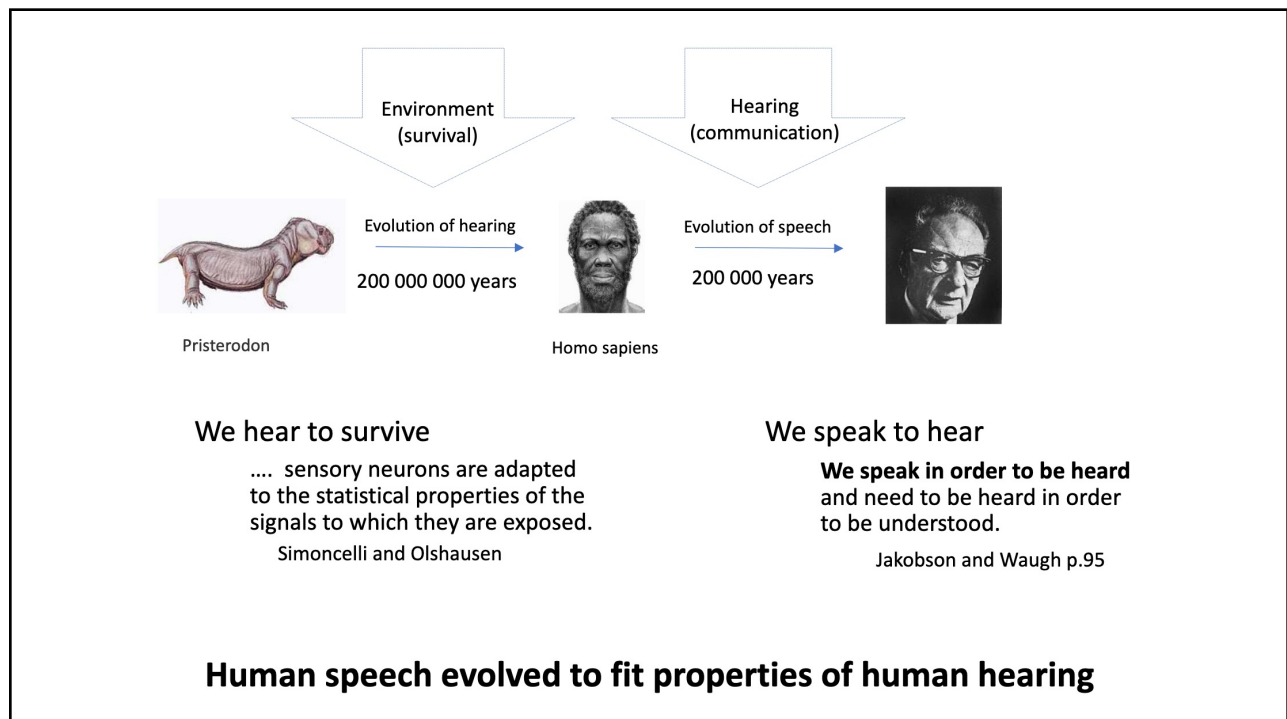


Speech basics

Instructor: Laureano Moro-Velazquez

Most slides from Hynek Hermansky

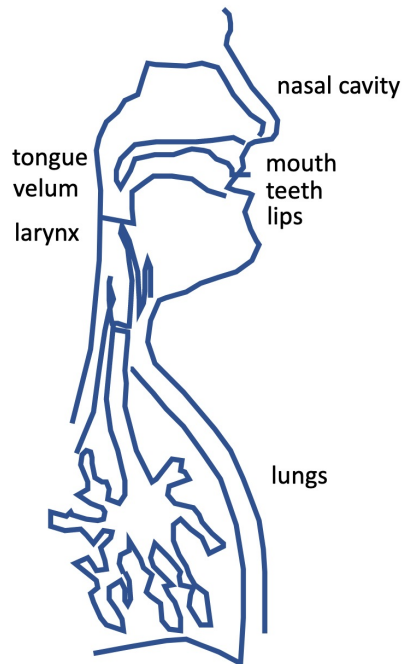
1



2

Speech generation

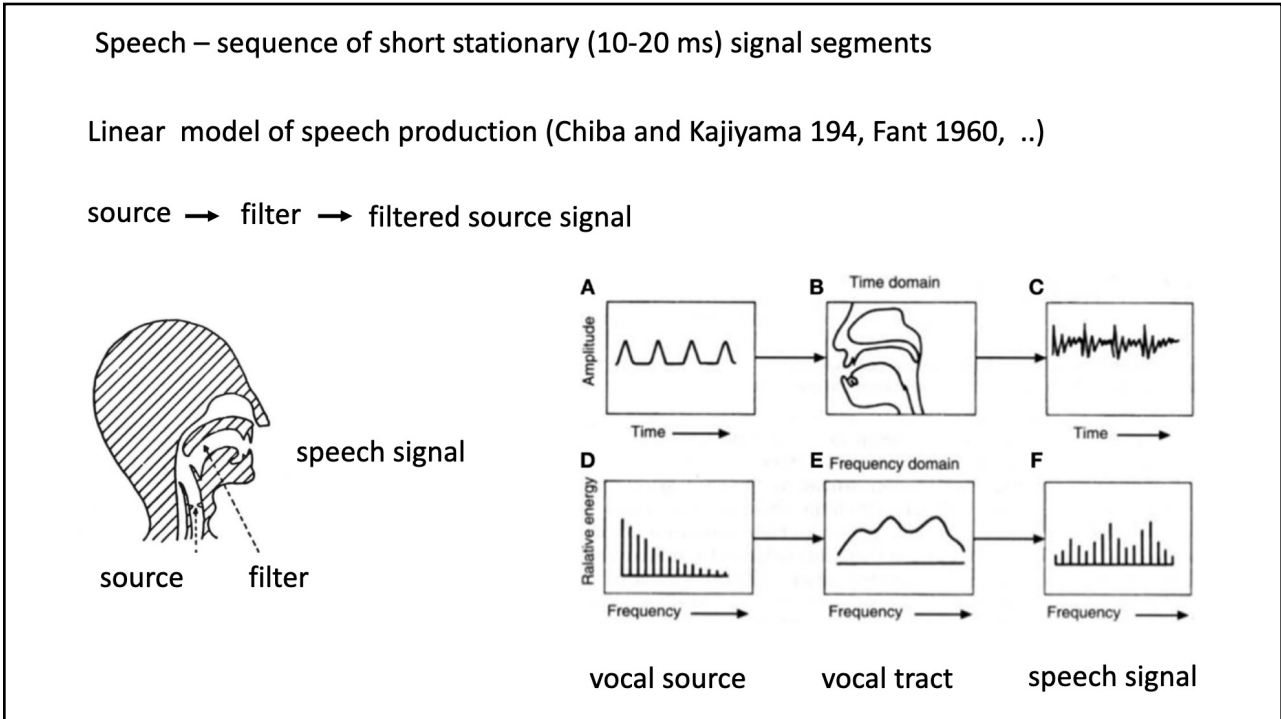
