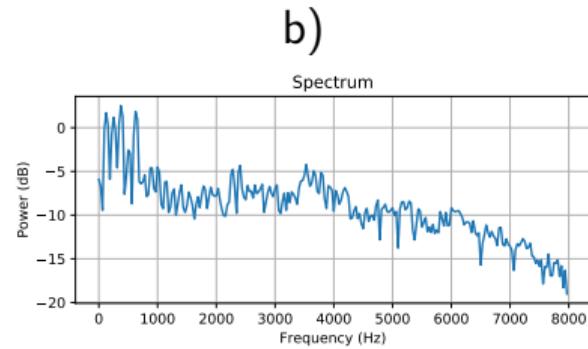
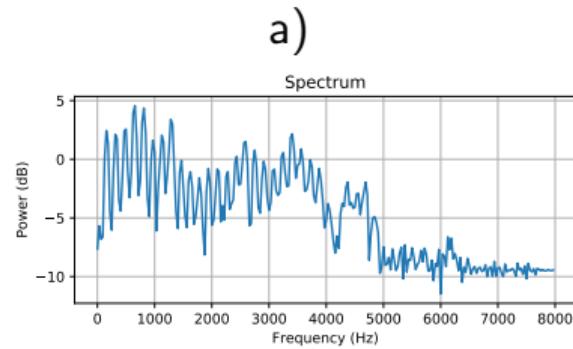
Question: which spectrogram depicts an /a/?



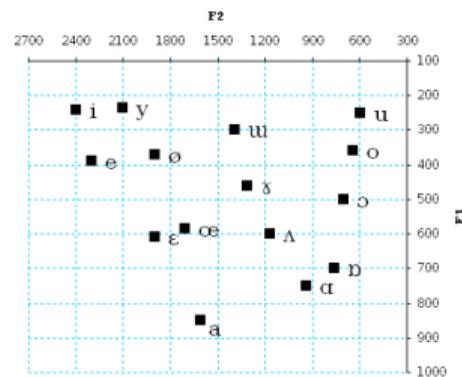
Average formant frequencies



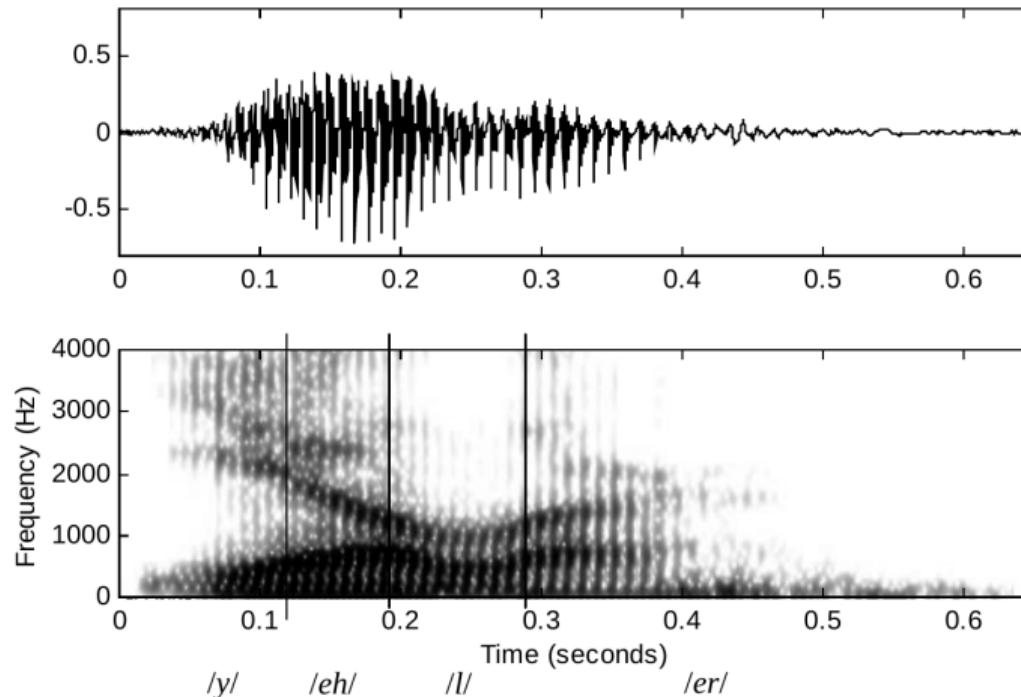
Question: which spectrogram depicts an /i/?



