

# Speech production and modeling

Simon Leglaive

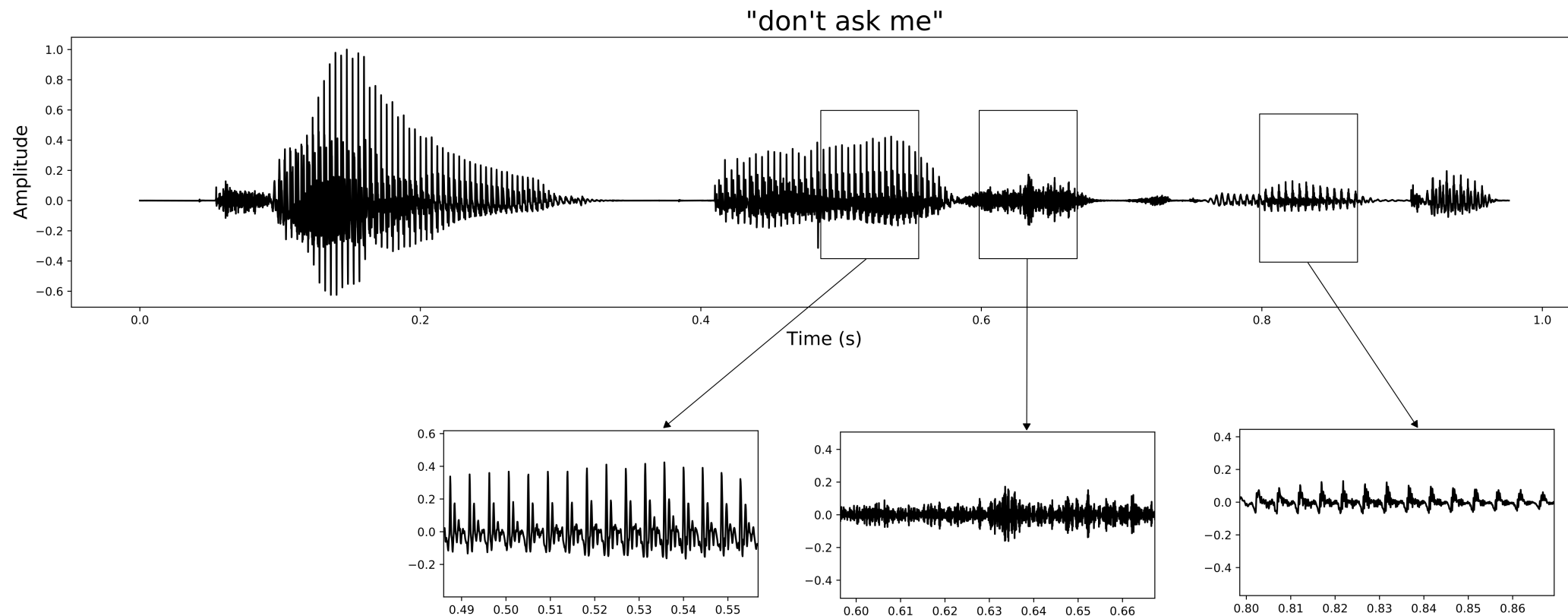
CentraleSupélec

# Today

- Speech production mechanisms
- Characteristics of speech signals

# Speech production

# Speech signal



Can you guess to what "speech sound" each bloc corresponds?

# Phonemes

Elementary speech sounds are called phonemes.

- 44 phonemes in English.
- 10-15 phonemes per second in normal English speech.
- We are going to see what are the key differences in the production of the different phonemes.

		monophthongs				diphthongs			
VOWELS	i:	ɪ	ʊ	u:	ɪə	eɪ			
	sheep	ship	good	shoot	here	wait			
	e	ə	ɜ:	ɔ:	ʊə	ɔɪ	əʊ		
	bed	teacher	bird	door	tourist	boy	show		
	æ	ʌ	ɑ:	ɒ	eə	aɪ	aʊ		
	cat	up	far	on	hair	my	cow		
CONSONANTS	p	b	t	d	tʃ	dʒ	k	g	
	pea	boat	tea	dog	cheese	June	car	go	
	f	v	θ	ð	s	z	ʃ	ʒ	
	fly	video	think	this	see	zoo	shall	television	
	m	n	ŋ	h	l	r	w	j	
	man	now	sing	hat	love	red	wet	yes	

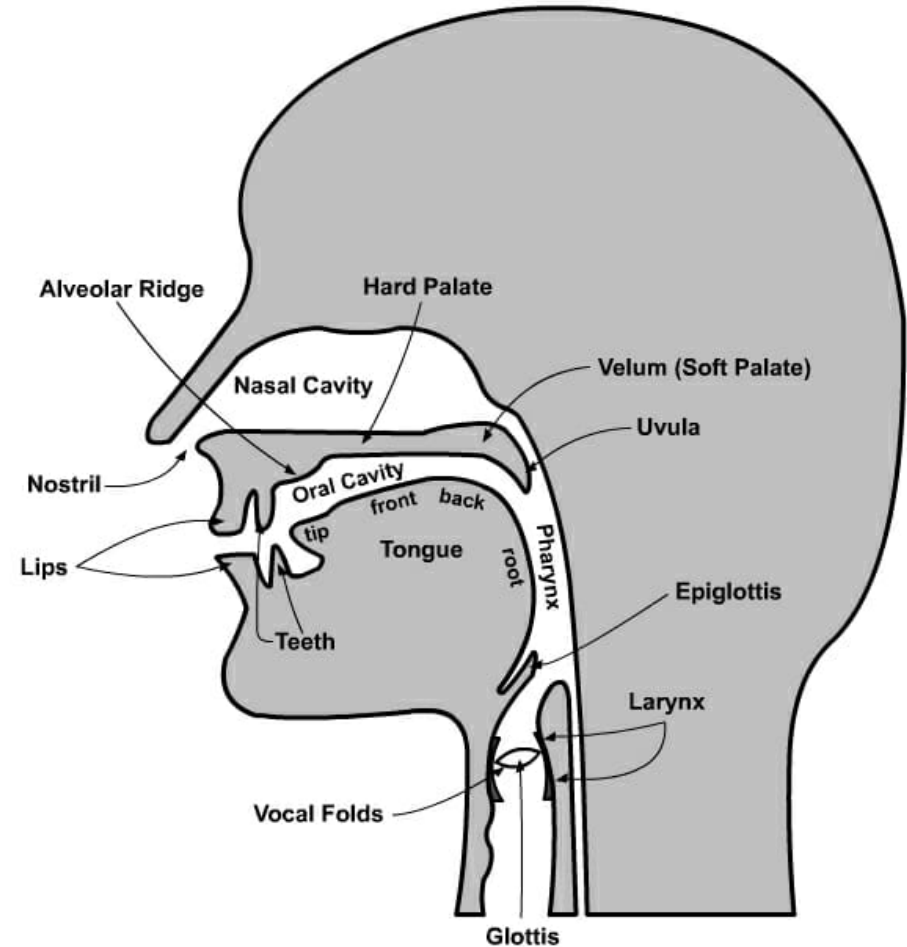
**Phonemic Chart**  
voiced  
unvoiced

The 44 phonemes of Received Pronunciation based on the popular Adrian Underhill layout

adapted by EnglishClub.com

# Speech production – the global view

- The energy comes from air expelled from the **lungs**.
- At the **larynx**, this airflow passes between the **vocal folds**.
- Then it goes through the **vocal tract**, which is made of **three cavities**:
  1. the pharynx
  2. the oral cavity
  3. the nasal cavity
- Finally, sound goes out of the mouth and nose openings.



