

## Lecture 2 & 3: Speech Analysis and Linguistics Speech Basics

**Instructor: Dr Chiranjeevi Yarra**

Speech lab, LTRC



Aug 03& 07, 2023

# Outline

- 1 Introduction
- 2 Short-time Analysis
  - Short-time Fourier Transform (STFT)
- 3 Production based knowledge
  - Speech Excitation
  - Vocal tract
- 4 Conclusion

## 1 Introduction

## 2 Short-time Analysis

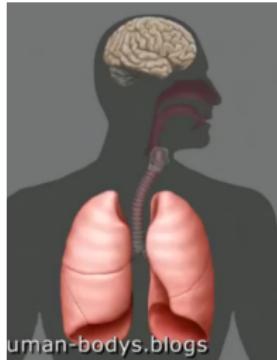
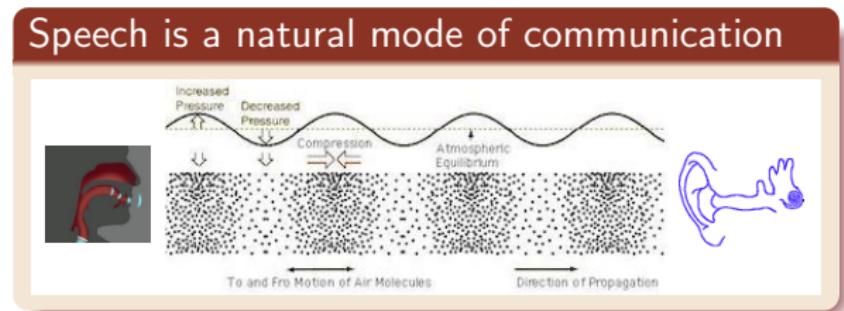
- Short-time Fourier Transform (STFT)

## 3 Production based knowledge

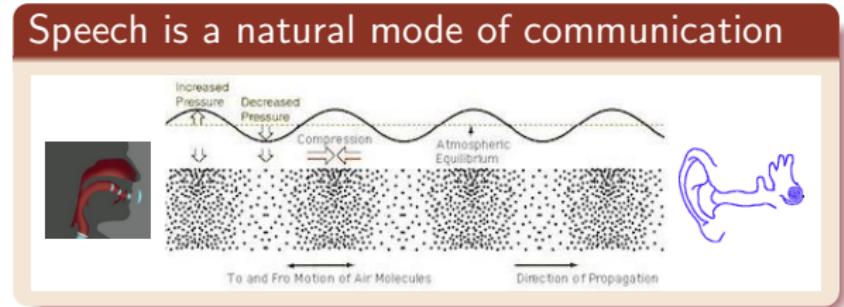
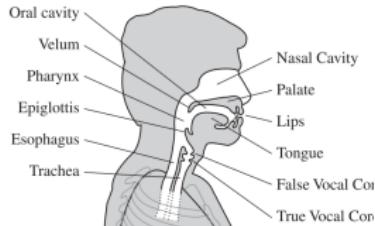
- Speech Excitation
- Vocal tract

## 4 Conclusion

# Speech Production Mechanism



# Speech Production Mechanism



## Non-stationary nature

### Illustrative example

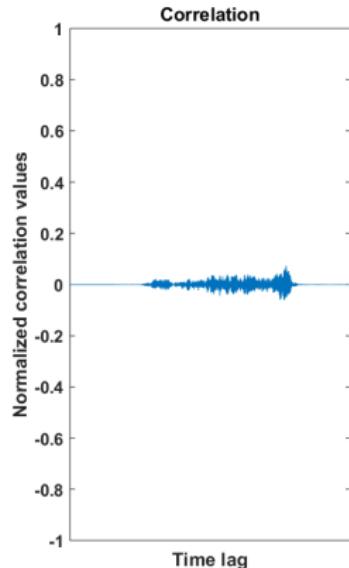
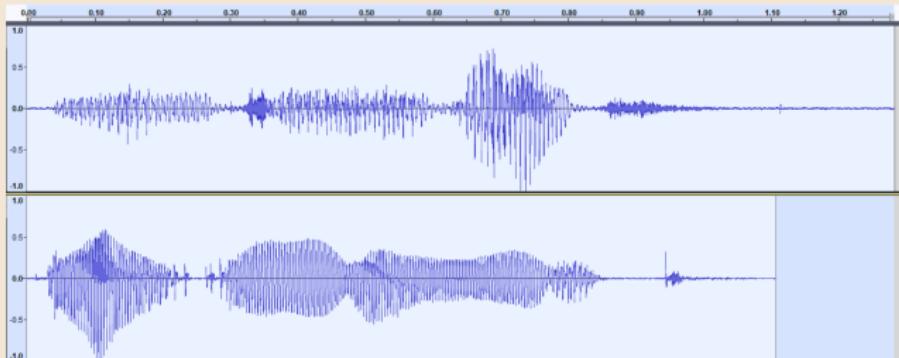
Sentence: "You enjoyed it"; Signal #1: , Signal #2:



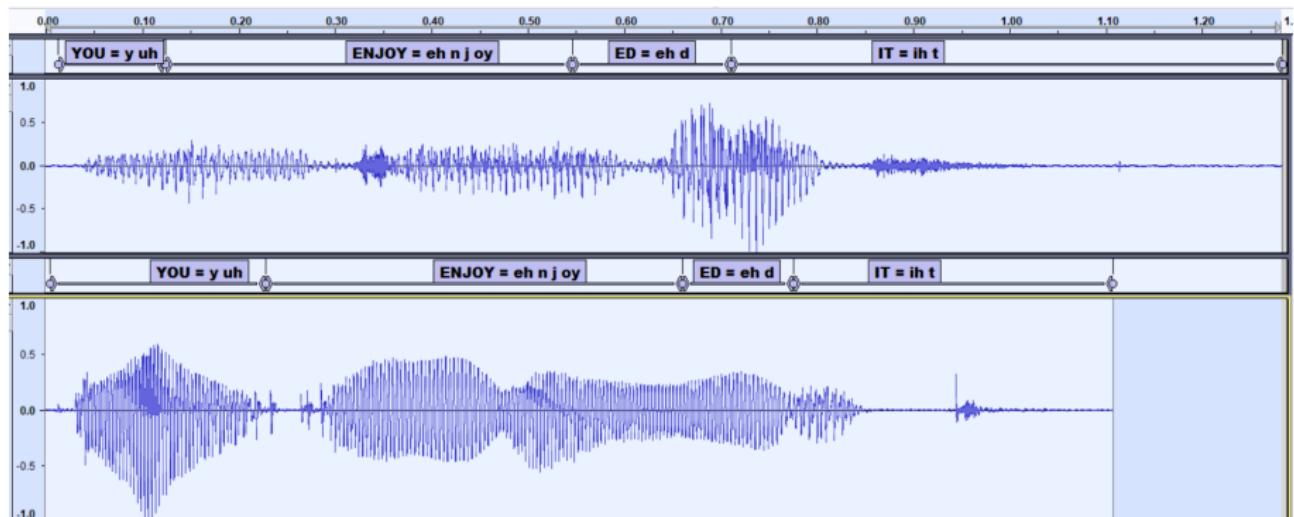
## Non-stationary nature

### Illustrative example

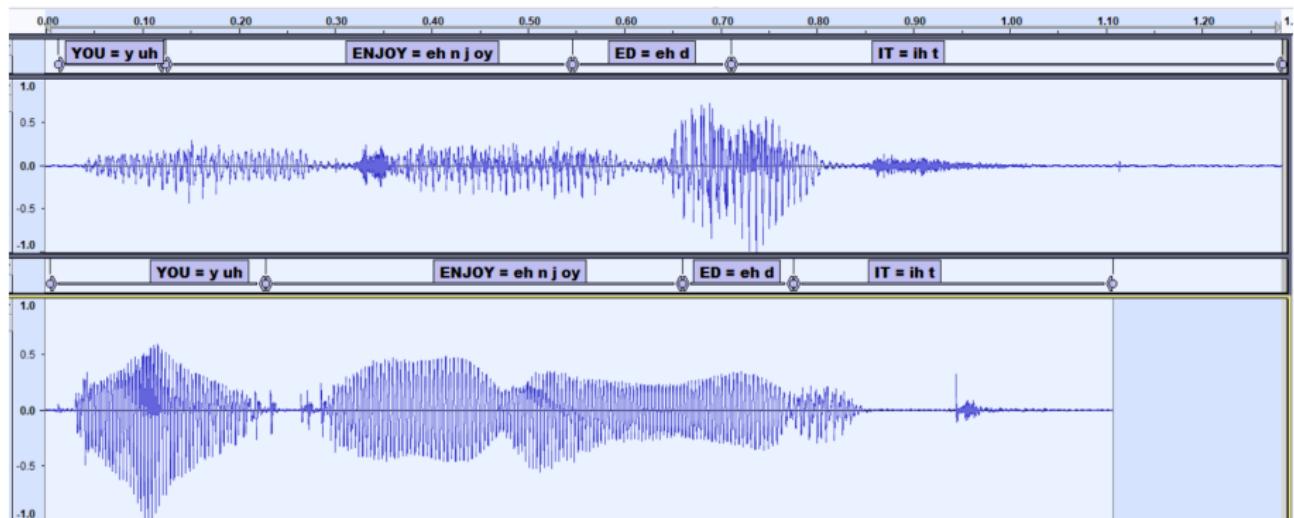
Sentence: "You enjoyed it"; Signal #1: , Signal #2:



# Time-varying Nature



# Unknown Segment Boundaries



# Summary

- 1 Non-stationary
- 2 Time-varying nature
- 3 Unknown Segment Boundaries

1 Introduction

2 Short-time Analysis

■ Short-time Fourier Transform (STFT)

3 Production based knowledge

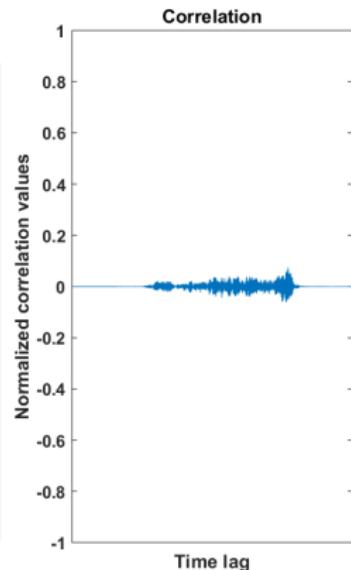
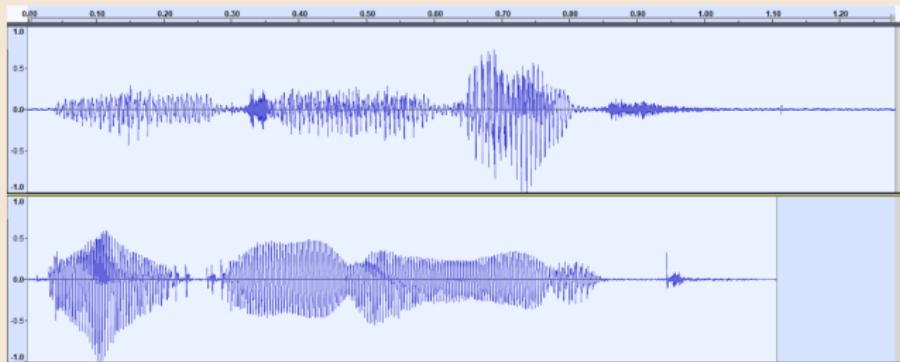
- Speech Excitation
- Vocal tract