Average formant frequencies

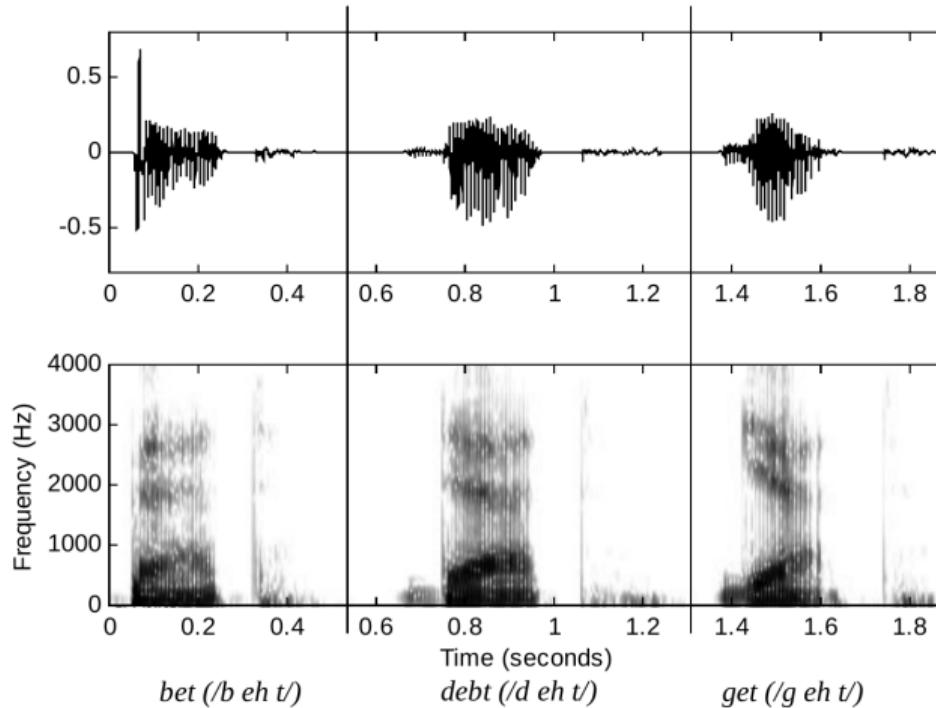


Coarticulation: vowels-semivowels



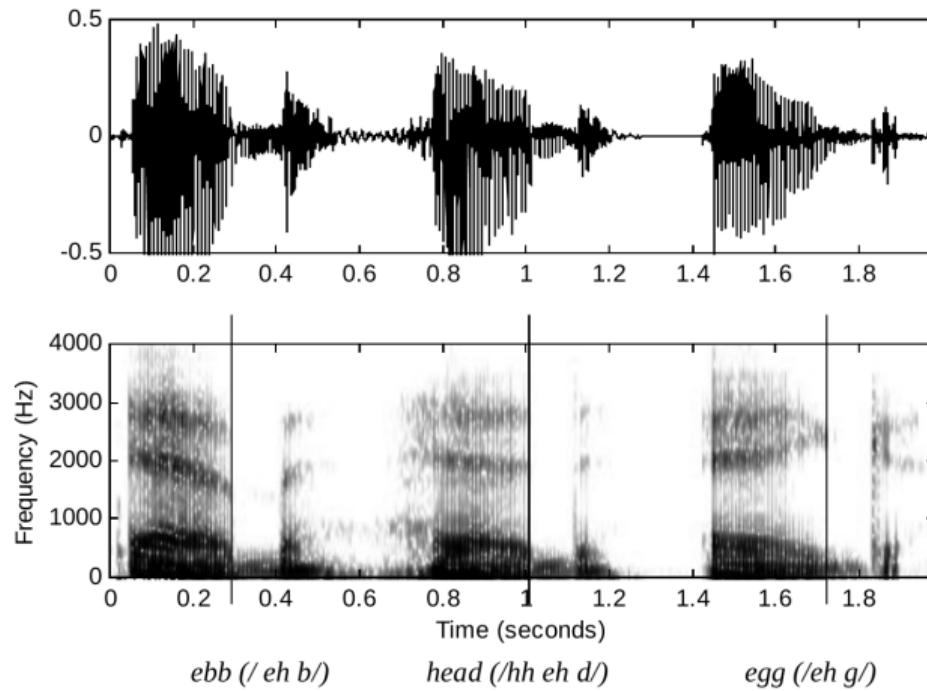
from Huang, Acero and Hon's book

Coarticulation: plosive-vowel



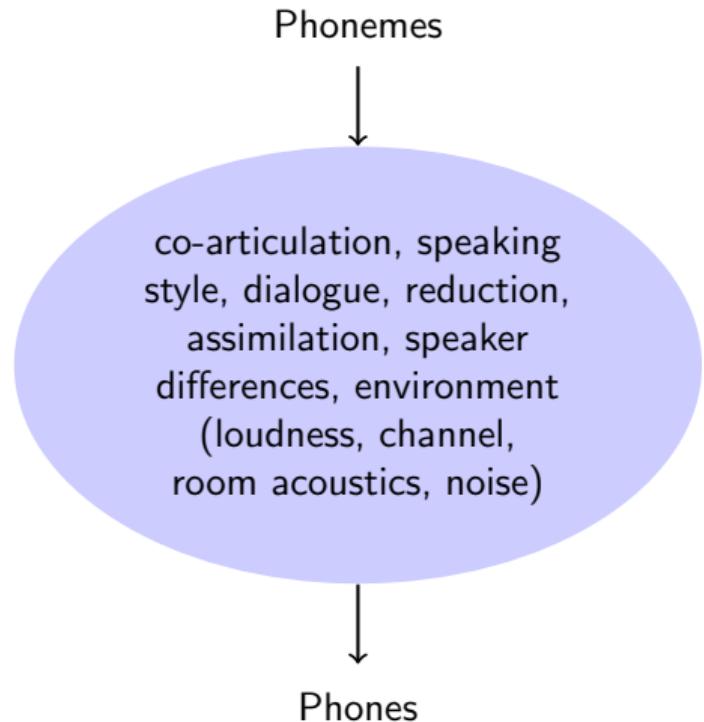