# Articulators

We consider as articulator any **mobile part of the vocal tract on which we can act voluntarily** and which is **used in the production of speech sounds**.



*Real-time MRI scan of a person talking.*

What are the three main speech articulators?

## Tongue

- Very mobile and flexible
- Very important for phonation

## Jaw

- Little degrees of freedoms and rigid
- Less important for phonation

## Lips

- Very mobile and flexible
- Important movements for phonation:
  - occlusion
  - protrusion
  - raising and lowering
  - stretching, raising and lowering of lip corners



# Speech sound sources

We distinguish 3 types of **sound sources**, which can be **combined** or **occur individually**:

- **Quasi-periodic source** resulting from the vibration of the **vocal folds**.

We say that the sound is **voiced**.

It can be **arbitrarily long** (in the limits of an exhalation).

- **Fricative noise source** produced by a **turbulent airflow** with a **constriction** in the vocal tract.

It can also be **arbitrarily long**.

- **Plosive noise source** produced by quick **occlusions** of the vocal tract and generating an **acoustic impulse**.

Here the **duration is short**.

# Voice production

PhysclipsWS > Human sound > 9.1 Voice Production

Voice production

## PHYSCLIPS Waves and Sound

### Introduction

1. Oscillations
2. Travelling Waves I
3. Travelling waves II
4. Sound
5. The Doppler Effect
6. Quantifying sound
7. Interference, consonance
8. Standing waves
9. Human sound

- 9.1 Voice production
- 9.2 Pitch & mechanisms
- 9.3 Resonances & formants
- 9.4 Outer ear
- 9.5 Middle ear
- 9.6 Inner ear
- 9.7 Pitch perception

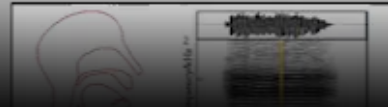
Respiration & the voice



Turbulence, sibilants & unvoiced speech



Voiced speech



## Vocal folds and pitch

- The vibration of the vocal folds defines the **pitch** of the speech signal (i.e. its fundamental frequency).
- Variations of pitch along time define the melody of the voice.

	Average pitch (Hz)	Pitch range (Hz)
Male	100 - 130	90 - 270
Female	150 - 300	120 - 360
Child	350 - 400	200 - 600

# Pitch and mechanisms

PhysclipsWS > Human sound > 9.2 Pitch and mechanisms

Pitch and mechanisms

## PHYSCLIPS Waves and Sound

### Introduction

1. Oscillations
2. Travelling Waves I
3. Travelling waves II
4. Sound
5. The Doppler Effect
6. Quantifying sound
7. Interference, consonance
8. Standing waves
9. Human sound
  - 9.1 Voice production
  - 9.2 Pitch & mechanisms
  - 9.3 Resonances & formants
  - 9.4 Outer ear
  - 9.5 Middle ear
  - 9.6 Inner ear
  - 9.7 Pitch perception

Mechanisms



tension & pressure



pitch & language

妈妈骂马

pitch frequency

$\lambda \gg L$



## Vocal tract and formants

- The three elementary sound sources are **modified by the vocal tract**, before propagating out of the phonatory system, through the mouth and nose openings.
- The vocal tract actually corresponds to an **acoustic filtering** of the source signal.
- The cavities in the vocal tract give rise to **resonances**, that are called the **formants**.
- By **modifying the shape** of the vocal tract, we change the acoustic filter and the associated resonances.
- We can **change the formants independently of the pitch**, or in signal processing terms, we can change the filter independently of the source

# Resonances and formants

PhysclipsWS > Human sound > 9.3 Resonances and formants

## Resonances and formants

### PHYSCLIPS Waves and Sound

#### Introduction

1. Oscillations
2. Travelling Waves I
3. Travelling waves II
4. Sound
5. The Doppler Effect
6. Quantifying sound
7. Interference, consonance
8. Standing waves
9. Human sound

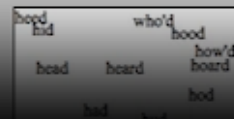
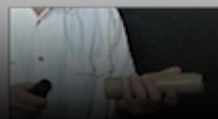
- 9.1 Voice production
- 9.2 Pitch & mechanisms
- **9.3 Resonances & formants**
- 9.4 Outer ear
- 9.5 Middle ear
- 9.6 Inner ear
- 9.7 Pitch perception