3



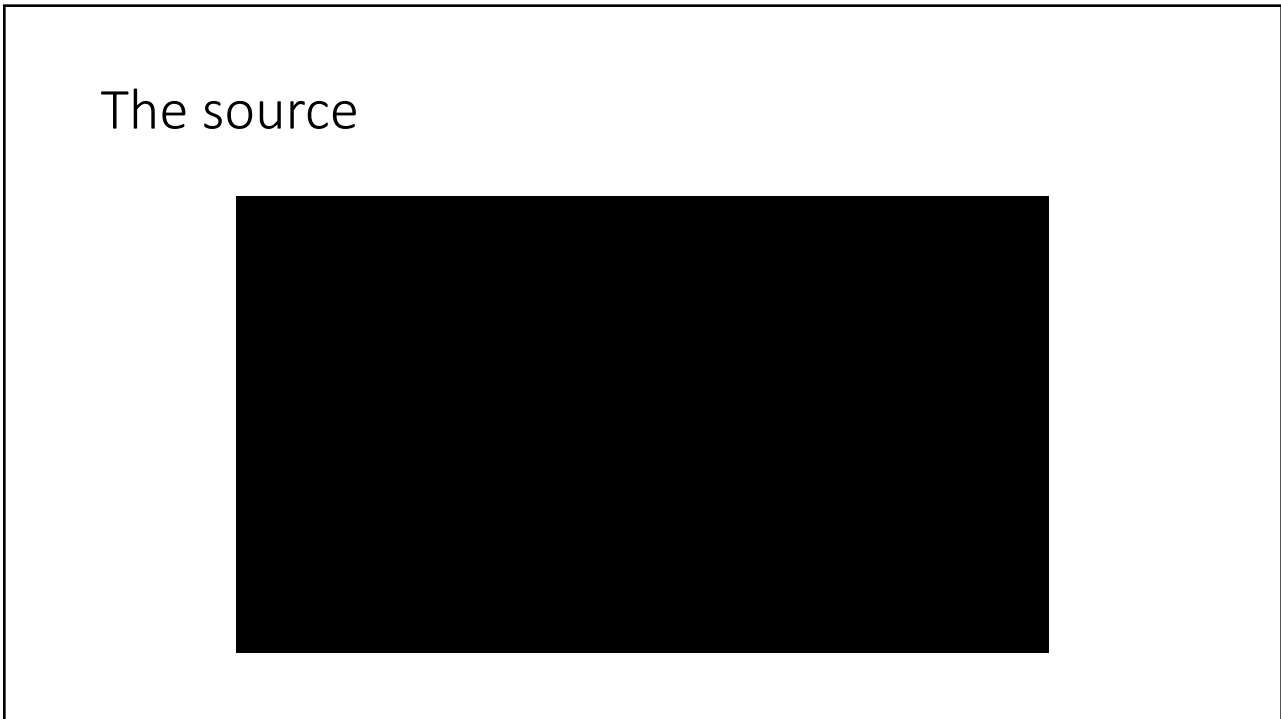
breathing
eating
biting

speaking?

4

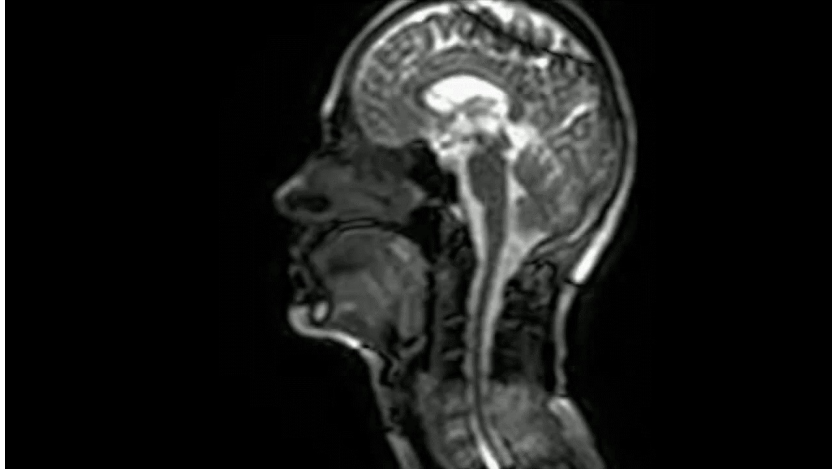


5



6

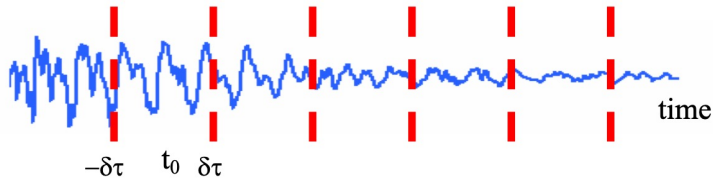
The vocal tract (the filter)