4 Conclusion

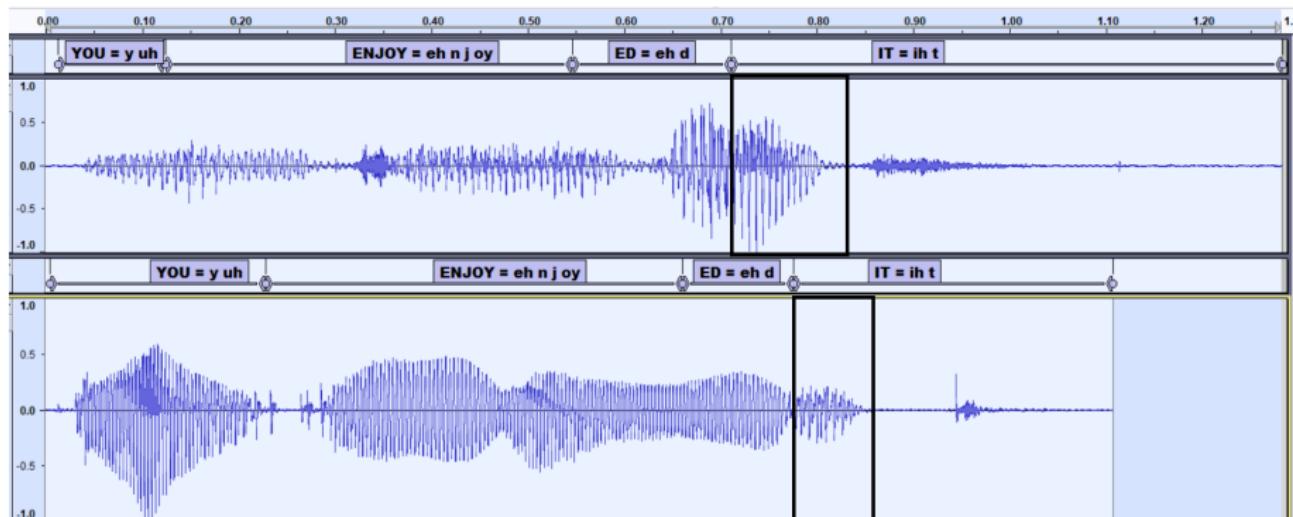
## Non-stationary nature

### Illustrative example

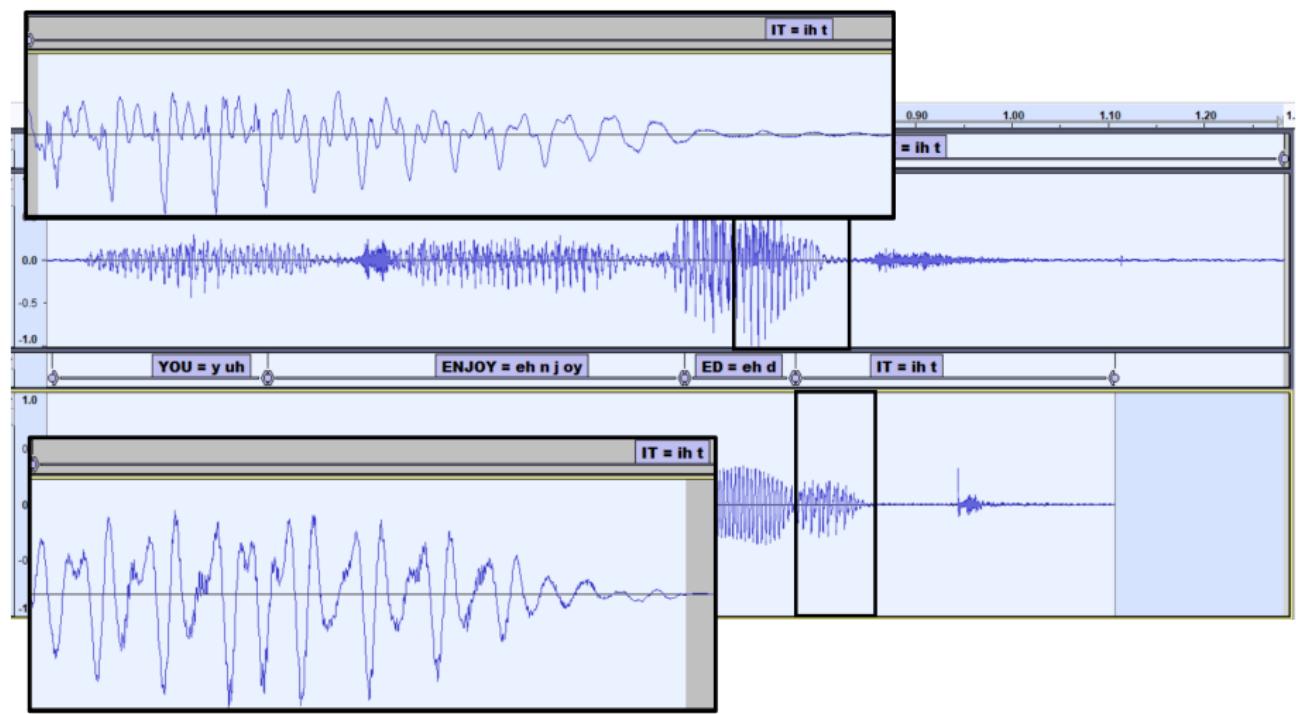
Sentence: "You enjoyed it"; Signal #1: Signal #2:



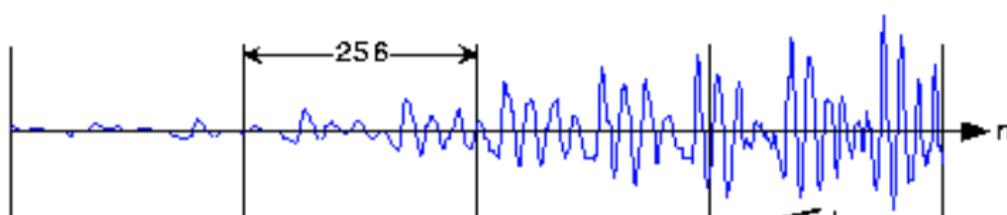
Non-stationary, but there is a structure in small segment



Non-stationary, but there is a structure in small segment



## Stationary within the segment



- Signal analysis within the short segments
- Short-time energy, Short-time zero crossing rate, Short-time Fourier Transform etc..