tract resonances

resonances & formants

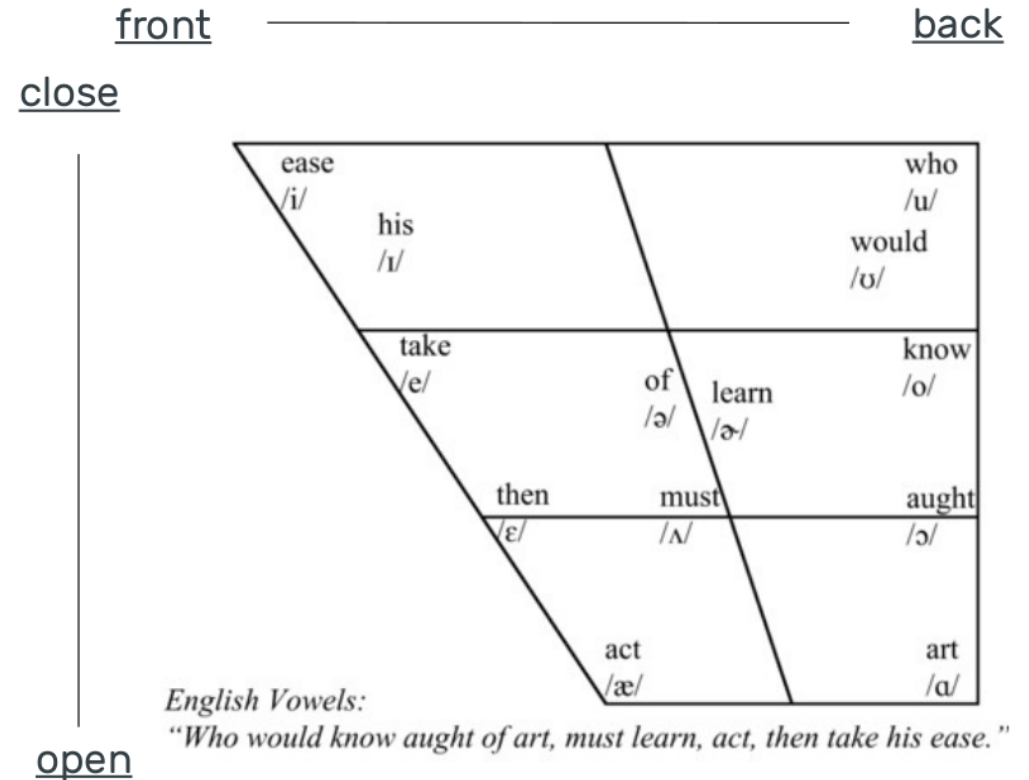
Vowel plane

Consonants



# Distinctive articulatory features of vowels

- **Opening** of the mouth
  - Opened vowel [a] in “hat”
  - Closed vowel [i] in “meet”
- **“Frontness”** of the tongue
  - Front vowel [i] in “meet”
  - Back vowel [u] in “boot”
- **Rounding** of the lips
  - Rounded vowel [ɔ] in “not”
  - Not rounded vowel [i] in “meet”
- **Nasalization**: sound comes out of the mouth only, or out of the mouth and nose.

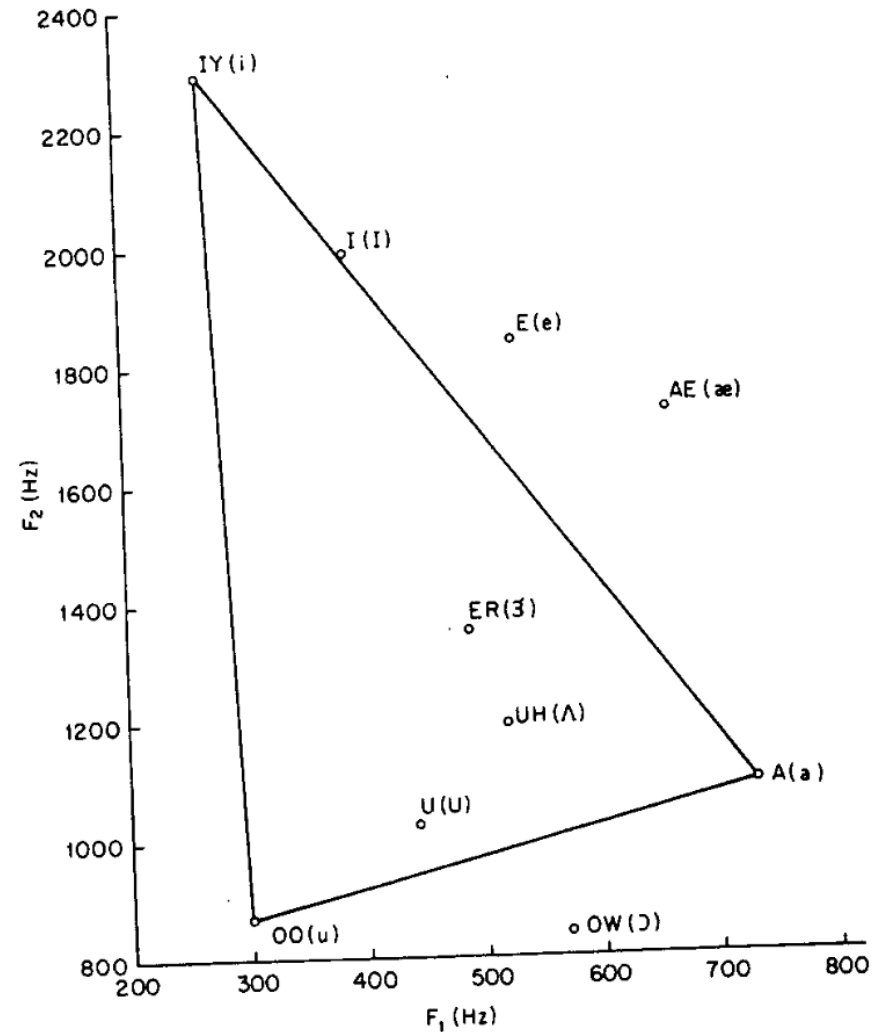


# Vowels and formants

We can distinguish between vowels using the **position of the first formants**

- high/low F1 ↔ opened/closed
- high/low F2 ↔ front/back
- high/low F3 ↔ not rounded/rounded lips

By moving articulators, the shape of the vocal tract varies, formants move in frequency, and vowels change.