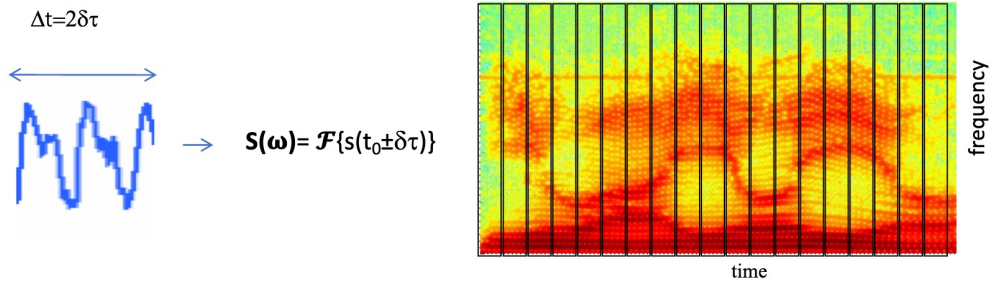
Video of rt-MRI of vocal tract during speech. (Freitas, A. C. et al, 2016)

7

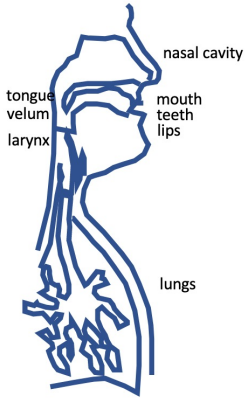
non-stationary speech signal $s(t)$



time-frequency representation of the signal (spectrogram)



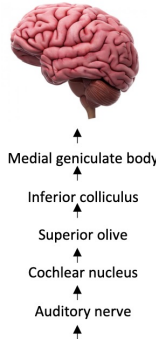
8



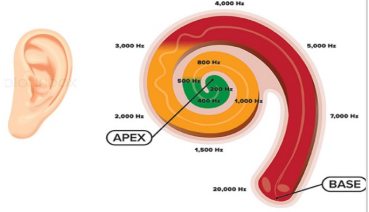
Redundant spread of information

- every change of the tract shape shows at all frequencies of speech spectrum
- tract shape changes do not happen very fast

brain



ear



- frequency selective (about 20 bands)
- sluggish (tenths of seconds)

9

INFORMATION ABOUT TRACT SHAPES DISTRIBUTED IN FREQUENCY

motor control

→

critical elements
(tongue, lips, velum)

→

shape of
the whole
vocal tract

→

→

→

→

→

→

spectrum of speech signal
(redundant contributions of movements of critical elements in different frequency bands)

