

DT2119

Speech and Speaker Recognition

Introduction

Giampiero Salvi

KTH/CSC/TMH giampi@kth.se

VT 2017

Outline

Course Organization

Introduction

- The Big Picture
- Challenges

Models of Speech Production

- Source/Filter Model: Vowel-like sounds
- Source/Filter Model, General Case

Outline

Course Organization

Introduction

The Big Picture
Challenges

Models of Speech Production

Source/Filter Model: Vowel-like sounds
Source/Filter Model, General Case

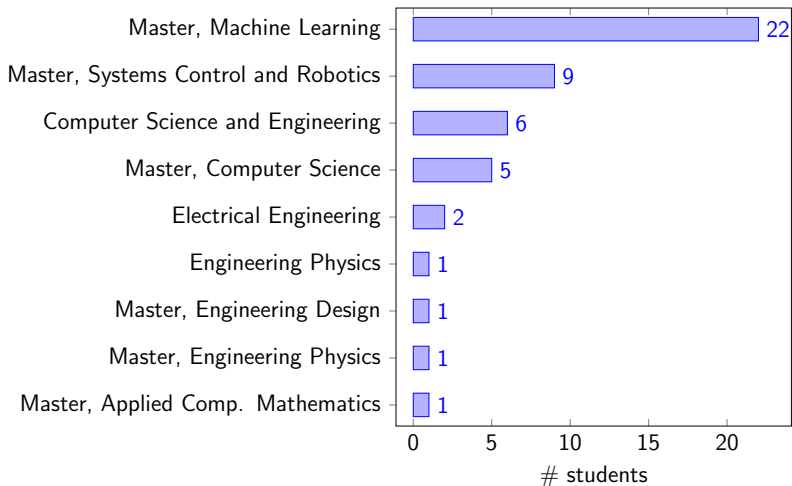
Who are we? (Contact Info)

Giampiero Salvi (giamp@kth.se)

TA: Kalin Stefanov (kalins@kth.se)

All communications handled through Canvas

Who are you?



data from expected participants

Course Objectives

after the course you should be able to:

- ▶ **implement** simple training and evaluation methods for speech recognition
- ▶ **train** and **evaluate** a speech recogniser using software packages
- ▶ **compare** different feature extraction and training methods
- ▶ **document** and **discuss** specific aspects related to speech and speaker recognition
- ▶ with the help of the literature, **review** and **criticise** other students' work in the subject

Course Objectives (research perspective)

- ▶ explore literature
- ▶ carry out experiments
- ▶ produce documentation
- ▶ provide feedback (peer review)
- ▶ accept and use feedback (revision)
- ▶ present results

Topics

Part	Topic	time (hours)
1	Introduction, Speech Signal, Features, Statistics	~ 4
2	Hidden Markov Models, Training and Decoding, Acoustic Models	~ 4-6
3	Deep Learning for ASR	~ 2
4	Decoding and Search Algorithms	~ 2
5	Language Models (Grammars)	~ 2
6	Noise robustness and Speaker Recognition	~ 2-4

Literature

- ▶ **Spoken Language Processing: A Guide to Theory, Algorithm, and System Development**

Xuedong Huang, Alex Acero, Hsiao-Wuen Hon, Prentice Hall

- ▶ 3 at KTH library,
- ▶ 6 at TMH library (against 300 SEK deposit)

- ▶ **Automatic Speech Recognition: A deep learning approach**

Dong Yu and Li Deng, Springer 2015

Available in PDF from SpringerLink (via KTH Biblioteket)

- ▶ **HTK manual** version 3.4
- ▶ selected research articles

Reading Instructions

These are indicative, check the schedule for more updated instructions

		pages	# pages
Part 1	(Spoken Language Structure)	(19–71)	(52)
	Digital Signal Processing	(201–273)	73
	Probability, Statistics and Inform. Theory	73–131	59
	Pattern Recognition	133–197	65
	Speech Signal Representations	275–336	62
Part 2	Hidden Markov Models	377–413	37
	Acoustic Modeling	415–475	61
	Environmental Robustness	477–544	68
Part 3	Deep Neural Nets and ASR	Ch4, Ch6 ¹	35
Part 4	Basic Search Algorithms	591–643	53
	(Large-Vocabulary Search Algorithms)	(645–685)	(41)
	(Applications and User Interfaces)	(919–956)	(38)
Part 5	Language Modeling	545–590	46
Part 6	Speaker Recognition literature		

Dong and Deng's book

Activities

In groups:

1. three **labs** (with presentation to TA)
2. **project work** and **report**
3. project **presentation** at final seminar

Individually:

1. three **quizzes** on Canvas
2. **review** other students' report

Lab 1: Speech Feature Extraction

- ▶ implement extraction for typical speech features
- ▶ analyse the features on speech data
- ▶ compare utterances with Dynamic Time Warping

Lab 2: Gaussian Hidden Markov Models

- ▶ implement the decoding algorithms for HMMs
- ▶ implement the training algorithms for HMMs
- ▶ test the algorithms on isolated digits

Lab 3: Continuous Speech Recognition and Deep Learning

- ▶ Extend the training and testing algorithms to continuous speech
- ▶ test the algorithms on the TIDIGIT database (connected digits)
- ▶ Optional: implement DNNs using Theano, compare with GMM-HMMS

Project report

- ▶ Suggest a title or choose a topic from a list
- ▶ Project report in form of research paper
- ▶ Suggested topics:

Own work and experiments after discussion with the teacher

Limitations in standard HMM and a survey of alternatives

Pronunciation variation and its importance for speech recognition

Language models for speech recognition

New search methods

Techniques for robust recognition of speech

Confidence measures in speech recognition

The role of prosody for speech recognition

Speaker variability and methods for adaptation

Grading Criteria: Prerequisite for Pass

In groups:

1. **present** the three labs
2. carry out **project work**
3. submit **report draft**
4. **present** at final seminar
5. submit **final report**

Individually:

1. carry out three **quizzes** on Canvas
2. **review** other students' report

Grading Criteria

Extended literature study: max grade C

- E:** complete literature study and present it
- D:** give extended feedback to other students
- C:** incorporate feedback you receive