└ Short-time Analysis

  └ Short-time Fourier Transform (STFT)

## 1 Introduction

## 2 Short-time Analysis

### ■ Short-time Fourier Transform (STFT)

## 3 Production based knowledge

- Speech Excitation
- Vocal tract

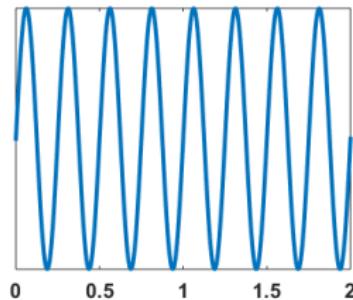
## 4 Conclusion

└ Short-time Analysis

  └ Short-time Fourier Transform (STFT)

# Fourier Transform (FT)

Sine signal:

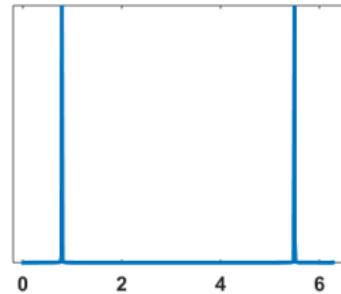
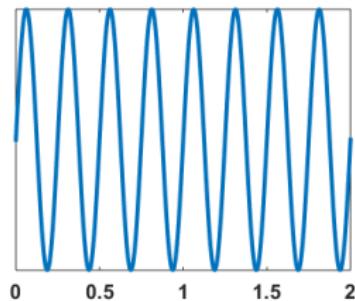


└ Short-time Analysis

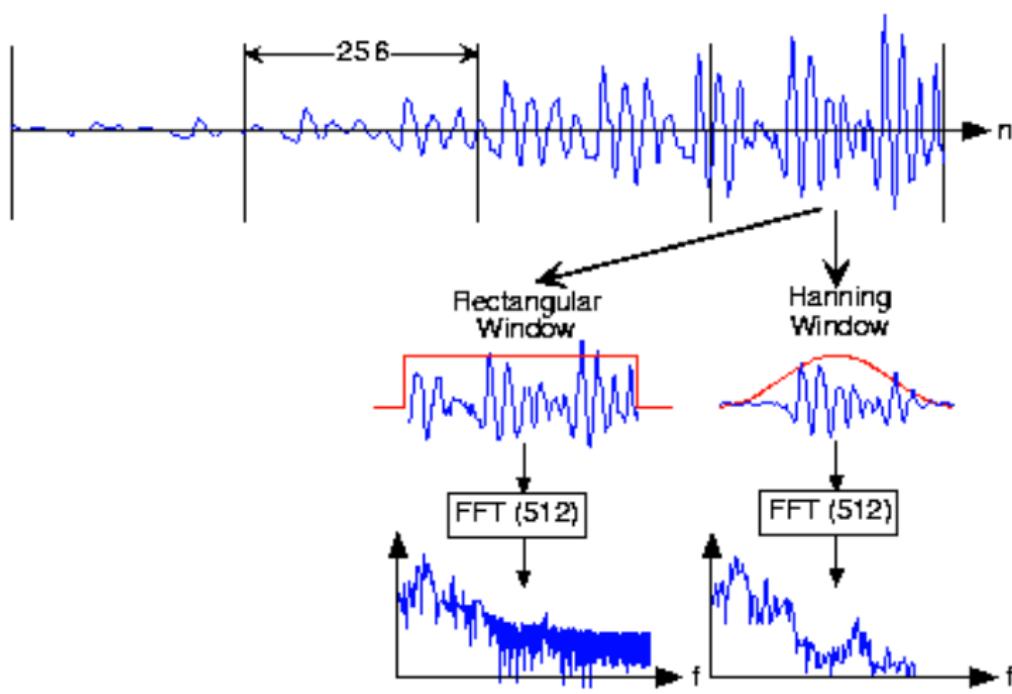
  └ Short-time Fourier Transform (STFT)