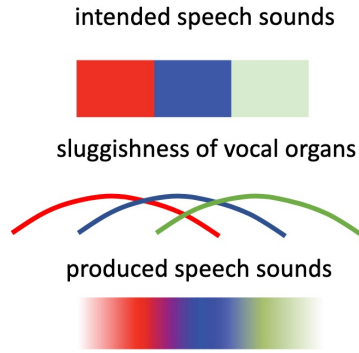
10

INFORMATION ABOUT TRACT SHAPES DISTRIBUTED IN TIME

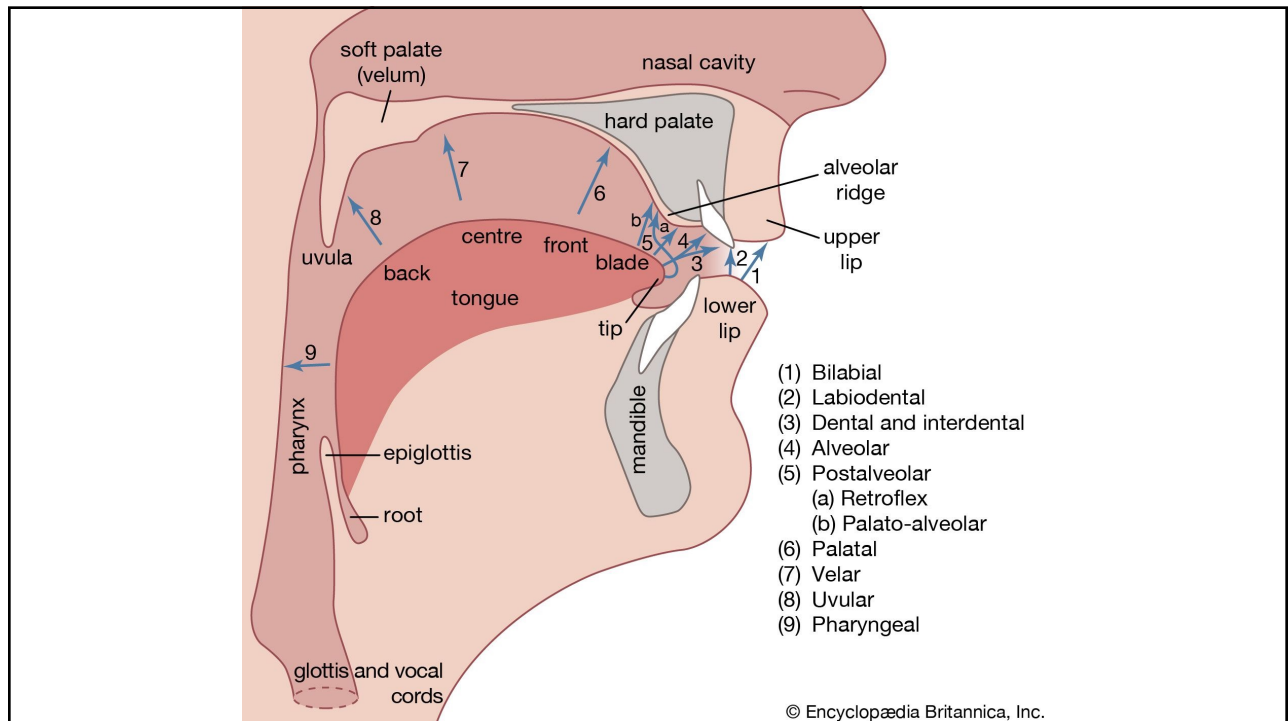


from Sri Narayanan

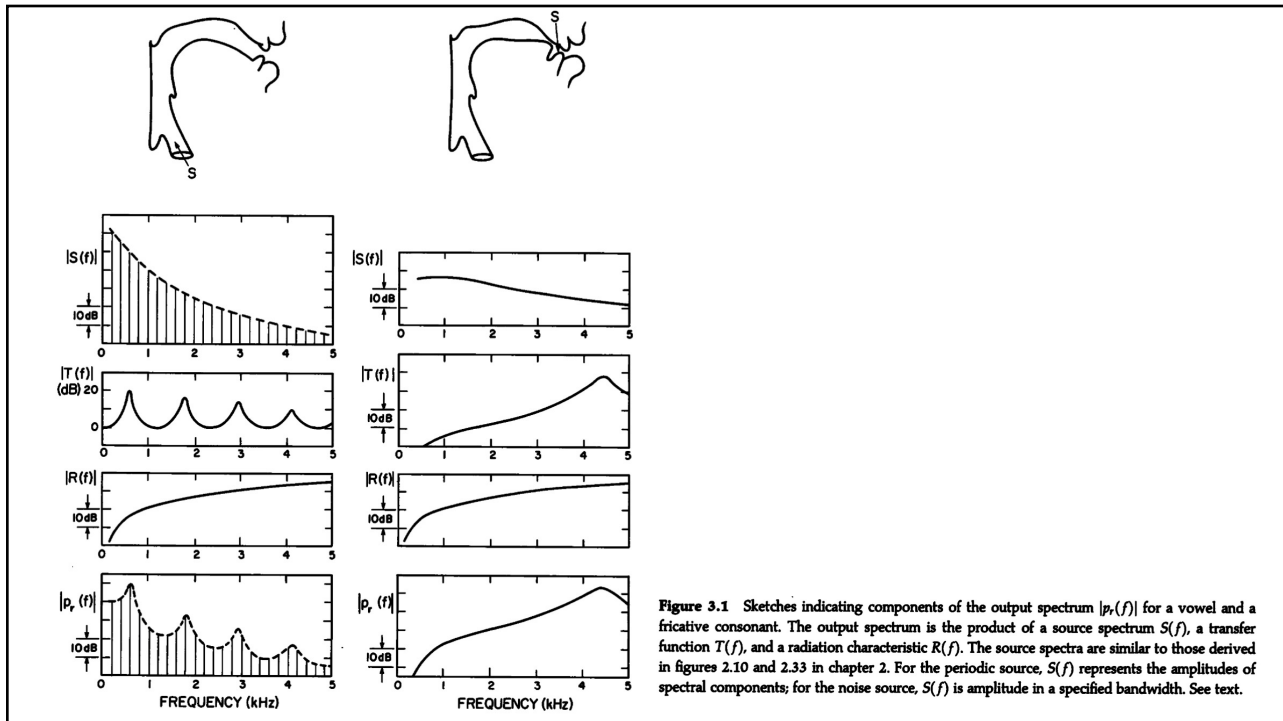
movements of vocal organs are rather sluggish



11



12



13

Articulation places THE INTERNATIONAL PHONETIC ALPHABET (revised to 2015)

CONSONANTS (PULMONIC)

© 2015 IPA

	Bilabial	Labiodental	Dental	Alveolar	Postalveolar	Retroflex	Palatal	Velar	Uvular	Pharyngeal	Glottal
Plosive	p b			t d		ʈ ɖ	c ɟ	k ɡ	q ɢ		ʔ
Nasal	m	ɱ		n		ɳ	ɲ	ŋ	ɴ		
Trill	ʙ			ʀ					ʁ		
Tap or Flap		ⱱ		ɾ		ɽ					
Fricative	ɸ β	f v	θ ð	s z	ʃ ʒ	ʂ ʐ	ç ʝ	x ɣ	χ ʁ	ħ ʕ	h ɦ
Lateral fricative				ɬ ɮ							
Approximant		ʋ		ɹ		ɻ	j	ɰ			
Lateral approximant				l		ɭ	ʎ	ʟ			

Symbols to the right in a cell are voiced, to the left are voiceless. Shaded areas denote articulations judged impossible.

Manner classes

14

Some phonetics

- **Articulation places:** points or areas in the vocal tract where there is a constriction (with or without contact) which has the most relevance in the generated sound.
- **Manner classes:** ways in which we articulate that produce consonant sounds.

15

Some phonetics: manner classes

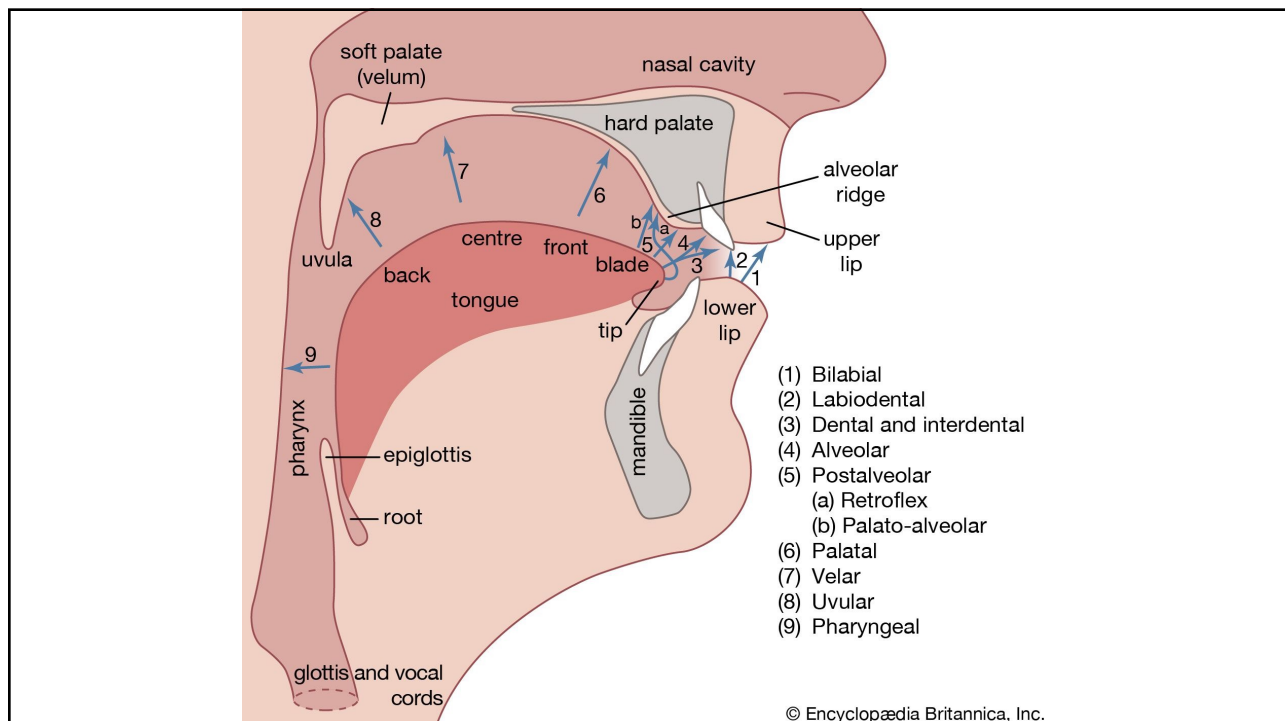
- **Plosive:** There is a constriction in the vocal tract and the airflow is interrupted for a short period of time (stop). Then, there is usually a release of the air that generates a sound.
- **Fricative:** The constriction of the vocal tract is very narrow but does not interrupt the airflow. This generates certain turbulences that lead to fricative sounds