Experimental project: max grade A

- C:** complete experiments and report and present it
- B:** give extended feedback to other studies
- A:** incorporate feedback you receive

Computational Resources at PDC

For Lab 3 and the Project

- ▶ PDC accounts will be created for all registered students
- ▶ alternatively, apply for an account at <https://www.pdc.kth.se/support/accounts/user>
- ▶ use `edu17.DT2119` when asked for time allocation
- ▶ 45 min introduction to PDC in one of the lectures

Time Organisation

Week 12 (March 20): Course start

Week 14 (April 6): Decide groups/project topics

Week 14 (Canvas): Present Lab 1

Week 16 (Canvas): Present Lab 2

Week 18 (Canvas): Present Lab 3

Week 20 (May 19): Submit first version of report

Week 21 (May 26): Submit review on report

Week 22 (May 31): Project presentations (posters)

Week 23 (June 7): Submit final report.

Canvas quizzes can be completed at your discretion

Part 1

Outline

Course Organization

Introduction

- The Big Picture
- Challenges

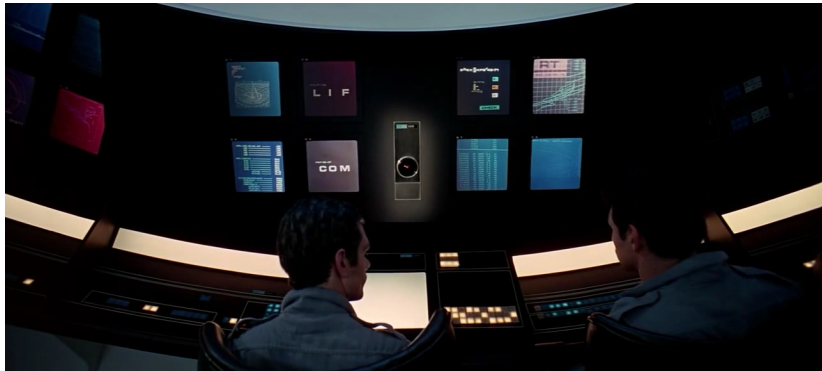
Models of Speech Production

- Source/Filter Model: Vowel-like sounds
- Source/Filter Model, General Case

Motivation

- ▶ 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)

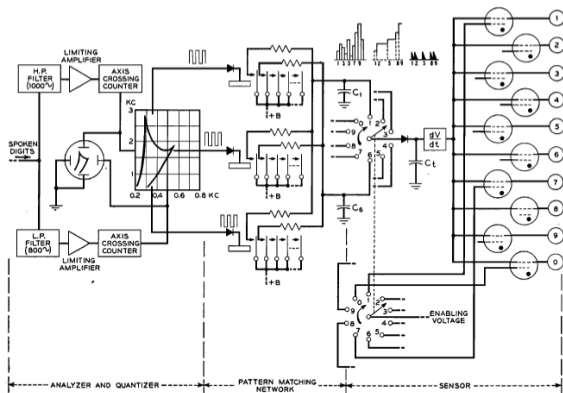
The dream of Artificial Intelligence



2001: A space odyssey (1968)

A very long endeavour

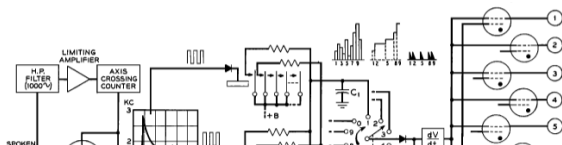
1952, Bell laboratories, isolated digit recognition,
single speaker, hardware based [1]



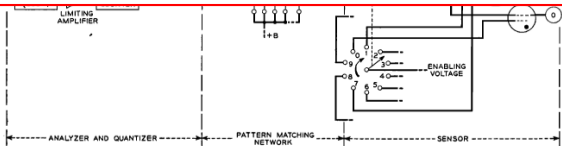
[1] K. H. Davis, R. Biddulph, and S. Balashek. "Automatic Recognition of Spoken Digits". In: *JASA* 24.6 (1952), pp. 637–642

A very long endeavour

1952, Bell laboratories, isolated digit recognition,
single speaker, hardware based [1]



An underestimated challenge:
60 years of bold announcements



[1] K. H. Davis, R. Biddulph, and S. Balashek. "Automatic Recognition of Spoken Digits". In: *JASA* 24.6 (1952), pp. 637–642

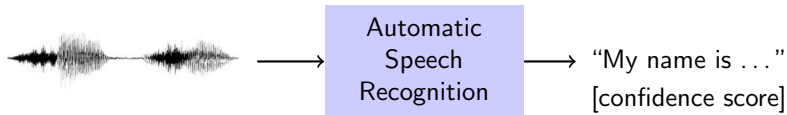
Today's Reality



I Now Pronounce You Chuck & Larry (2007)

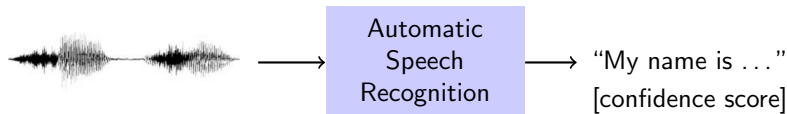
The ASR Goal (for this course)

Convert speech into text



The ASR Goal (for this course)

Convert speech into text



CC Please tell me your name

LV Larry Valentine

CC I'm sorry, I didn't quite get that

LV Larry Valentine

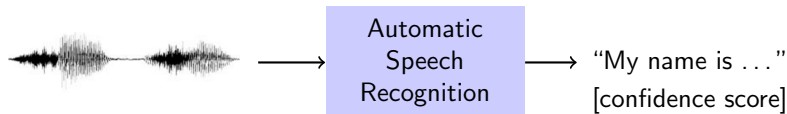
CC You said "Berry Schmallenpine" ... is that right?

LV Schmallenpine?!?!

CC You said "Schmallenpine" ... is that right?