from Huang, Acero and Hon's book

Coarticulation: vowel-plosive

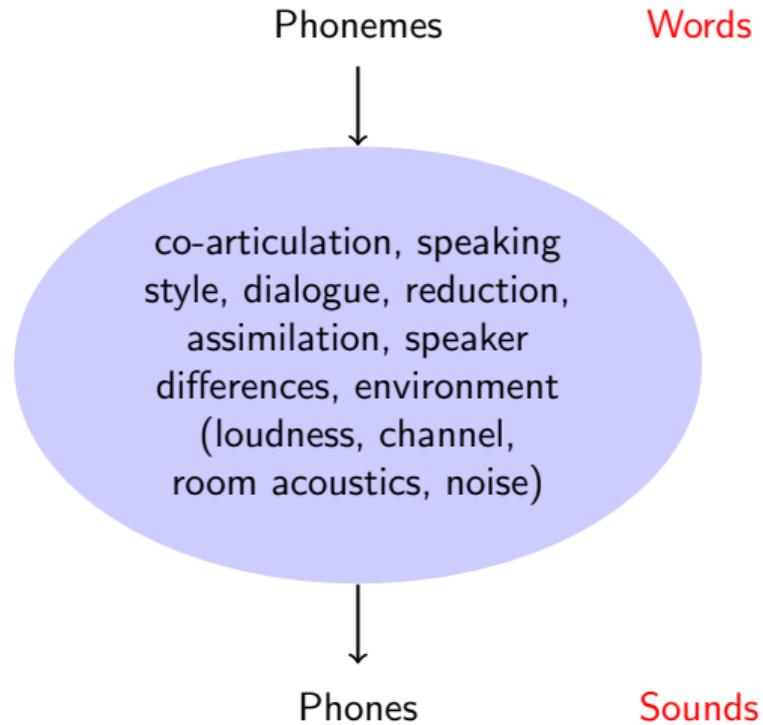


from Huang, Acero and Hon's book

Phonology vs Phonetics



Phonology vs Phonetics



Outline

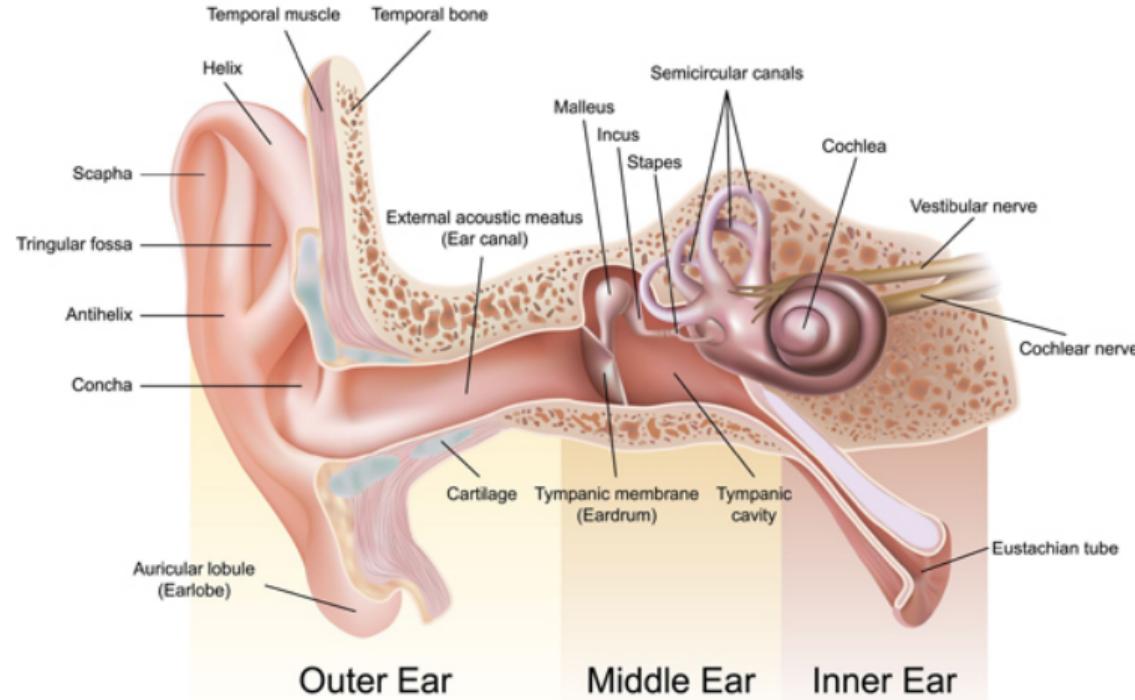