16

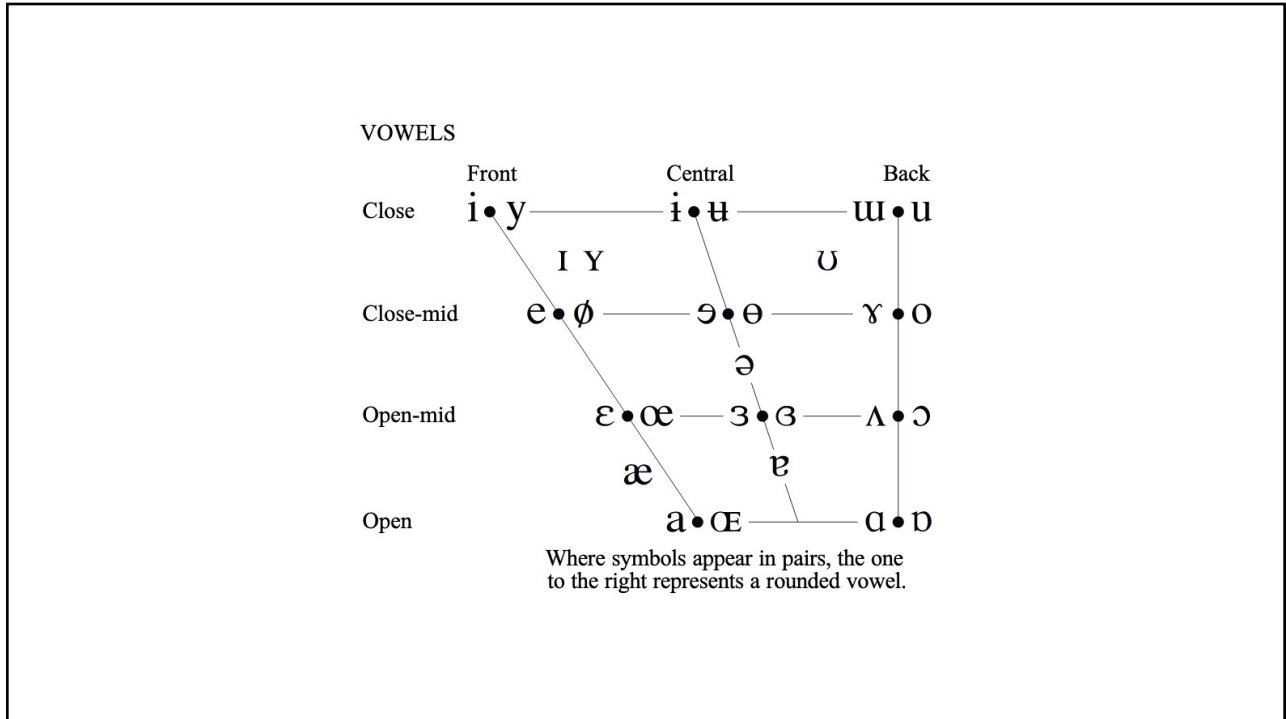
Some phonetics: manner classes

- Nasal: The soft palate is lowered, and the airflow goes through the nasal cavity. Nasal sounds are usually voiced (the source is on).
- Trill: The articulator (tongue, lips...) vibrate against other parts of the mouth with multiple flaps while the airstream is flowing (like in the word Ramon, in Spanish).
- Flap: Is similar to trill, but in this case there is only one short flap, like /r/ in the word radar in English).