The ASR Goal (for this course)

Convert speech into text



CC Please tell me your name

LV Larry Valentine

CC I'm sorry, I didn't quite get that

LV Larry Valentine

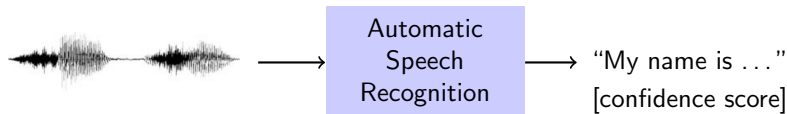
CC You said "Berry Schmallenpine" ... is that right?

LV Schmallenpine?!?!

CC You said "Schmallenpine" ... is that right?

The ASR Goal (for this course)

Convert speech into text



CC Please tell me your name

LV Larry Valentine

CC I'm sorry, I didn't quite get that

LV Larry Valentine

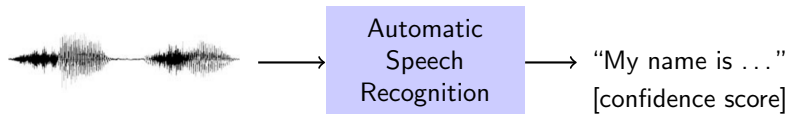
CC You said "Berry Schmallenpine" ... is that right?

LV Schmallenpine?!?!

CC You said "Schmallenpine" ... is that right?

The ASR Goal (for this course)

Convert speech into text



CC Please tell me your name

LV Larry Valentine

CC I'm sorry, I didn't quite get that

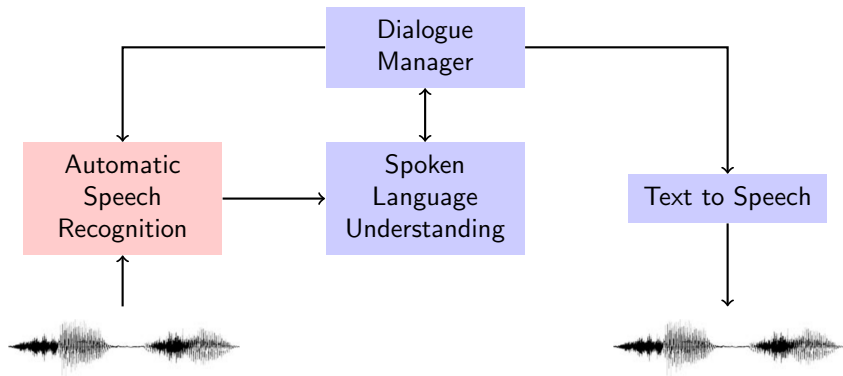
LV Larry Valentine

CC You said "Berry Schmallenpine" ... is that right?

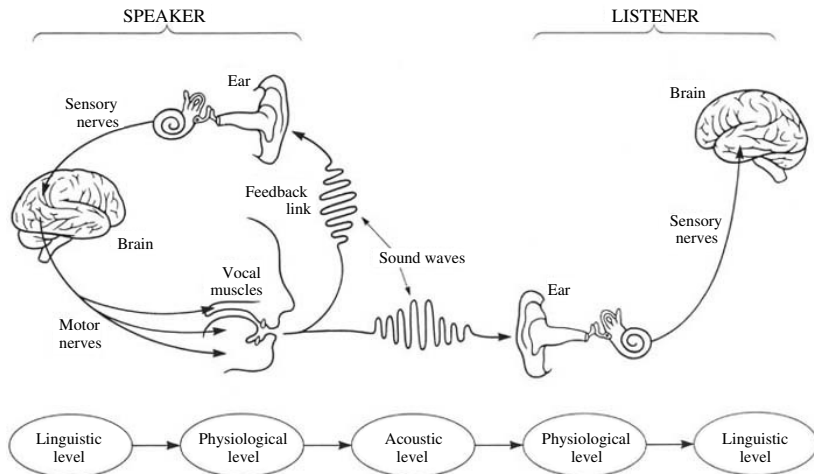
LV Schmallenpine?!?!

CC You said "Schmallenpine" ... is that right?

ASR in a Broader Context



The Speech Chain



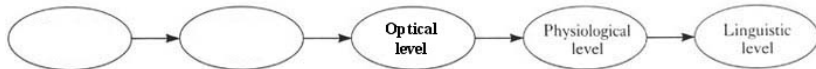
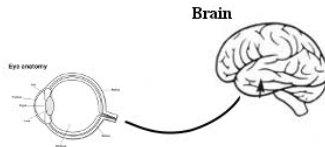
Peter Denes, Elliot Pinson, 1963

ASR versus Computer Vision

SCENE



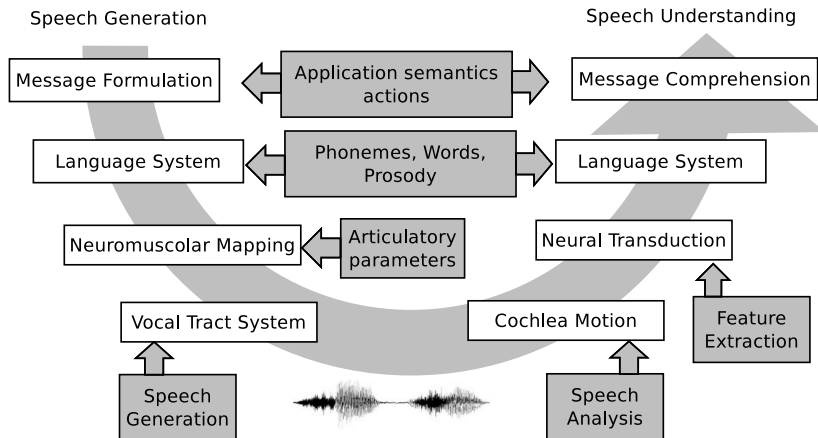
OBSERVER



ASR versus Computer Vision

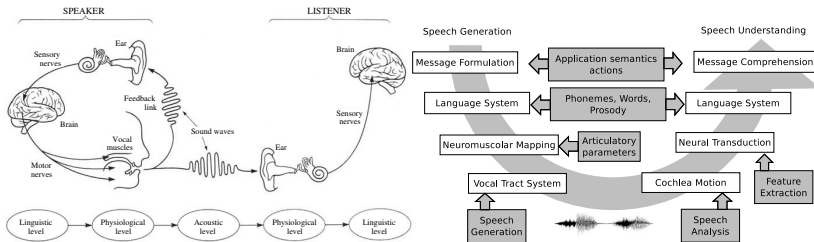
Property	ASR	Computer Vision
signal originates from:	cognition + physics	physics
persistence:	disappears as soon as heard	continually available (active perception)
across countries:	different languages	same objects
type of interaction:	two-way	one-way