1 Speech Production

- Source/Filter Model

2 Speech Perception

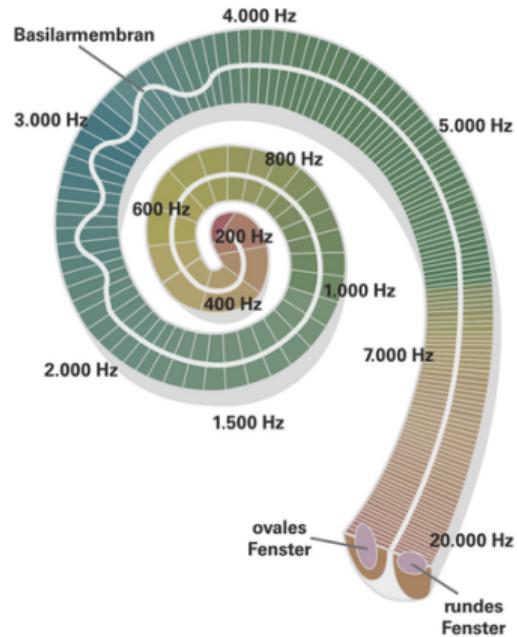
3 Challenges

Anatomy of the ear



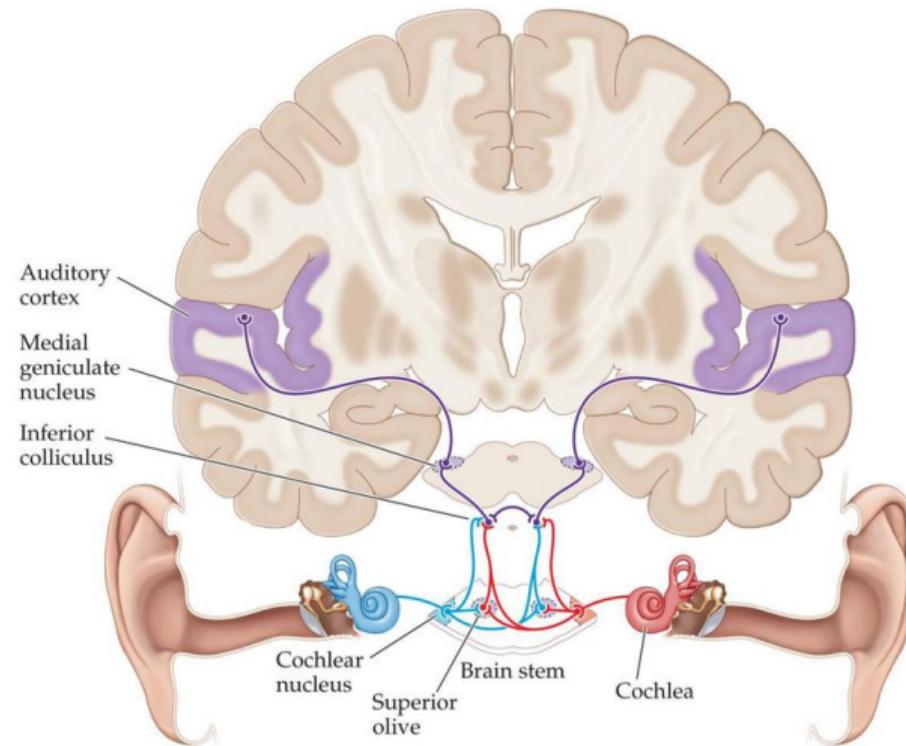
from <https://drjillgordon.com/how-hearing-works>