**Fig. 3.5** The vowel triangle.

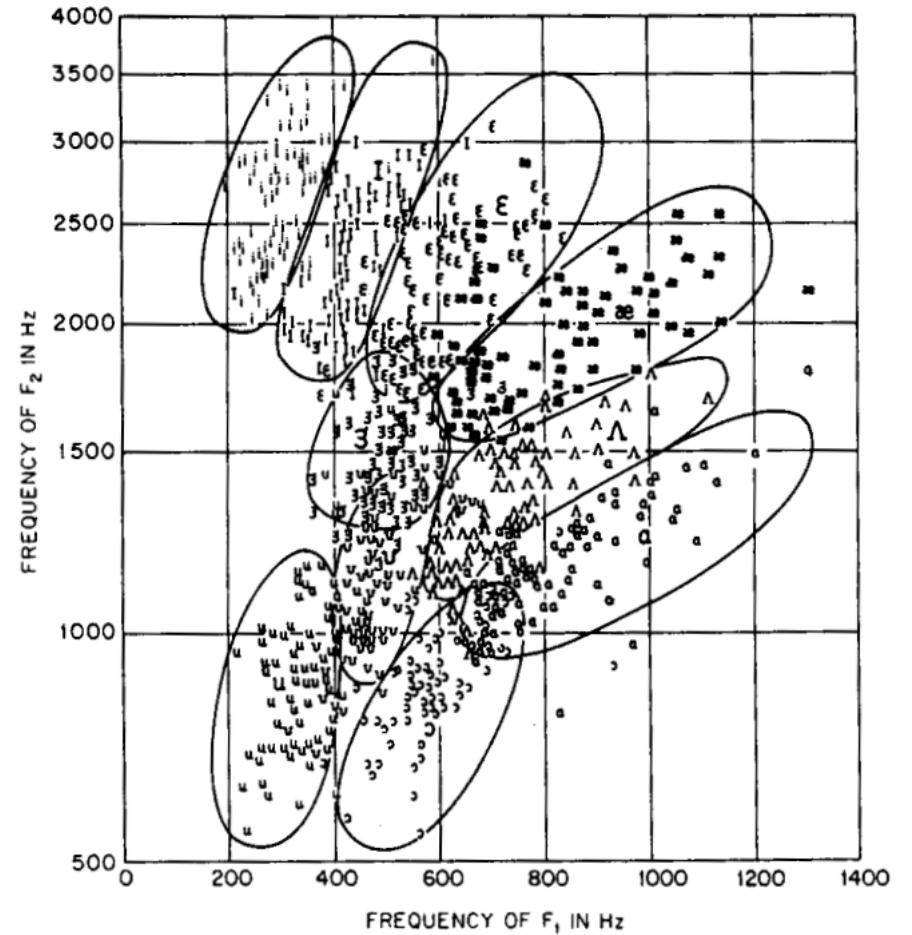


# Vowels clustering in the formants space

**Table 3.2** Average Formant Frequencies for the Vowels. (After Peterson and Barney [11].)

FORMANT FREQUENCIES FOR THE VOWELS					
Typewritten Symbol for Vowel	IPA Symbol	Typical Word	F <sub>1</sub>	F <sub>2</sub>	F <sub>3</sub>
IY	i	(beet)	270	2290	3010
I	ɪ	(bit)	390	1990	2550
E	ɛ	(bet)	530	1840	2480
AE	æ	(bat)	660	1720	2410
UH	ʌ	(but)	520	1190	2390
A	ɑ	(hot)	730	1090	2440
OW	ɔ	(bought)	570	840	2410
U	u	(foot)	440	1020	2240
OO	ʊ	(boot)	300	870	2240
ER	ɜ	(bird)	490	1350	1690

*male speakers*



**Fig. 3.4** Plot of second formant frequency versus first formant frequency for vowels by a wide range of speakers. (After Peterson and Barney [11].)

*male and children speakers*

# Consonants

## Fricatives

- fricative noise source
- voiced [v, z, ʒ] or unvoiced [f, θ, s]
- locally stationary

## Plosives

- plosive noise source
- voiced [b, d, g] or unvoiced [p, t, k]
- highly non-stationary

## Nasal

- voiced
- sound comes mostly from the nose
- examples: [m, n]

## Liquids

- voiced
- the vocal tract changes rapidly, especially using the tongue
- examples: [l, r]

# Consonants

## Fricatives

- fricative noise source
- voiced [v, z, j] or unvoiced [f, s, ch]
- locally stationary

## Plosives

- plosive noise source
- voiced [b, d, g] or unvoiced [p, t, k]
- highly non-stationary

## Nasal

- voiced
- sound comes mostly from the nose
- examples: [m, n]

## Liquids

- voiced
- the vocal tract changes rapidly, especially using the tongue
- examples: [l, r]

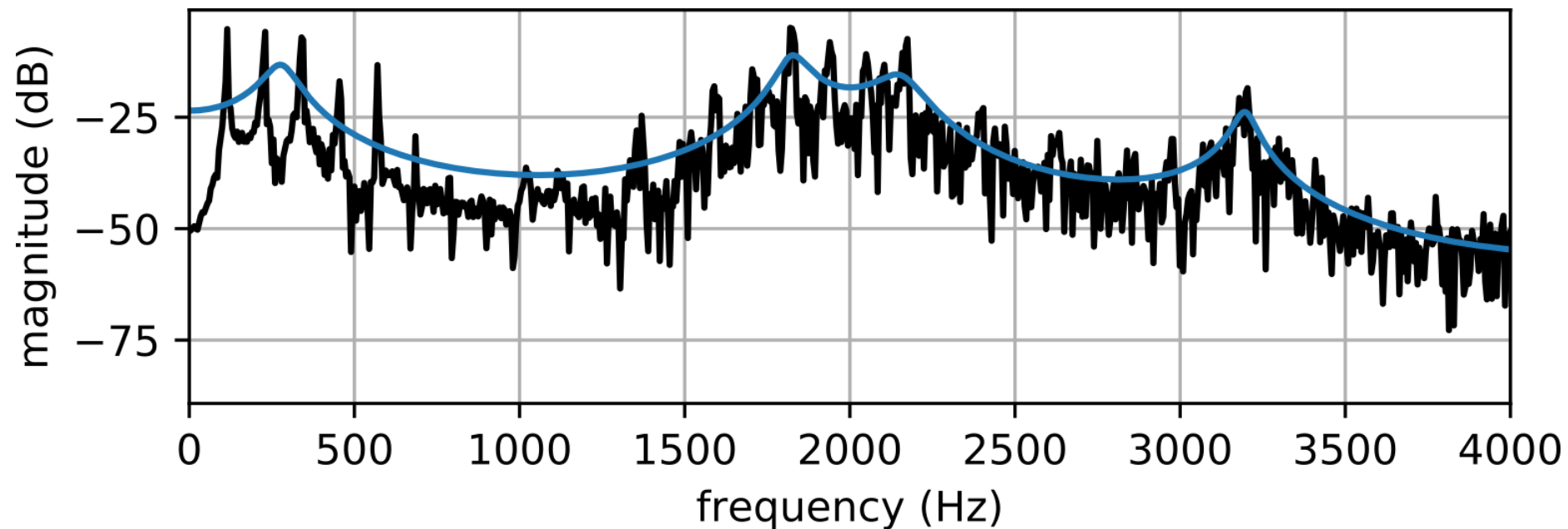
Go to <https://app.wooclap.com/CXIOJL>

# Prosody

- **Prosody is on top of the flow of phonemes.**
- Prosodic variables:
  - **pitch** (fundamental frequency)
  - **speech rate** (number of speech units, e.g. phonemes, per second)
  - **loudness** (intensity)
  - **timbre** (spectral characteristics such as amplitude of harmonics)
- Different combinations of these variables are exploited for intonation and accentuation.
- Prosody may reflect various features of the speaker or the utterance:
  - the identity of the speaker
  - the emotional state of the speaker
  - the form of the utterance (statement, question, or command)