The Speech Chain (from the book)

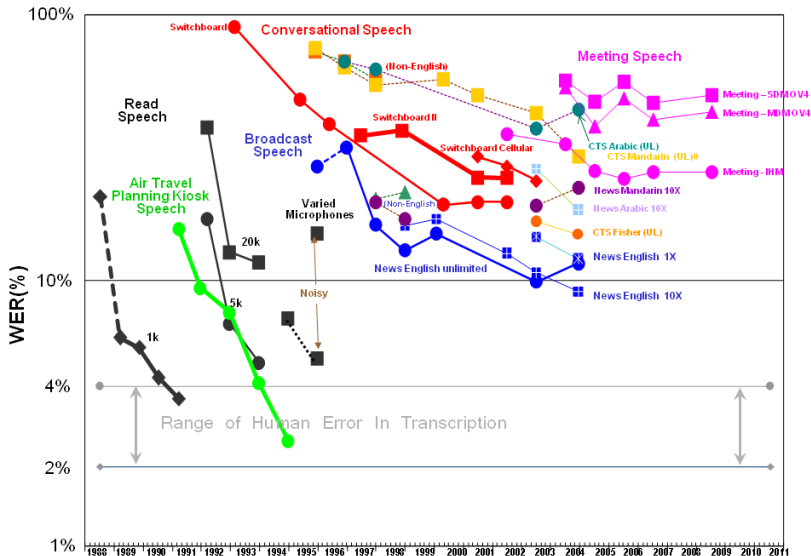


Not covered in this course:

- ▶ multimodality
- ▶ interaction (bi-directional)
- ▶ incrementality
- ▶ non-verbal communication



NIST STT Benchmark Test History – May. '09



<http://www.itl.nist.gov/iad/mig/publications/ASRhistory/>

Main variables in ASR

Speaking mode isolated words vs continuous speech

Speaking style read speech vs spontaneous speech

Speakers speaker dependent vs speaker independent

Vocabulary small (<20 words) vs large ($>50\,000$ words)

Robustness against background noise

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

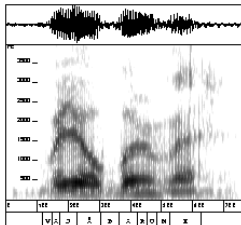
Environment

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

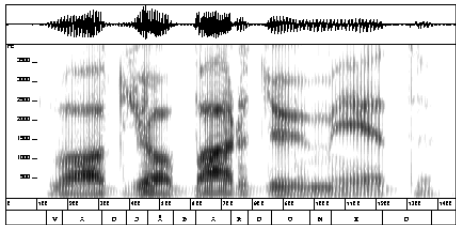
Listener

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

Example: spontaneous vs hyper-articulated



Va jobbaru me



Vad jobbar du med

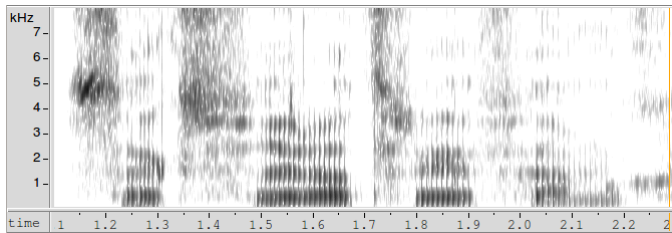
“What is your occupation”
(“What work you with”)

Examples of reduced pronunciation

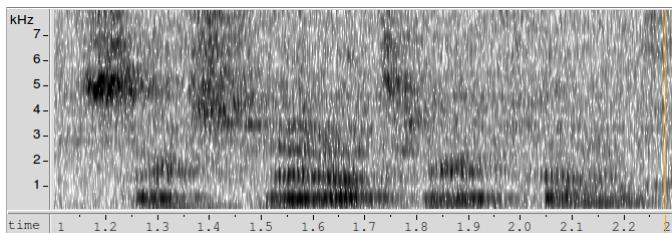
Spoken	Written	In English
Tesempel	Till exempel	for example
åhamba	och han bara	and he just
bafatt	bara för att	just because
javende	jag vet inte	I don't know

Microphone distance

Headset



2 m distance



Applications today

Call centers:

- ▶ traffic information
- ▶ time-tables
- ▶ booking...

Accessibility

- ▶ Dictation
- ▶ hand-free control (TV, video, telephone)

Smart phones

- ▶ Siri, Android...

Smart speakers

- ▶ Amazon Echo...

Outline

Course Organization

Introduction

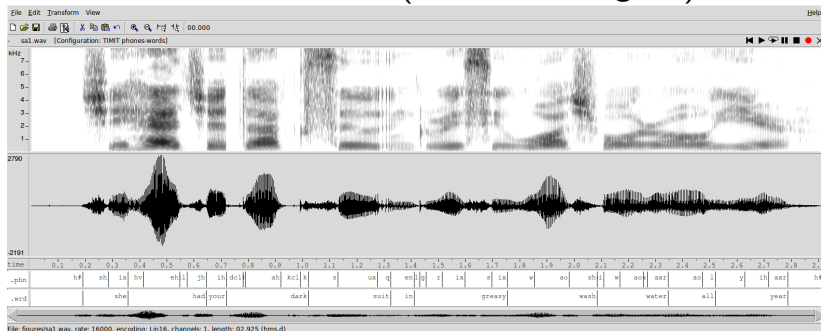
The Big Picture
Challenges

Models of Speech Production

Source/Filter Model: Vowel-like sounds
Source/Filter Model, General Case

Speech Examples

TIMIT database (American English)