Cochlea



from https://tu-dresden.de/ing/elektrotechnik/ias/aha/forschung/akustik/psychoakustik?set_language=en

Auditory Pathways

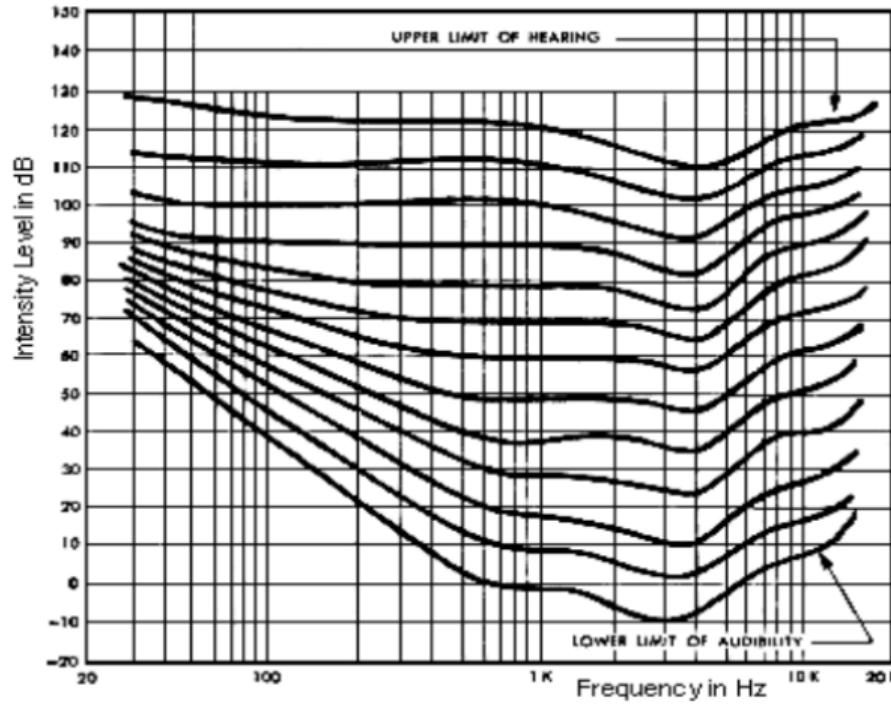


SENSATION & PERCEPTION 5e, Figure 9.20
© 2018 Oxford University Press

Physical vs perceptual attributes

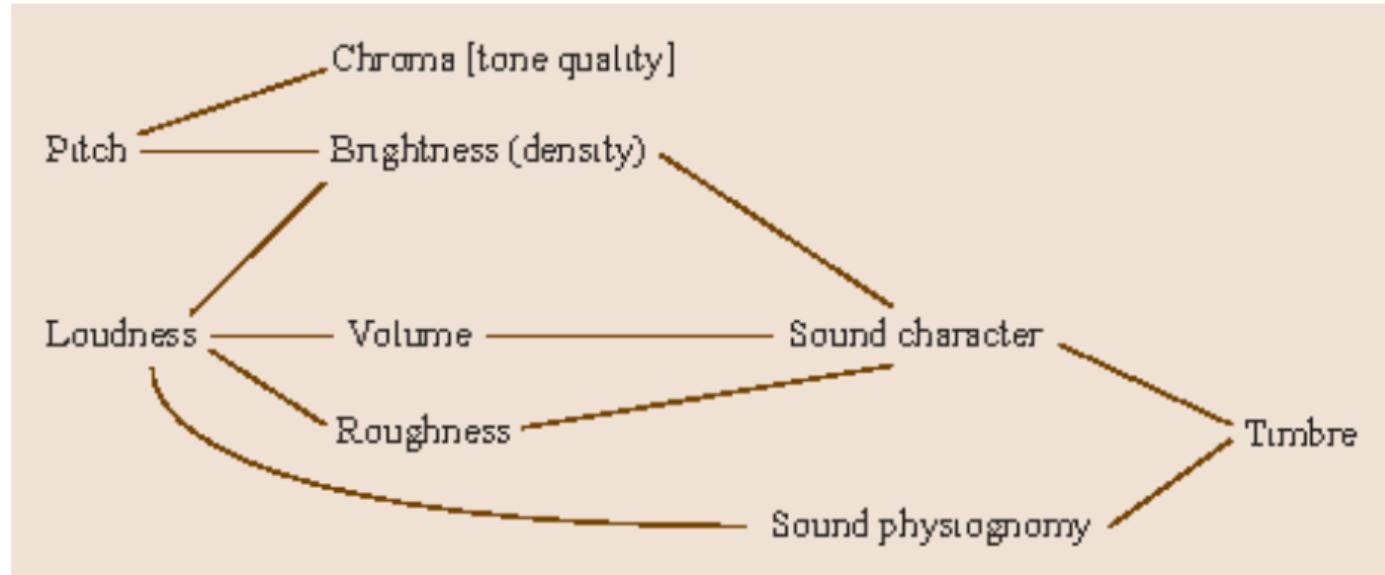
Physical Quantity	Perceptual Quantity
Intensity (SPL)	Loudness
Fundamental frequency (f_0)	Pitch
Spectral shape	Timbre
Onset/offset time	Timing
Phase difference in binaural hearing	Location