17



18

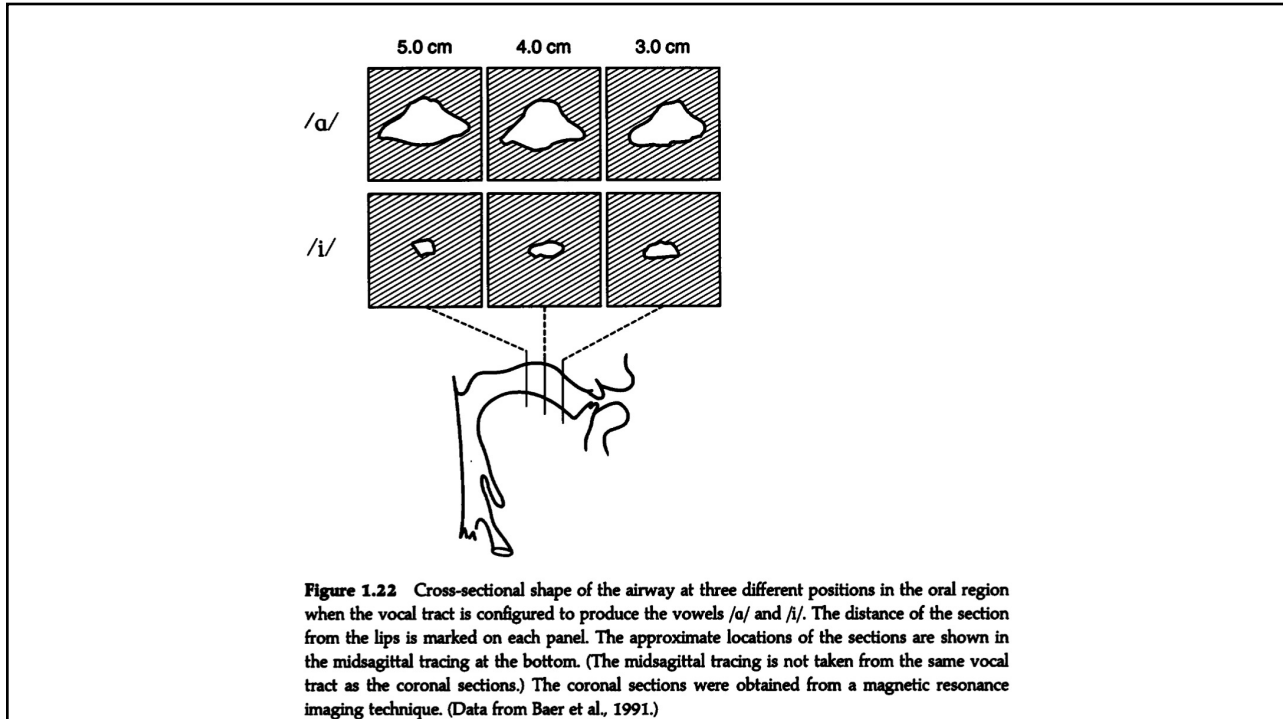


19

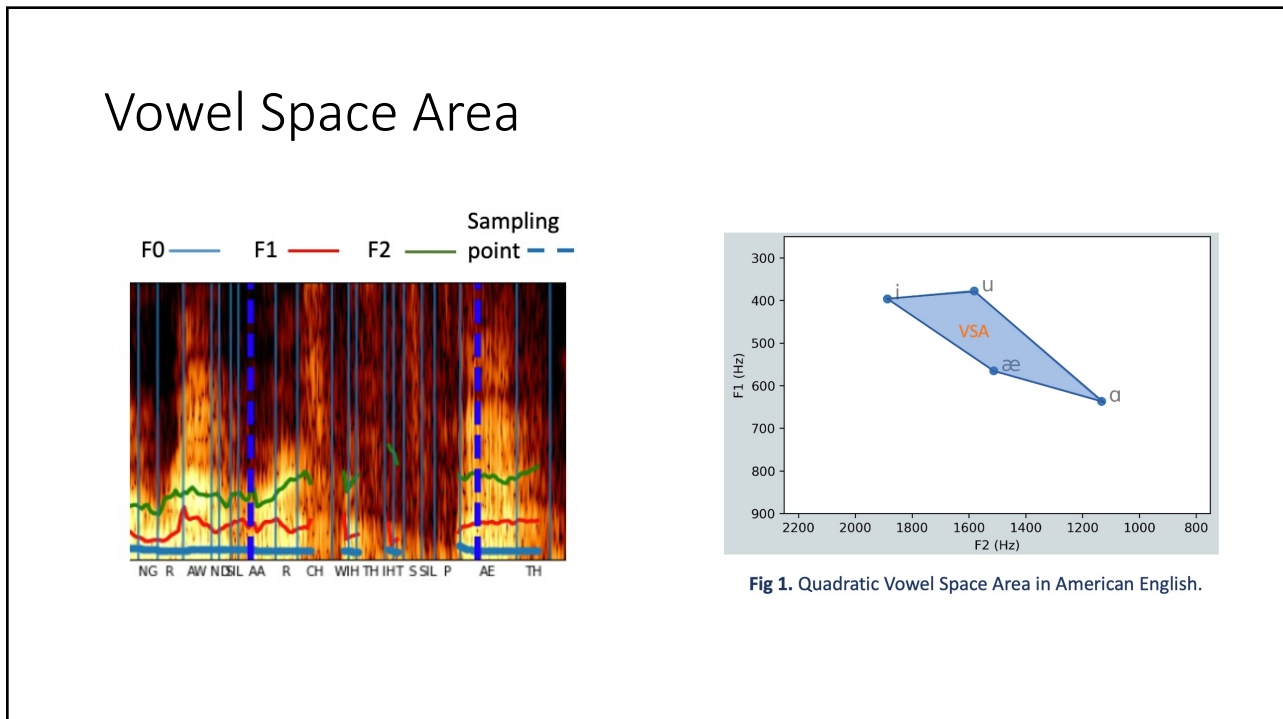
Formants

"heed" "hod" "who'd"

20



21

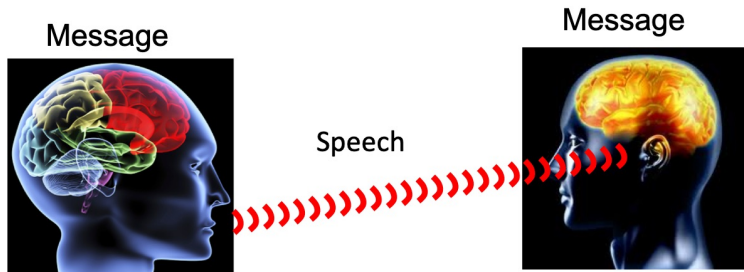


22

The message

23

Human Speech



24

Messages

- Only a limited number of speech sounds can be produced and distinguished
- Many things need to be said

Compositionality: meaning of a complex expression is determined by the meanings of its constituent expressions and the rules used to combine them (*Wikipedia*)

Create words as ordered sequences of speech sounds (phonemes).

file /fil/

life /lif/

Create phrases as ordered sequences of words.

Tom chased horse.

Horse chased Tom.