example of “clean” speech

Elements of Signal Processing

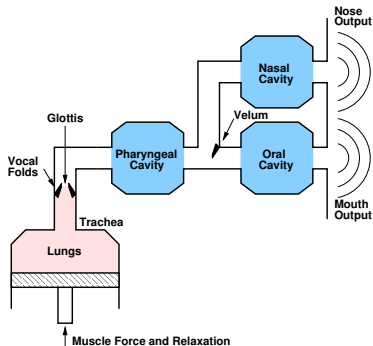
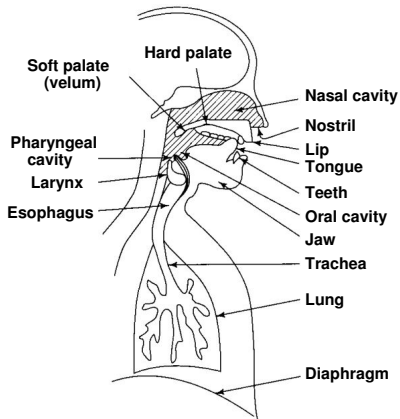
- ▶ continuous/digital signals
- ▶ Linear and Time Invariant (LTI) systems
- ▶ impulse response and convolution
- ▶ Fourier transform and transfer function
- ▶ sampling theorem
- ▶ short-time Fourier transform

(Chapter 5 in the book)

Speech Examples

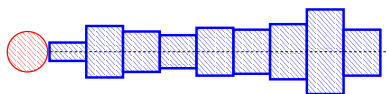
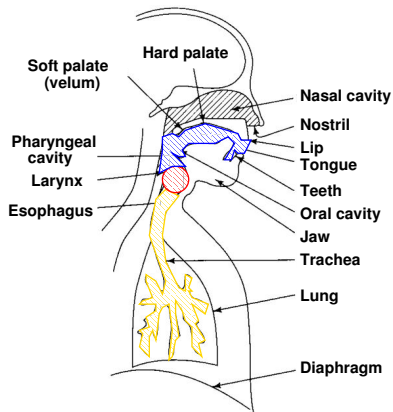
live examples





Physiology



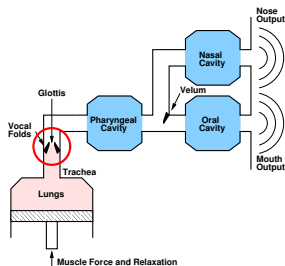
Source/Filter Model, Vowel-like sounds

Vowels

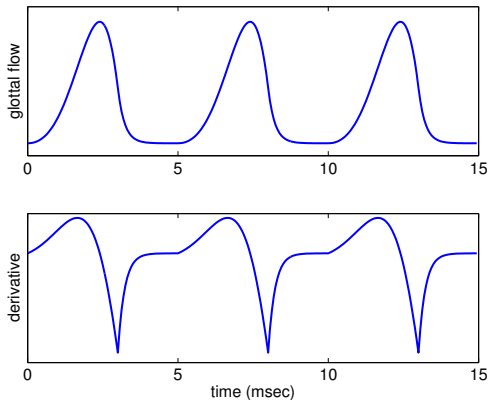


-  Source (periodic)
-  Front Cavity
-  Back Cavity
-  Back Cavity (2nd approx.)

Glottal Flow

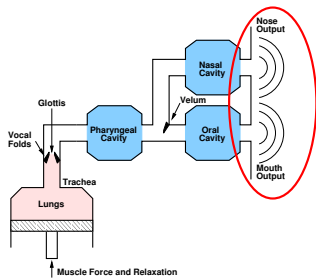


Liljencrants–Fant glottal model



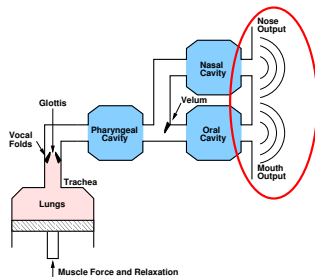
$$G(z) = \frac{1}{(1 - \beta z)^2}, \quad \beta < 1$$

Radiation from the Lips/Nose



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

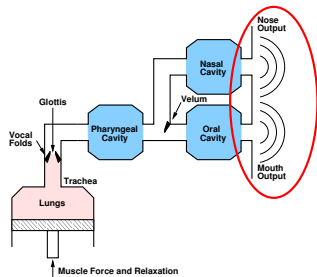
Radiation from the Lips/Nose



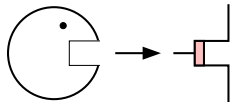
Approx. with a piston in a rigid sphere: solved but not in closed form



Radiation from the Lips/Nose

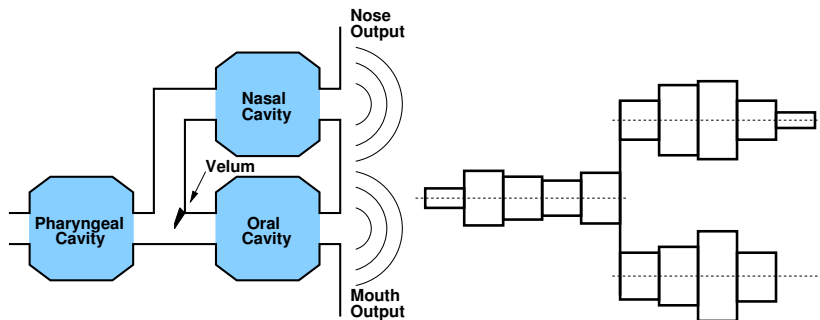


2nd approx: piston in an infinite wall

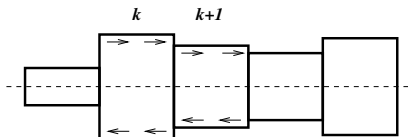


$$R(z) \approx 1 - \alpha z^{-1}$$

Tube Model of the Vocal Tract



Tube Model (cntd.)