Equal-loudness curves (pure tones)



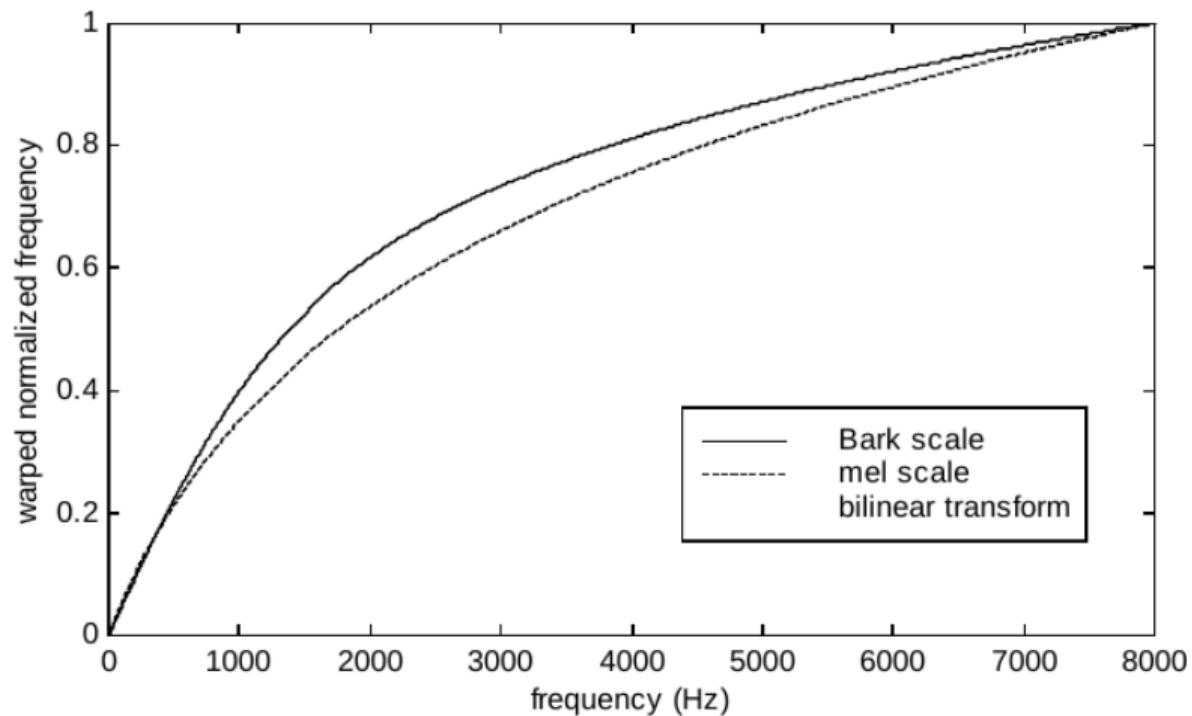
from Huang, Acero and Hon's book

Loudness perception (complex sounds)