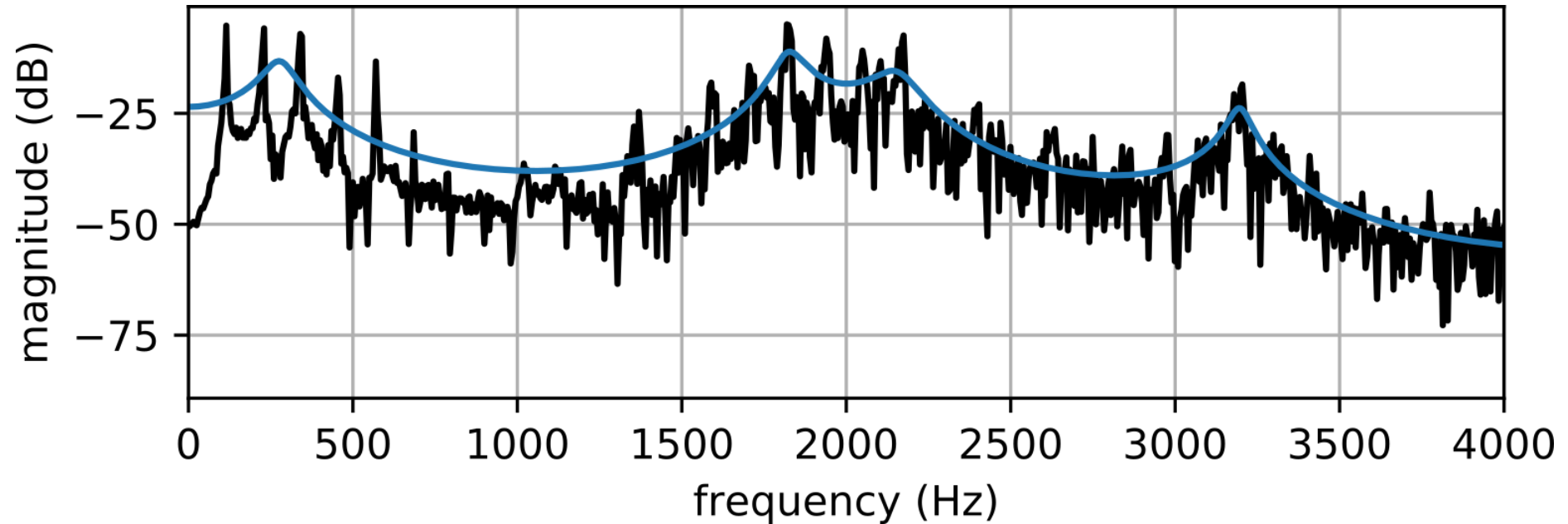
# Spectrum/spectrogram reading

## The spectral envelope



- Black curve: power spectrum (in dB) of the recording of a vowel, computed with the DFT.
- Blue curve: **spectral envelope** showing the formant resonances, computed with linear predictive coding (will be discussed in the lab session).

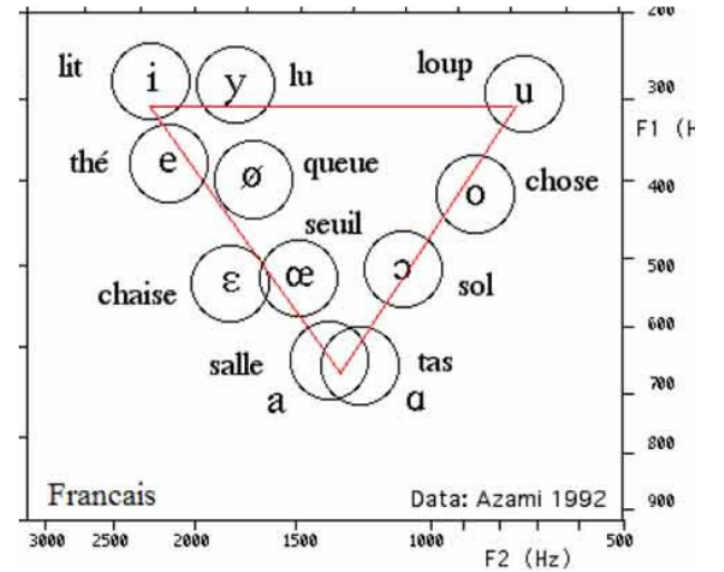
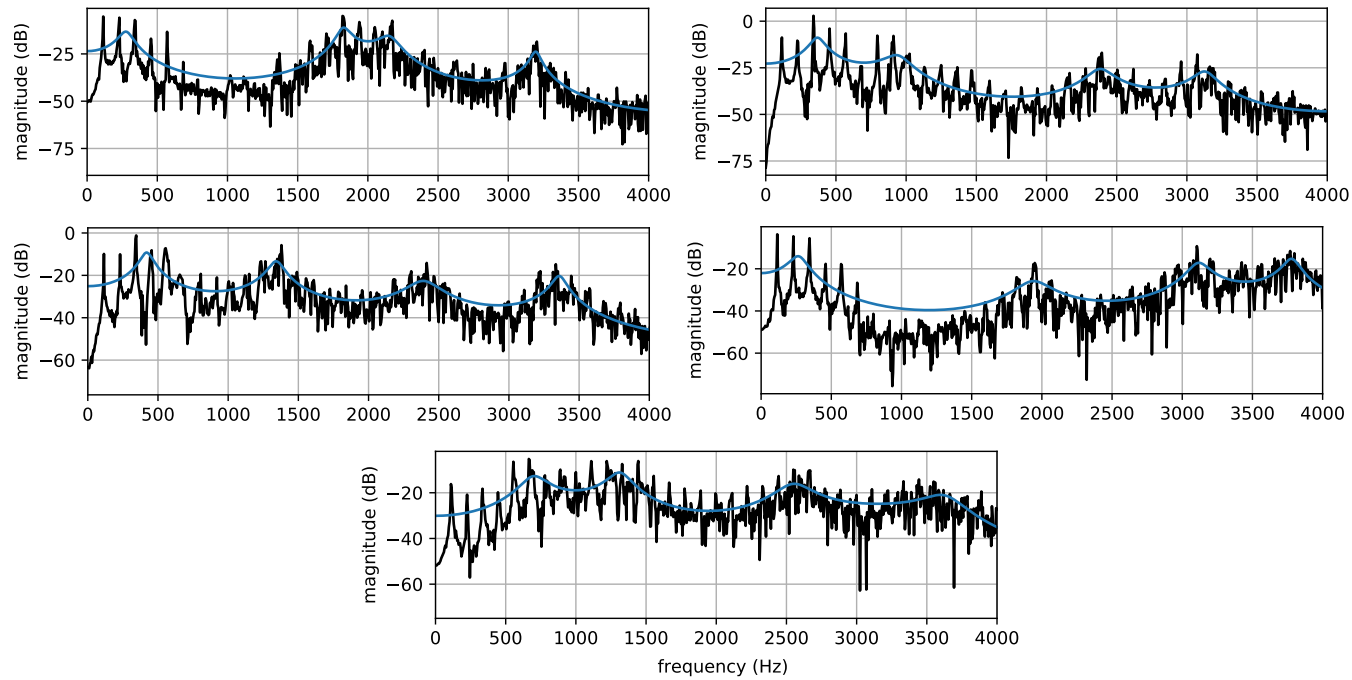
# The spectral envelope



- Black curve: power spectrum (in dB) of the recording of a vowel, computed with the DFT.
- Blue curve: **spectral envelope** showing the formant resonances, computed with linear predictive coding (will be discussed in the lab session).

Male or female speaker?