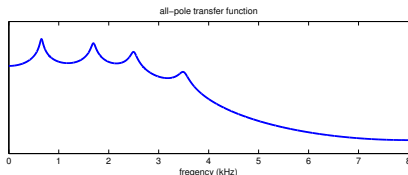
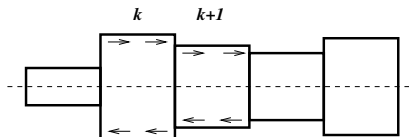


- ▶ assume planar wave propagation and lossless tubes
- ▶ solve pressure $p(x, t)$ and velocity $u(x, t)$ in each tube according to wave equation
- ▶ impose continuity of pressure and velocity at the junctions

⇒ all-pole transfer function (N = number of tubes)

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

Tube Model (cntd.)

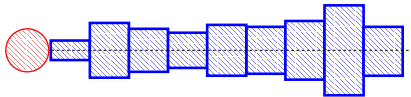
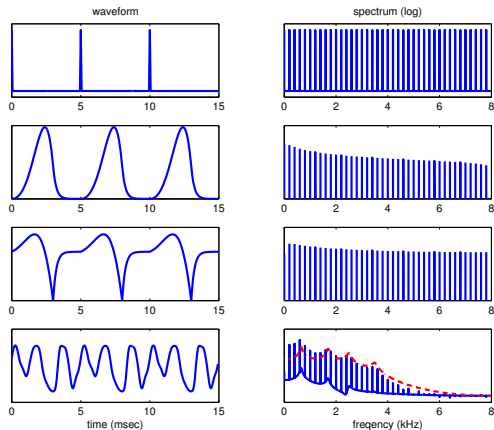


- ▶ assume planar wave propagation and lossless tubes
- ▶ solve pressure $p(x, t)$ and velocity $u(x, t)$ in each tube according to wave equation
- ▶ impose continuity of pressure and velocity at the junctions

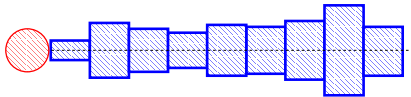
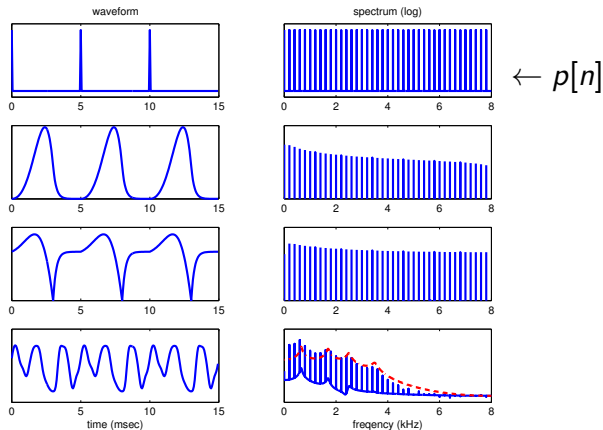
⇒ all-pole transfer function (N = number of tubes)

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

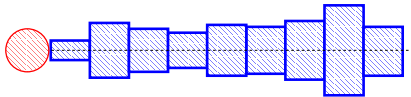
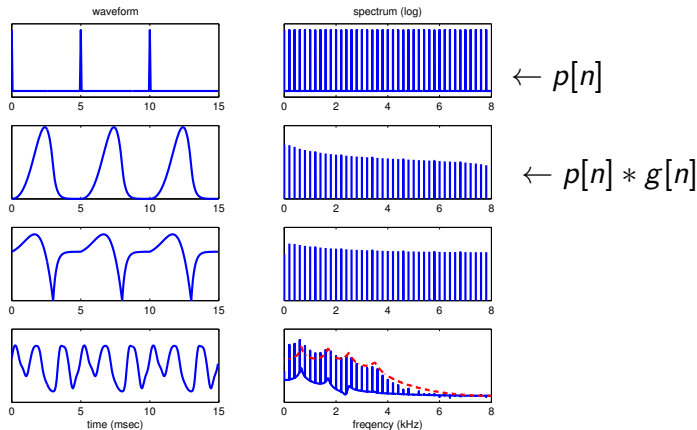
Source/Filter Model: vowel-like sounds



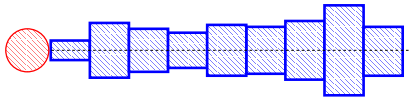
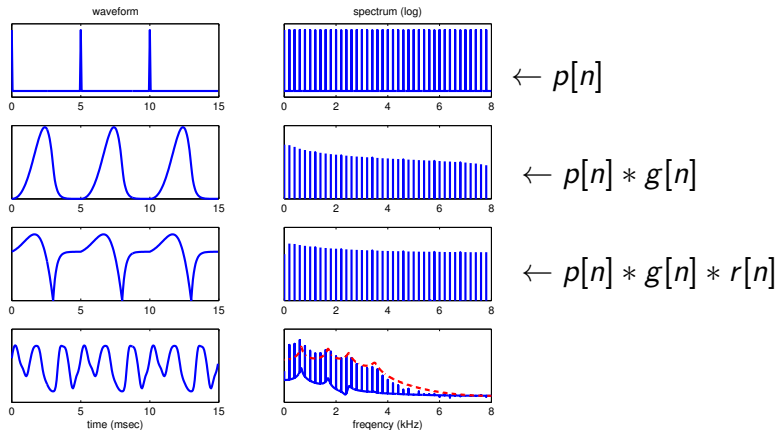
Source/Filter Model: vowel-like sounds



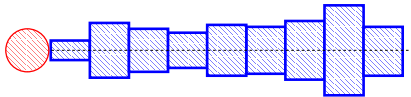
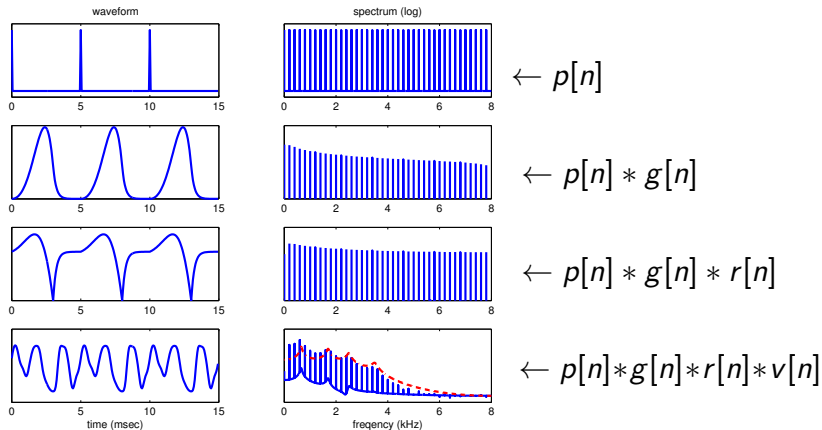
Source/Filter Model: vowel-like sounds