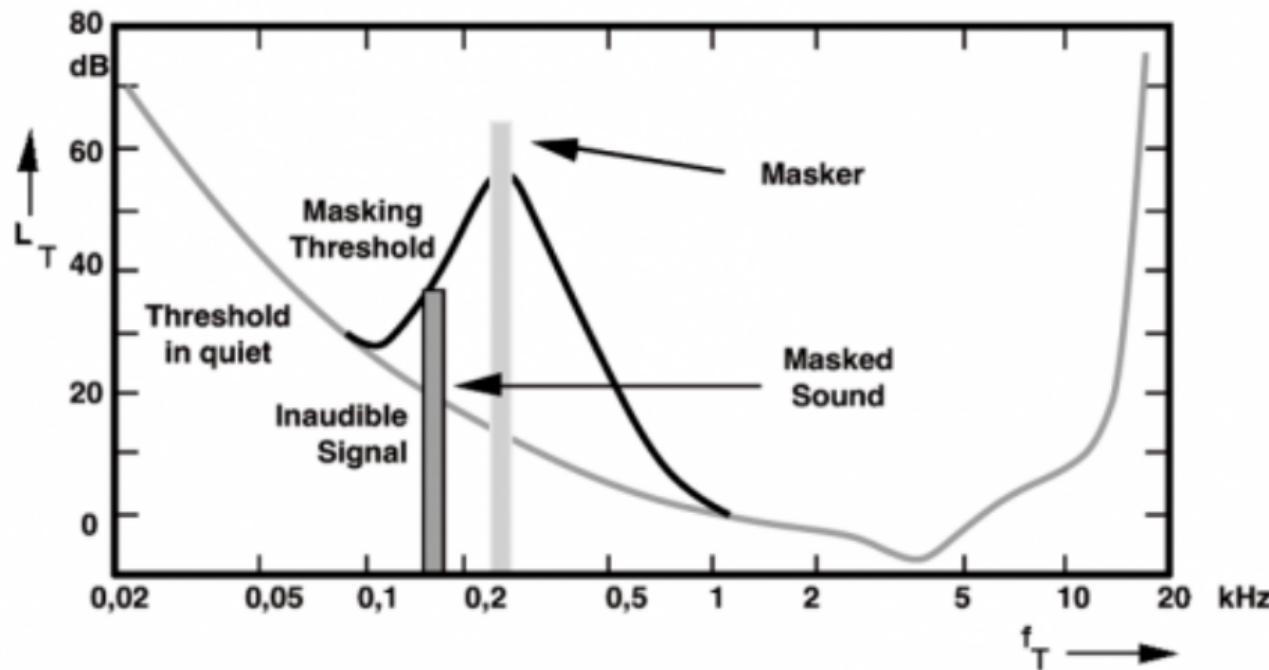
from https://link.springer.com/chapter/10.1007/978-3-662-55004-5_33