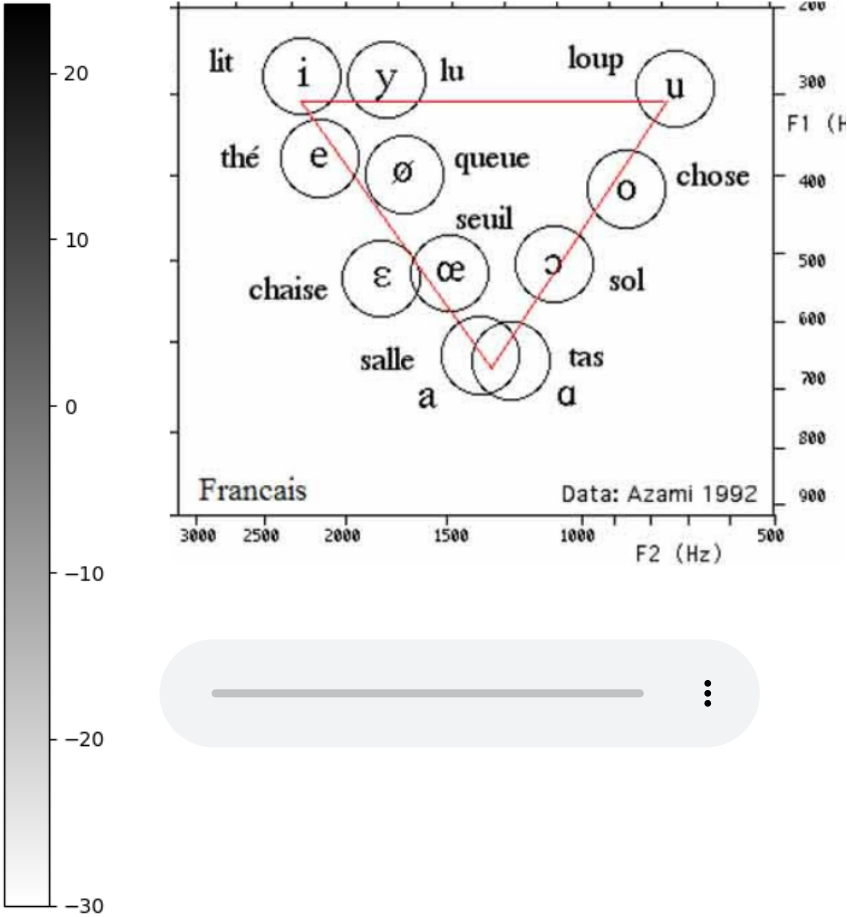
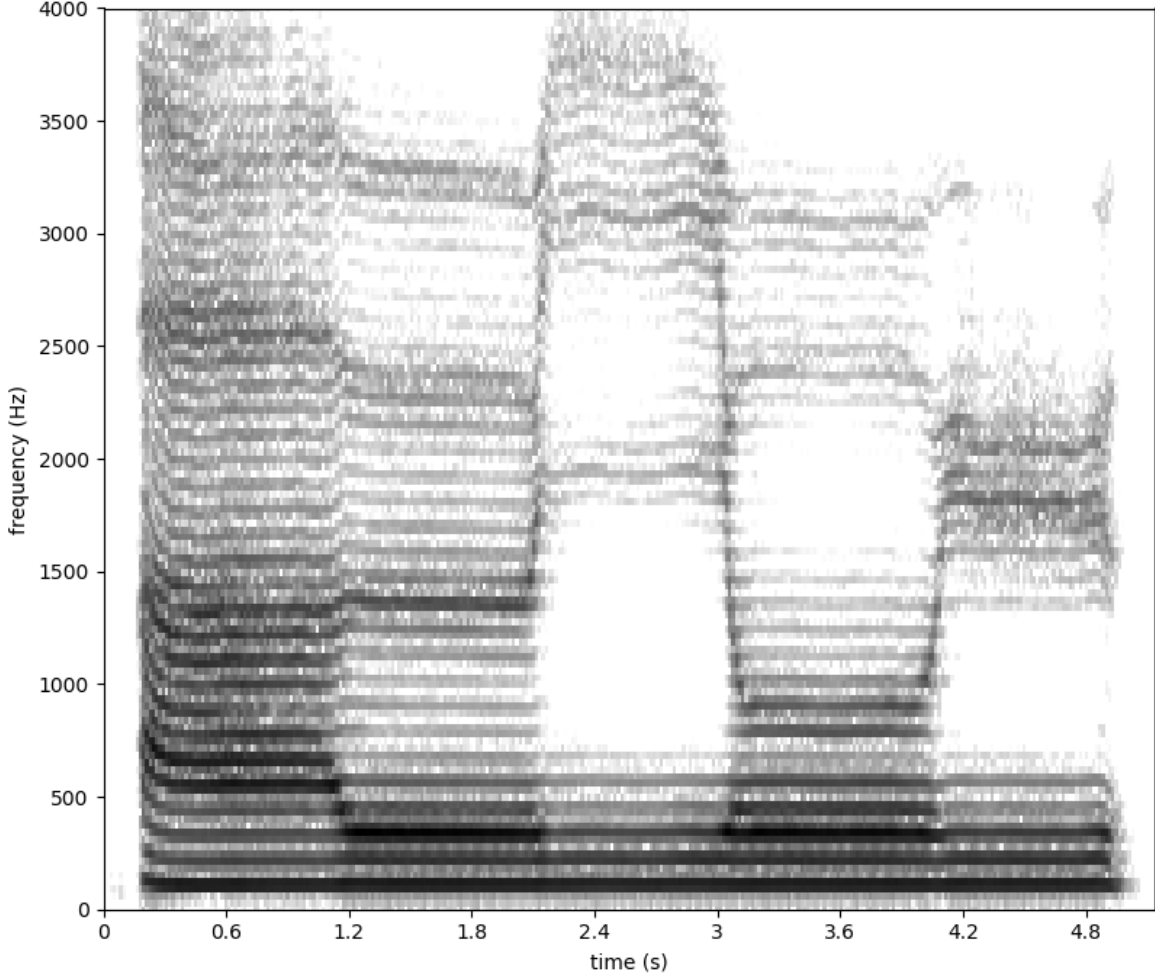