Source/Filter Model: vowel-like sounds

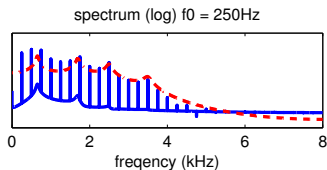
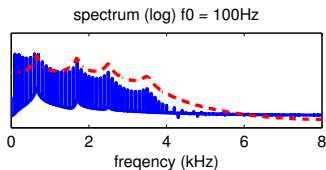


Source/Filter Model: vowel-like sounds



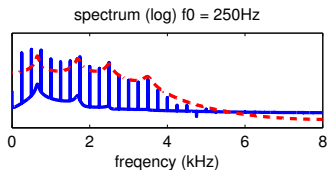
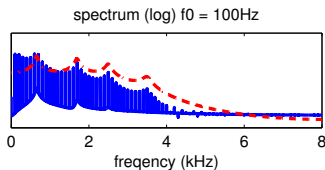
F_0 and Formants

- Varying F_0 (vocal fold oscillation rate)

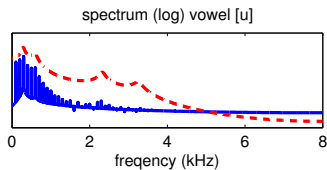
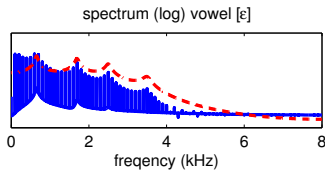


F_0 and Formants

- Varying F_0 (vocal fold oscillation rate)

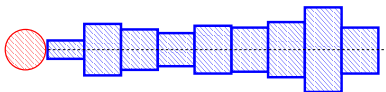
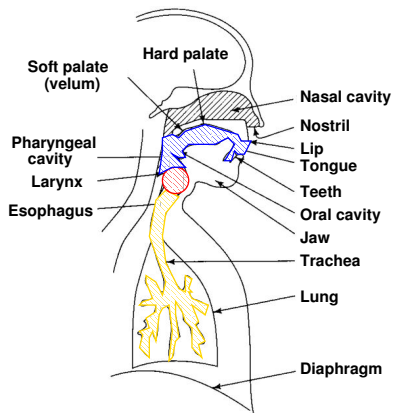




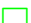

- Varying Formants (vocal tract shape)



Source/Filter Model, General Case

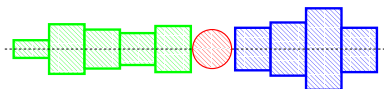
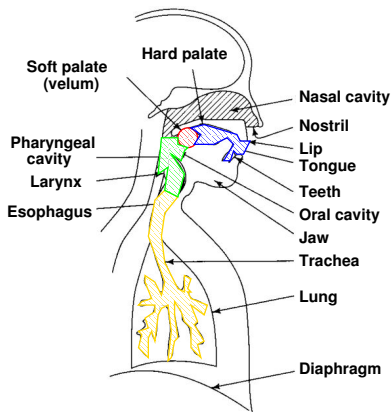
Vowels



-  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)

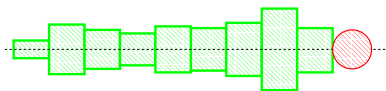
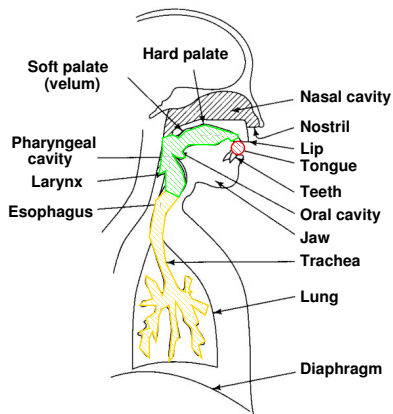
□ Front Cavity

□ Back Cavity

□ Back Cavity (2nd approx.)

Source/Filter Model, General Case

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



□ Source (noise or impulsive)

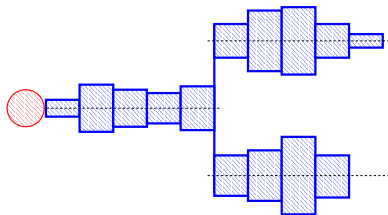
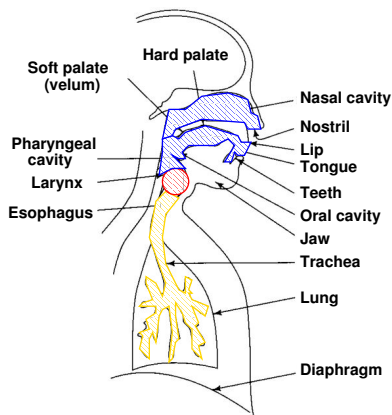
□ Front Cavity

□ Back Cavity

□ Back Cavity (2nd approx.)