25

Prior probabilities of different letters in English alphabet

Letter	Relative frequency	Letter	Relative frequency
e	12.702%	m	2.406%
t	9.056%	w	2.360%
a	8.167%	f	2.228%
o	7.507%	g	2.015%
i	6.966%	y	1.974%
n	6.749%	p	1.929%
s	6.327%	b	1.492%
h	6.094%	v	0.978%
r	5.987%	k	0.772%
d	4.253%	j	0.153%
l	4.025%	x	0.150%
c	2.782%	q	0.095%
u	2.758%	z	0.074%



Samuel Morse
(self-portrait)

Morse code

e - single dot

z - dot and three dashes

26

In 1939, Ernest Vincent Wright published a 267-page novel, Gadsby, in which **no use is made of the letter E**. Here is a paragraph from the novel:

Upon this basis I am going to show you how a bunch of bright young folks did find a champion; a man with boys and girls of his own; a man of so dominating and happy individuality that Youth is drawn to him as is a fly to a sugar bowl. It is a story about a small town. It is not a gossip yarn; nor is it a dry, monotonous account, full of such customary "fill-ins" as "romantic moonlight casting murky shadows down a long, winding country road." Nor will it say anything about tinklings lulling distant folds; robins carolling at twilight, nor any "warm glow of lamplight" from a cabin window. No. It is an account of up-and-doing activity; a vivid portrayal of Youth as it is today; and a practical discarding of that worn- out notion that "a child don't know anything."

27

How "efficient" is a given code?

all letters are equally probable (zero order)

Entropy

$H(s) = 4.74$ bit

$$H(s) = - \sum_{i=1}^n p_i \cdot \log(p_i)$$

Respecting relative frequencies of letters (first order)

$H(s) = 4.28$ bit

26 letters and space of English alphabet

$$H(s) = - \sum_{i=1}^{27} 1/27 \cdot \log(1/27) \\ = -\log(1/27) = 4.74 \text{ bit}$$

Respecting relative frequencies of combinations of three letters (third order)
 $H(s) = 2.77$ bit

Letters in real text (estimate)
 $H(s) \sim 0.6-1.3$ bit

Shannon
Prediction and Entropy of Printed English
BSTJ 1951

28

	Phoneme	Frequency Percentage				
		<i>per cent</i>				
The Relative Frequency of Phonemes in General-American English Hayden 1950	ə	9.96	n	7.95	f	1.61
	ɪ	9.75	t	7.59	y	1.20
	æ	3.09	r	7.10	g	1.14
	ɛ	2.03	s	4.89	h	1.11
	e	1.94	l	3.65	ʃ	0.87
	a	1.80	ʔ	3.35	ŋ	0.80
	i	1.66	d	3.21	č	0.53
	u	1.52	k	2.98	j	0.50
	o	1.49	m	2.87	θ	0.44
	a ⁱ	1.46	z	2.36	w	0.37
	ɔ	1.02	v	2.33	ʒ	0.03
	ʊ	0.99	p	2.25		
	a ^u	0.64	w	1.77		
	o ⁱ	0.06	b	1.65		
			37.4			62.6

29

Phonemes

Perceptually distinct speech elements that could distinguish one words from another

Graphemes

Letters and combinations of letters representing speech sounds (phonemes)

Rotokas language – East of New Guinea, 11 phonemes, 12 symbols, 1 symbol per sound

Taa language – Botswana (Africa), ~ 200 phonemes , 20-22 symbols, up to 6 symbols per sound

English

~45 phonemes, 27 symbols,

~ 250 graphemes, up to 5 symbols per sound

30

40 speech sounds (phonemes) in American English
 24 consonants
 19 vowels and diphthongs

vowels – mouth open
 consonants - mouth not so open

typical syllable

cvc
 onset – nucleus – coda
 cv
 onset – nucleus

/l/,/r/,/w/,/y/ - semivowels
 produced with open mouth
 can stand as nucleus in syllable

31

Phones, phonemes and allophones

- **Allophones:** different realizations of a phone, depending on the dialect or other domain changes. These do not change the word meaning when they change.
- **Phone:** we usually call phone to a specific segment that contains a distinct sound, but it does not have to be critical for to the meaning of a word. A phone can be a phoneme or part of it.

32

Phones, phonemes and allophones

- If in a word you change a phoneme, you will change the meaning of a word. If you change a phone, you might not change the meaning of that word.
- The phoneme is the mental realization, the phone is the sound representation of a phoneme.

33

Words

- ordered combinations of speech sounds
- represent objects, ideas, actions, relationships, qualities, e.t.c., **as agreed on by a particular society (language)**
- new words constantly invented and old words changing their meanings
- learned using interventions and rewards from other human beings
- particular word meanings often depend on context

34

Word sequences (sentences, phrases,..)

- Words organized into larger units (sentences, phrases,..) using rules of the language (syntax, grammar)
- Order also carries information
 - John beats Frank. Frank beats John.
 - I went home and had a dinner. I had a dinner and went home.

35

Relative frequencies of words in written English [%]

7.31	the	.58	not	.31	their	.20	time	.15	these
3.99	of	.58	at	.30	there	.20	up	.14	two
3.28	and	.57	this	.30	were	.20	do	.14	very
2.92	to	.54	are	.30	so	.20	out	.13	before
2.12	a	.52	we	.29	my	.19	can	.13	great
2.11	in	.51	his	.26	if	.19	than	.13	could
1.34	that	.50	but	.25	me	.18	only	.13	such
1.21	it	.47	they	.25	what	.18	she	.13	first
1.21	is	.46	all	.25	would	.17	made	.12	upon
1.15	I	.45	or	.24	who	.16	other	.12	every
1.03	for	.45	which	.23	when	.16	into	.12	how
.84	be	.44	will	.23	him	.16	men	.12	come
.83	was	.43	from	.22	them	.16	must	.12	us
.78	as	.41	had	.22	her	.16	people	.12	shall
.77	you	.39	has	.21	war	.16	said	.11	should
.72	with	.36	one	.21	your	.16	may	.11	then
.68	he	.33	our	.21	any	.15	man	.11	like
.64	on	.33	an	.21	more	.15	about	.11	well
.61	have	.32	been	.21	now	.15	over	.11	little
.60	by	.32	no	.20	its	.15	some	.11	say

In spoken language most frequency word is pronoun "I"
 Telephone conversations 5%
 Schizophrenics 8.4%

36

Predictability and unpredictability

- 100 % predictable message has no information value
 - When knowing exactly what will be said, no need to listen
- Speech is to large extent predictable since it follows rules
 - Grammar, use of words, word order, ...
- The predictability allows for easier communication

To communicate effectively, the right balance between predictability and unpredictability need to be maintained.