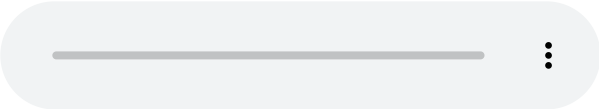
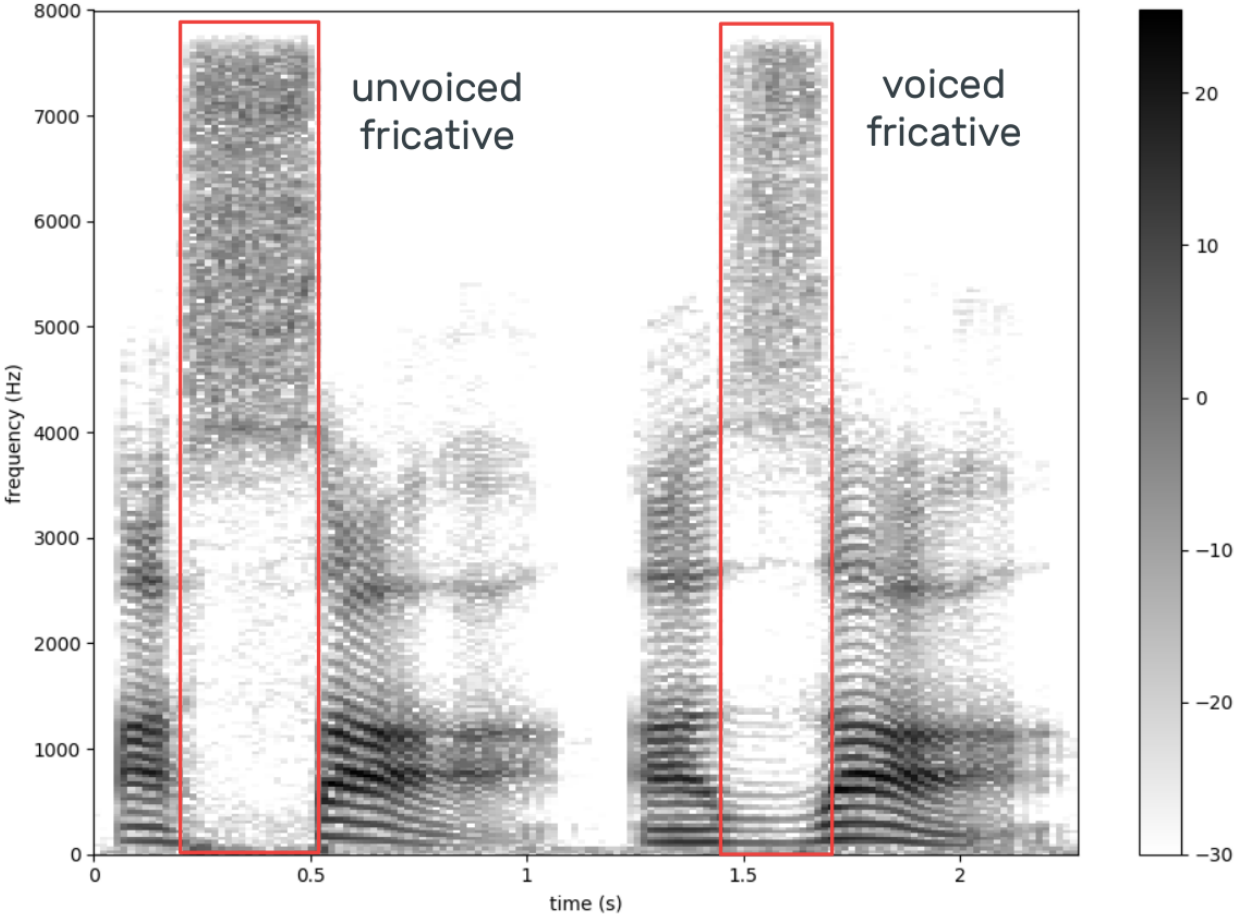
Go to <https://app.wooclap.com/CXIOJL> and find the vowel that corresponds to each spectrum, using the above French vocal triangle.

# Spectrogram reading - "aeiou"

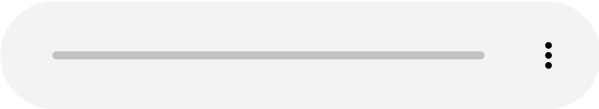
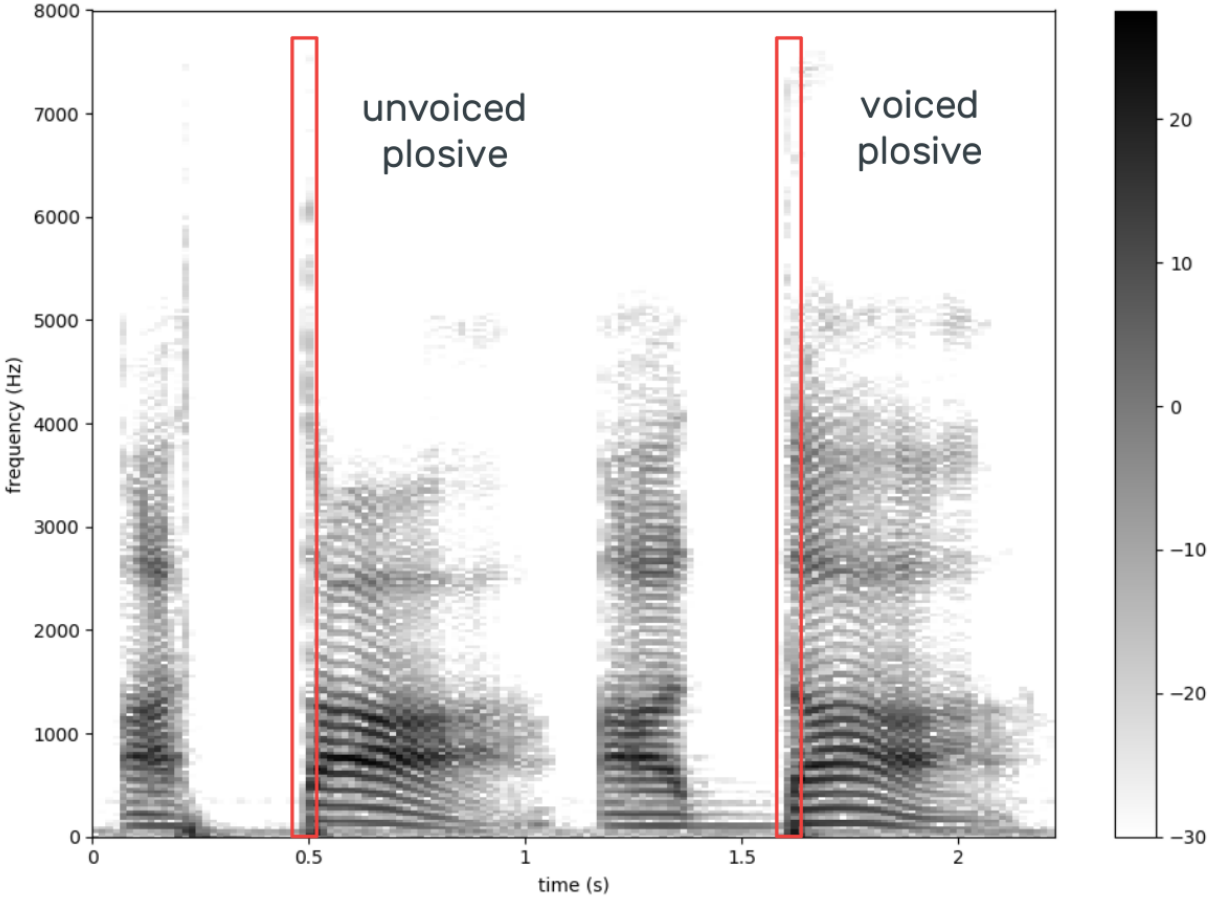
We could have done the same from a spectrogram representation.