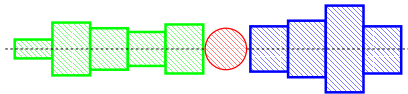
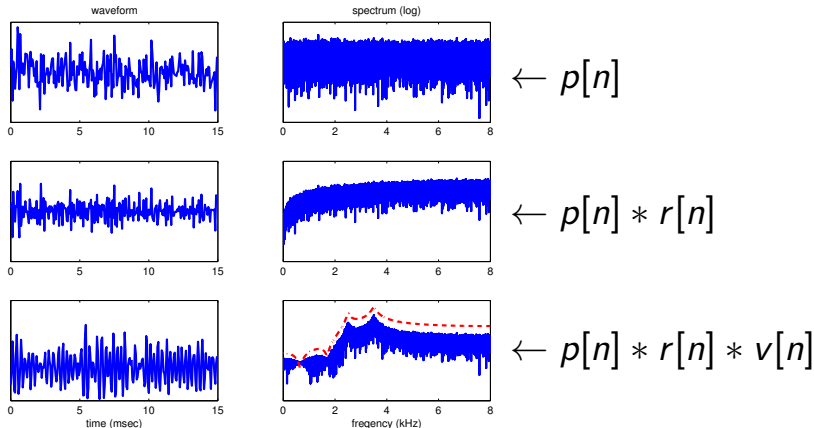
Source/Filter Model, General Case

Nasalised Vowels

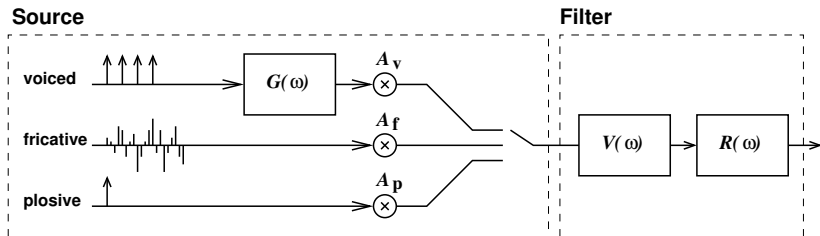


- Source (periodic)
- Front Cavity
- Back Cavity
- Back Cavity (2nd approx.)

Source/Filter Model: fricative sounds



Complete Source/Filter Model



IPA Chart: Consonants

THE INTERNATIONAL PHONETIC ALPHABET (2005)

CONSONANTS (PULMONIC)

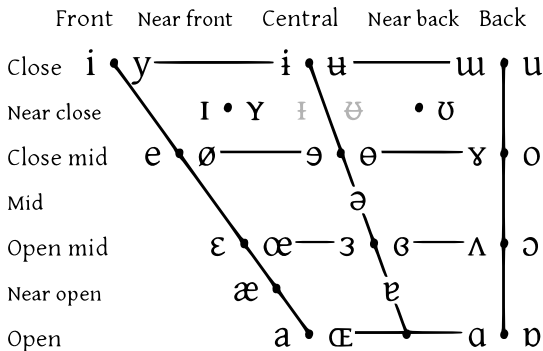
	LABIAL		CORONAL				DORSAL			RADICAL		LARYNGEAL
	Bilabial	Labio-dental	Dental	Alveolar	Palato-alveolar	Retroflex	Palatal	Velar	Uvular	Pharyngeal	Epi-glottal	Glottal
Nasal	m	ɱ	n			ɳ	ɲ	ŋ	ɴ			
Plosive	p b	ɸ β	t d			ʈ ɖ	c ɟ	k ɡ	q ɢ			
Fricative	ɸ β	f v	θ ð	s z	ʃ ʒ	ʂ ʐ	ç ʝ	x ɣ	χ ʁ	ħ ʕ	ħ ʕ	h ɦ
Approximant		ʋ	ɹ			ɻ	j	ɰ				
Trill	ʙ		r						ʀ		ʀ̤	
Tap, Flap		ɹ̥	ɾ			ɽ						
Lateral fricative			ɬ ɮ			ɮ̚	ɬ̺	ɮ̺				
Lateral approximant			l			ɭ	ʎ	ʟ				
Lateral flap			ɭ			ɭ̚						

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

IPA Chart: Vowels

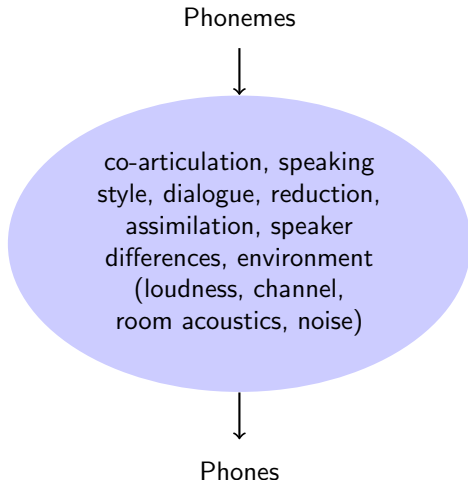
THE INTERNATIONAL PHONETIC ALPHABET (2005)

VOWELS

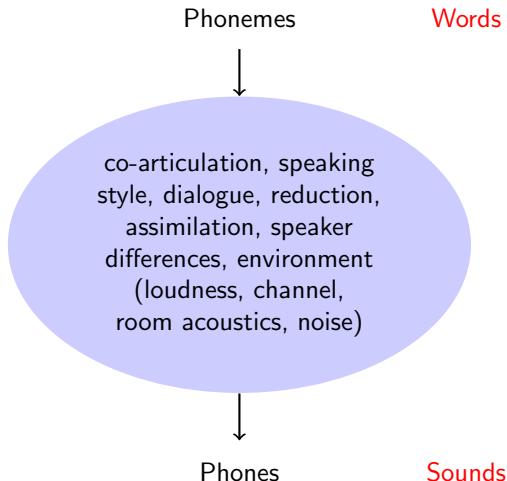


Vowels at right & left of bullets are rounded & unrounded.

Phonology vs Phonetics

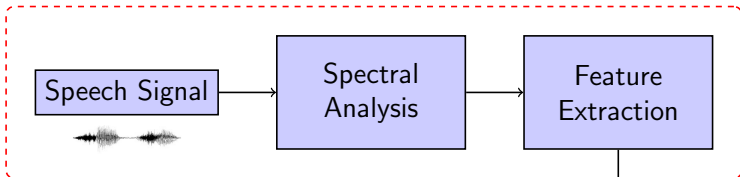


Phonology vs Phonetics

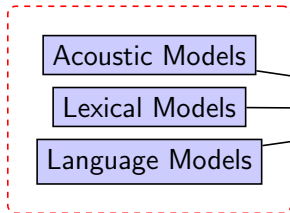


Components of ASR System

Representation



Constraints - Knowledge



Decoder