Bark and Mel frequency scales



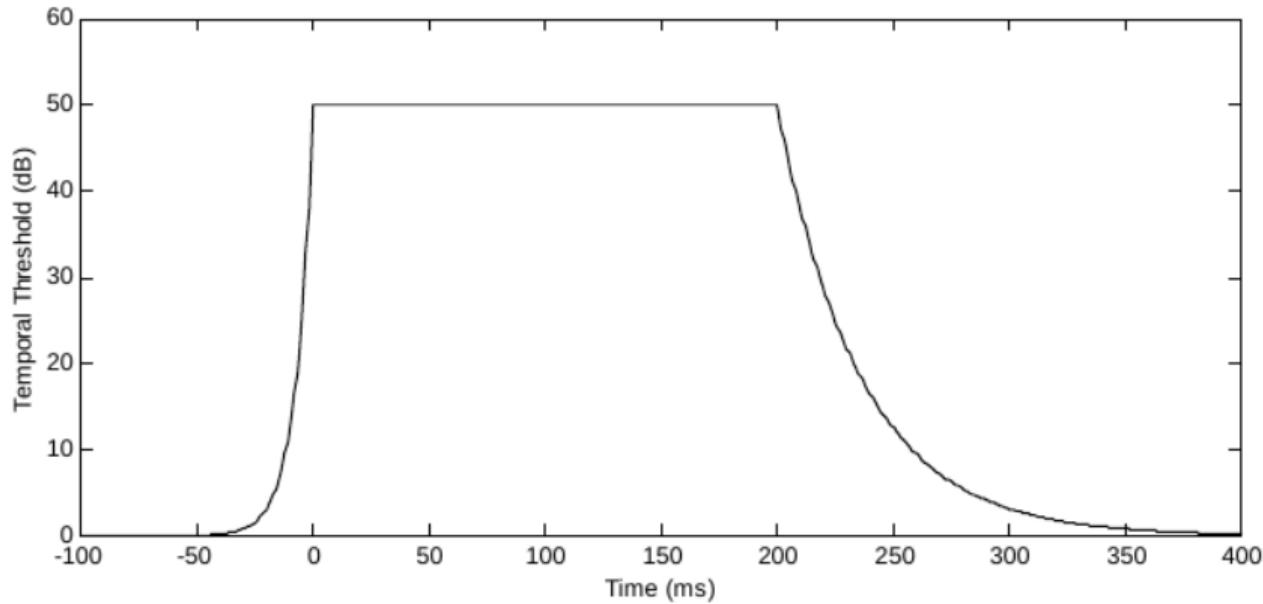
from Huang, Acero and Hon's book

Masking in frequency



from <http://hephaestusaudio.com/delphi/category/psychoacoustics/>

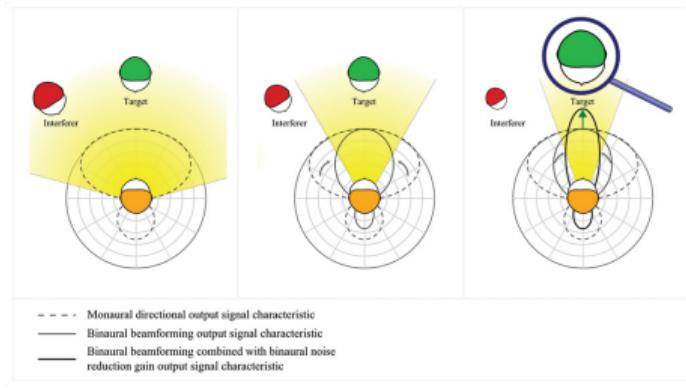
Masking in time



from Huang, Acero and Hon's book

Binaural hearing

- mainly used for localization
- main cue: inter-aural time difference (ITD)
- main cue: inter-aural level difference (ILD)
- filtering from ear pinna
- head-related transfer function HRTF



Check <http://auditoryneuroscience.com/spatial-hearing/binaural-cues>

Outline

1 Speech Production

- Source/Filter Model

2 Speech Perception

3 Challenges

Challenges — Variability

Between speakers

- Age
- Gender
- Anatomy
- Dialect

Within speaker

- Stress
- Emotion
- Health condition
- Read vs Spontaneous
- Adaptation to environment (Lombard effect)
- Adaptation to listener

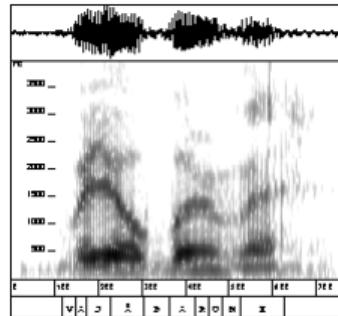
Environment

- Noise
- Room acoustics
- Microphone distance
- Microphone, telephone
- Bandwidth

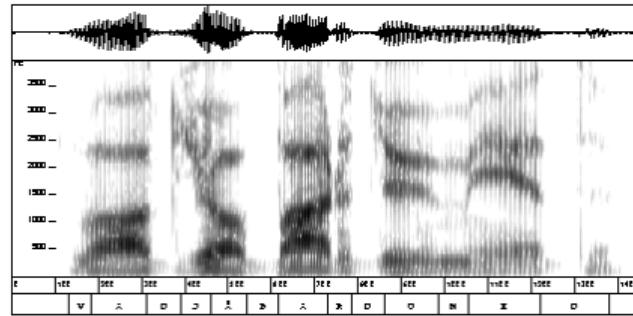
Listener

- Age
- Mother tongue
- Hearing loss
- Known / unknown
- Human / Machine

Example: spontaneous vs hyper-articulated



Va jobbaru me



Vad jobbar du med

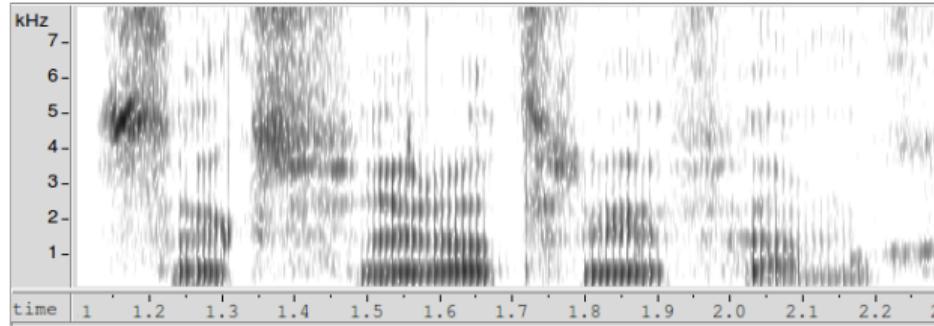
“What is your occupation”
 (“What work you with”)

Examples of reduced pronunciation

Spoken	Written	In English
Tesempel	Till exempel	for example
åhamba	och han bara	and he just
bafatt	bara för att	just because
javende	jag vet inte	I don't know

Microphone distance

Headset



2 m distance