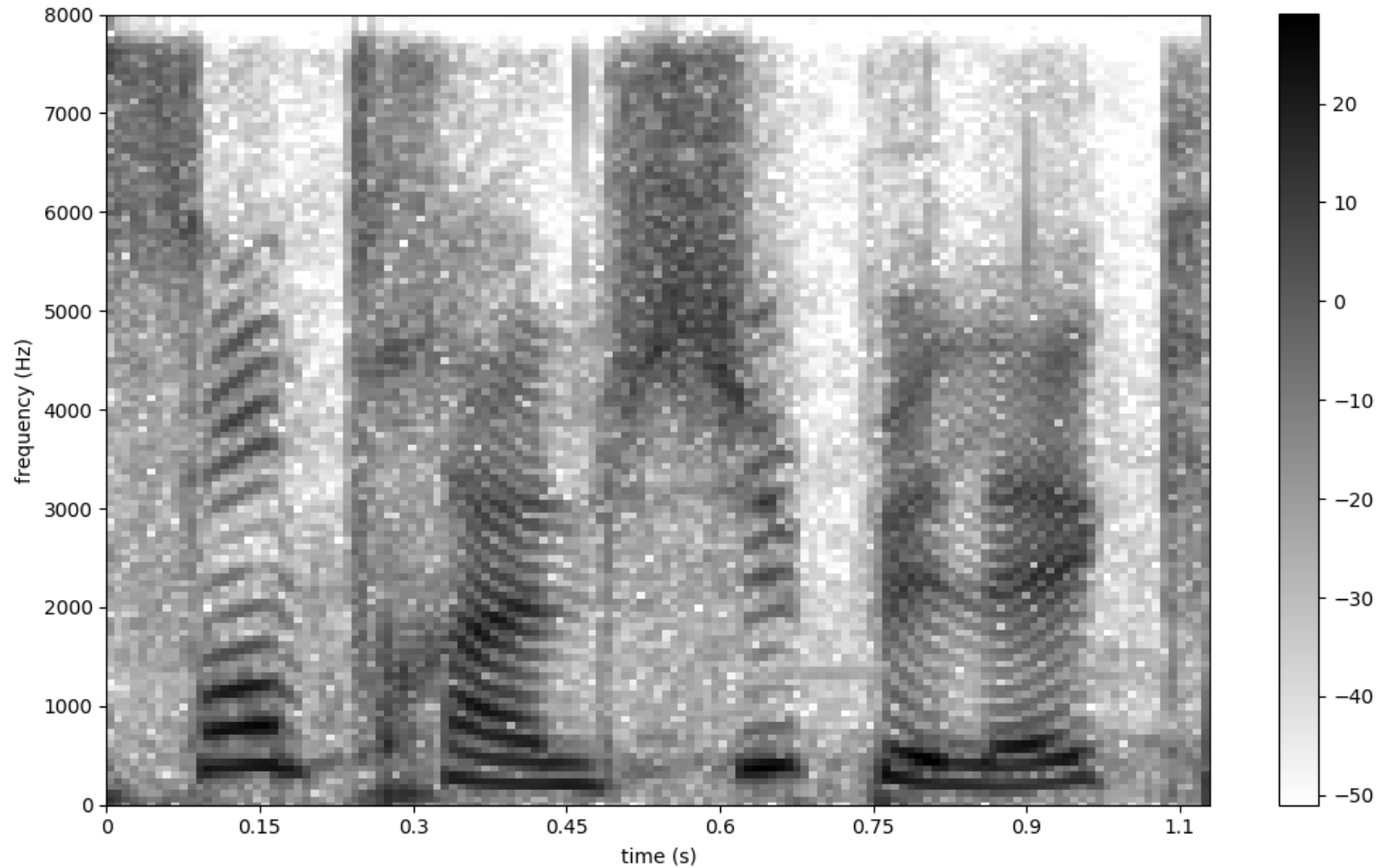


# Spectrogram reading - "assa - azza"



# Spectrogram reading - "atta - adda"





With a bit of practice you could be able to decode this mystery spectrogram.

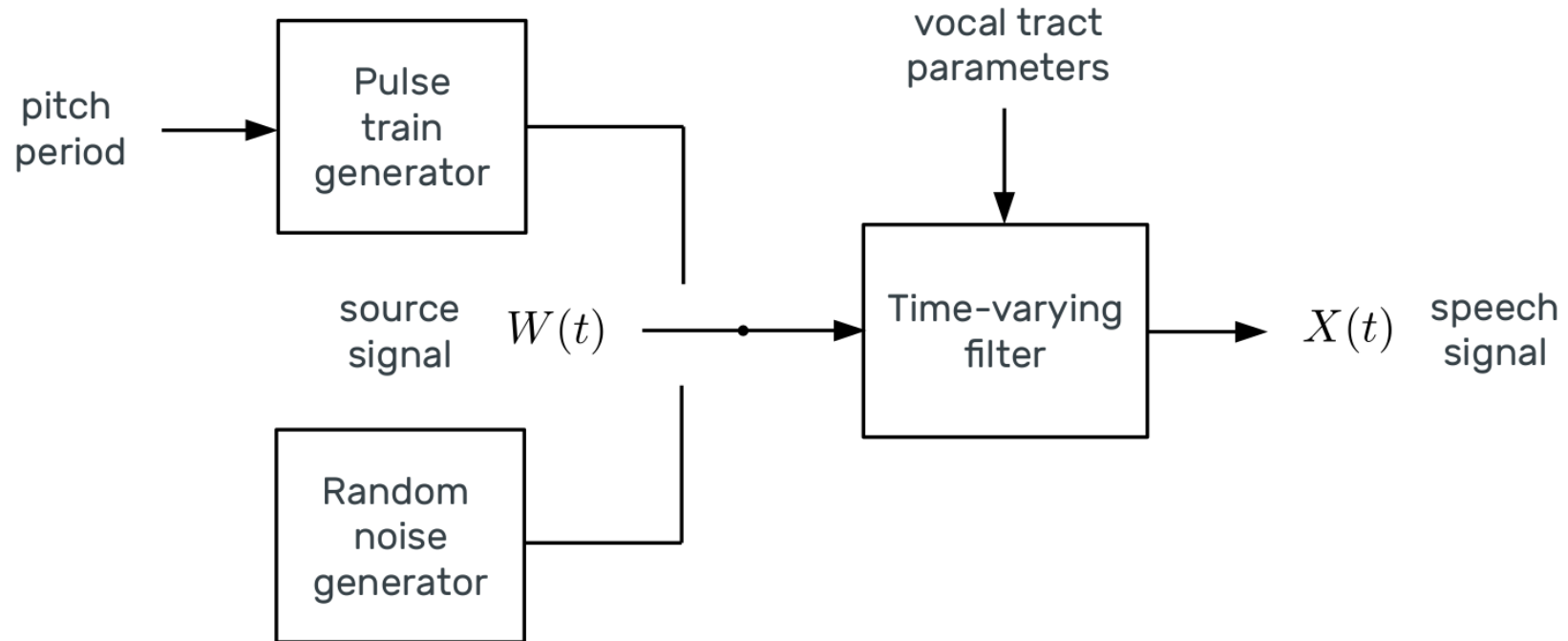
1 bonus point if you decode the message 😊.

## Further reading

Introduction to voice acoustics by Joe Wolfe, Emeritus Professor at the University of New South Wales (Sydney, Australia):

<https://newt.phys.unsw.edu.au/jw/voice.html>

# Lab session



Analysis, transformation and synthesis of speech signals with the **source-filter model** and **linear predictive coding**

# **Solution to the wooclap**