37

How predictable is language? - Claude Shannon

1. Think about the English sentence
2. Ask people to think about the first letter in the sentence
3. When correct, tell them, mark it by “-” and ask for the second letter
4. When incorrect, tell them the correct one and ask for the second letter
5. Go on until the end of the sentence

(1) THE ROOM WAS NOT VERY LIGHT A SMALL OBLONG

(2) ----R00-----NOT-V-----I-----SM----OBL-----

(1) READING LAMP ON THE DESK SHED GLOW ON

(2) REA-----0-----D----SHED-GLO--0--

(1) POLISHED WOOD BUT LESS ON THE SHABBY RED CARPET

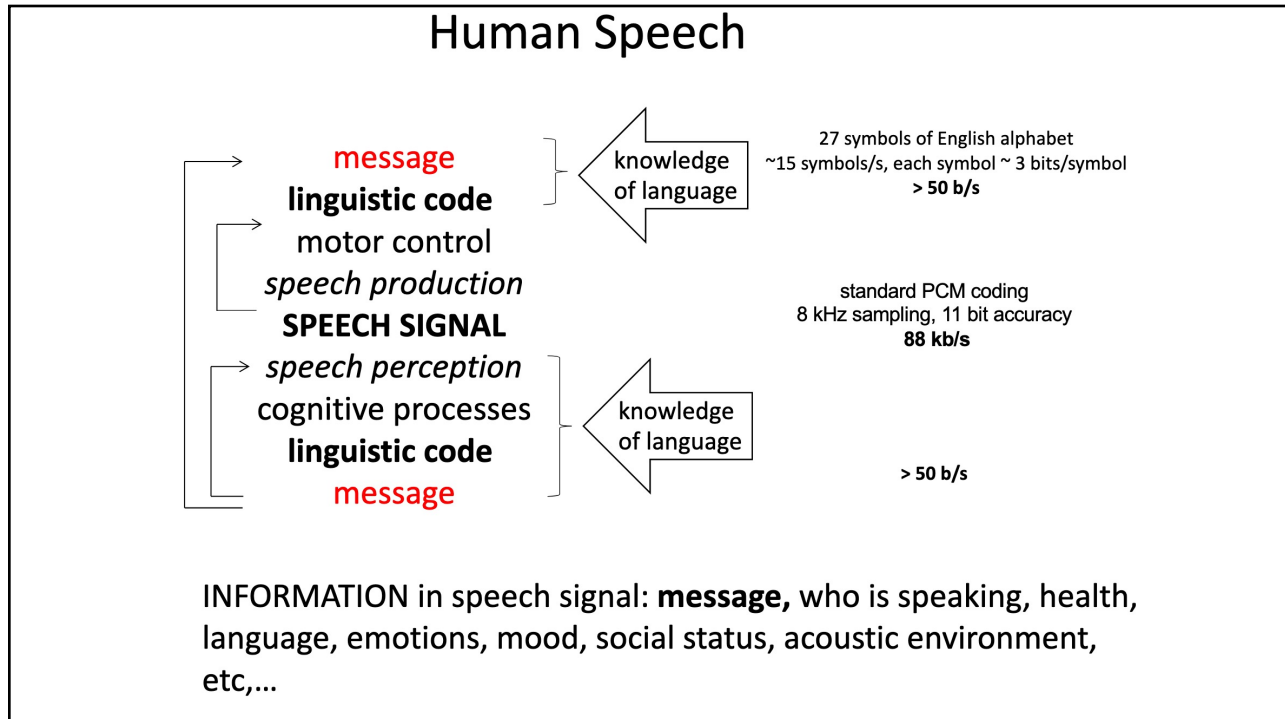
(2) P-L-S-----0---BU--L-S--0-----SH-----RE--C-----

69% of letters guessed correctly

Both line (1) and (2) contain the same information

- The line (1) can be guessed from the info in the line (2) – by the identical twin 😊

38



39

Why speech?

- **Profit**
 - searching large speech databases, transcription, voice control,...
 - ***voice will do to touch what touch did to keyboards.***
 - Mooly Eden, senior vice president Intel
- **Important spin-offs**
 - Digital signal processing
 - Sequence classification (Hidden Markov Models)
 - financial predictions
 - human DNA matching
 - action recognition
 - Image processing techniques

Spoken language is one of the most amazing accomplishments of human race.

40

Human Language Technologies

A brief look

43



44

Are We There Yet ?

- Repetition, fillers, hesitations, interruptions, unfinished and non-grammatical sentences, new words, dialects, emotions, ...
- Hands-free operation in noisy and reverberant environments,...

Alleviate need for large amounts of annotated training data

- Robustness to speech distortions, which do not seriously impact human speech communication
- Dealing with new unexpected lexical items
- Unsupervised learning/adaptation?

45

How to Get There ?

Fred Jelinek



Speech recognition
...a problem of maximum likelihood decoding
information and communication theory, machine learning, large data,....

Roman Jakobson



We speak, in order to be heard, in order to be understood
human communication, speech production, perception, neuroscience, cognitive science,..

Gordon Moore



The complexity for minimum component costs has increased at a rate of roughly a factor of two per year...

John Pierce



..devise a clear, simple, definitive experiments. So a science of speech can grow, certain step by certain step.

Signal processing,
information theory,
machine learning, ...

&

neural information processing,
psychophysics, physiology,
cognitive science, phonetics and
linguistics, ...

Engineering and Life Sciences together !

46