# Fourier Transform (FT)

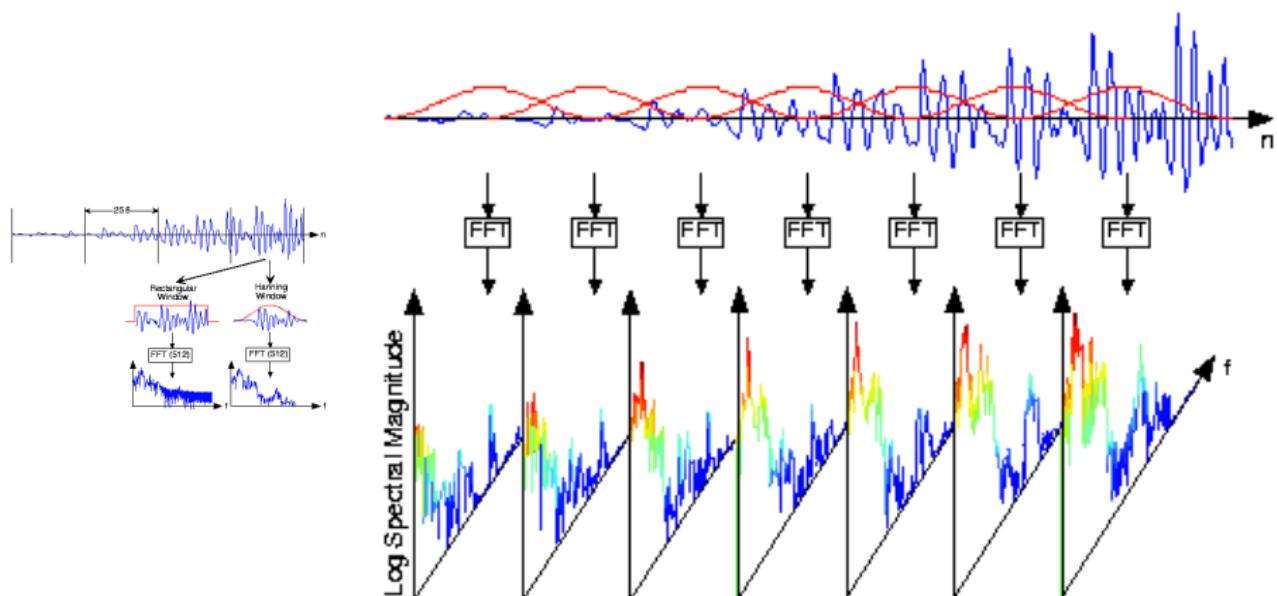
Sine signal:



# Short-time Fourier Transform (STFT)



# Short-time Fourier Transform (STFT)

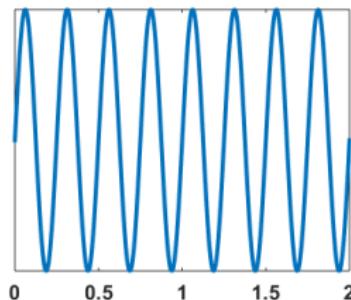


└ Short-time Analysis

  └ Short-time Fourier Transform (STFT)

# Short-time Fourier Transform (STFT)

Sine signal:

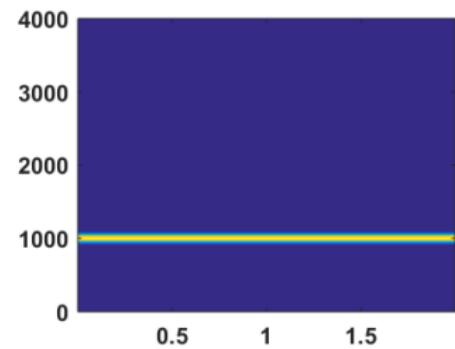
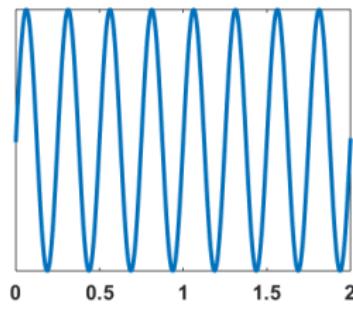


└ Short-time Analysis

  └ Short-time Fourier Transform (STFT)

# Short-time Fourier Transform (STFT)

Sine signal:

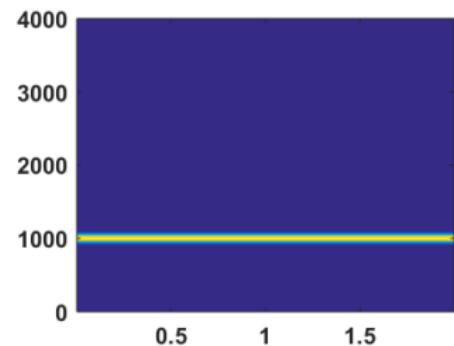
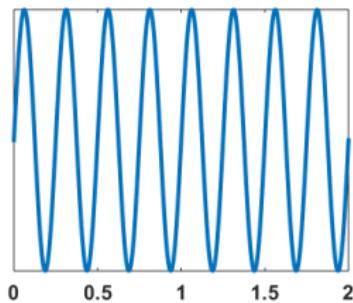


└ Short-time Analysis

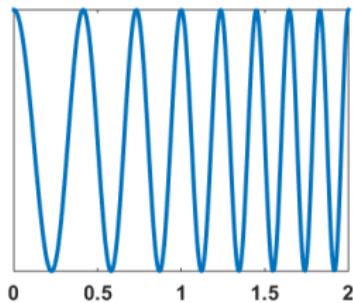
  └ Short-time Fourier Transform (STFT)

## Short-time Fourier Transform (STFT)

Sine signal:



Chirp Signal:

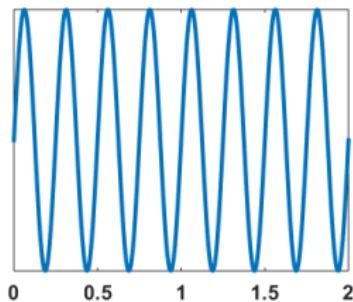


└ Short-time Analysis

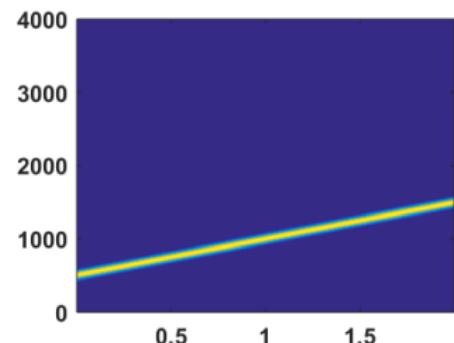
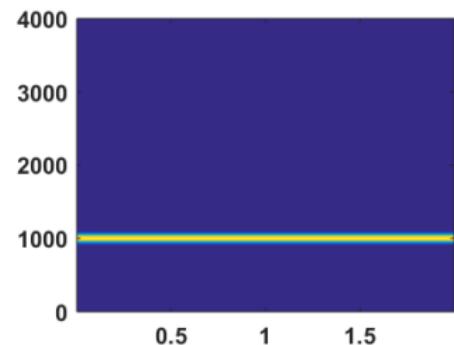
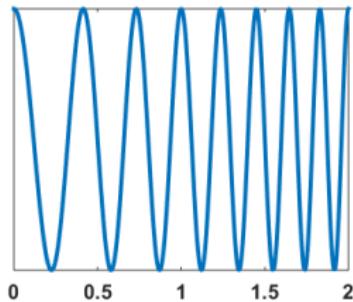
  └ Short-time Fourier Transform (STFT)

## Short-time Fourier Transform (STFT)

Sine signal:



Chirp Signal:

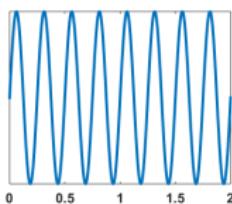


└ Short-time Analysis

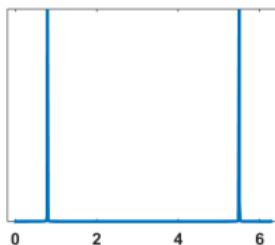
  └ Short-time Fourier Transform (STFT)

# Why STFT?

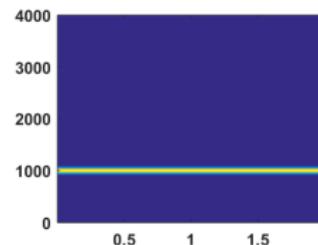
Sine signal:



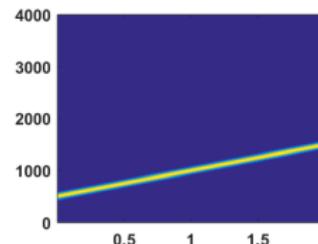
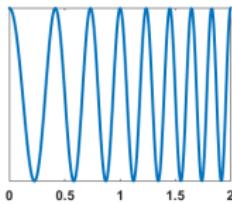
DFT:



STFT:



Chirp Signal:

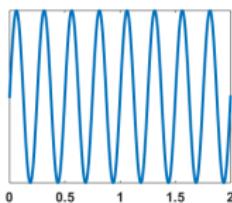


└ Short-time Analysis

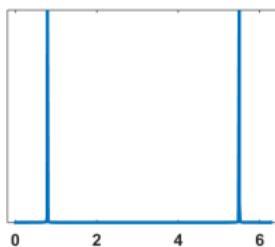
  └ Short-time Fourier Transform (STFT)

# Why STFT?

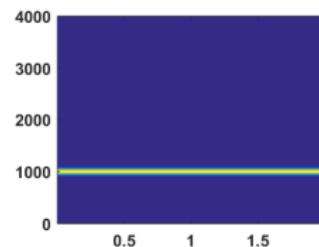
Sine signal:



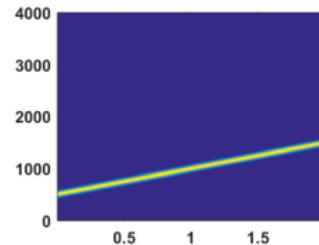
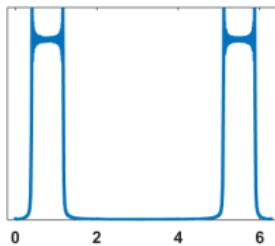
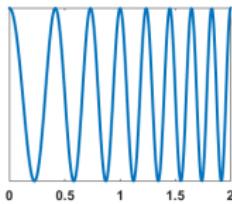
DFT:



STFT:



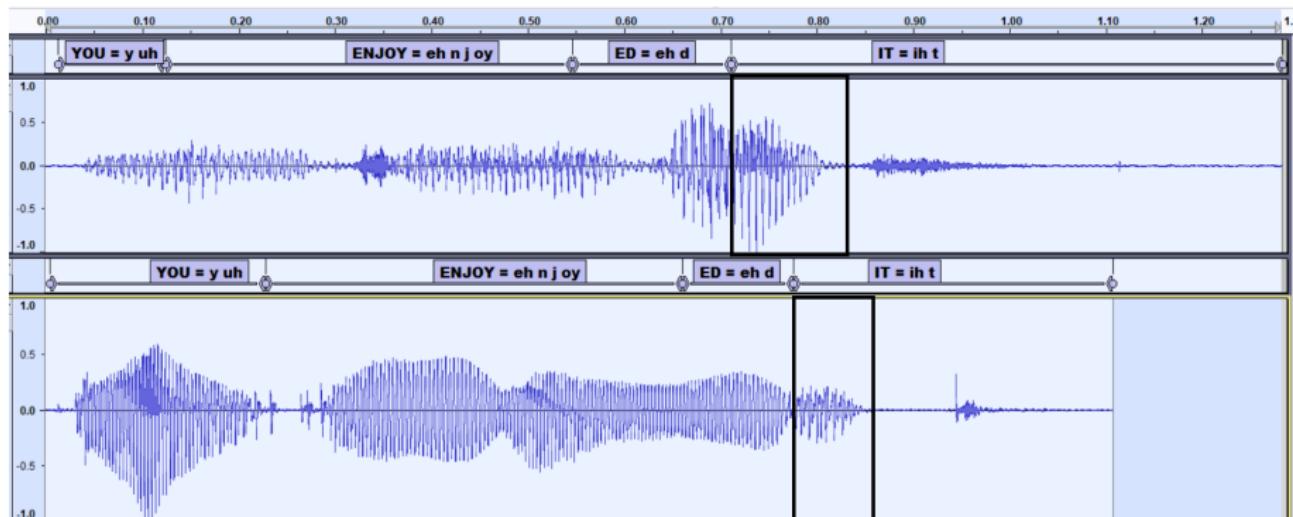
Chirp Signal:



└ Short-time Analysis

  └ Short-time Fourier Transform (STFT)

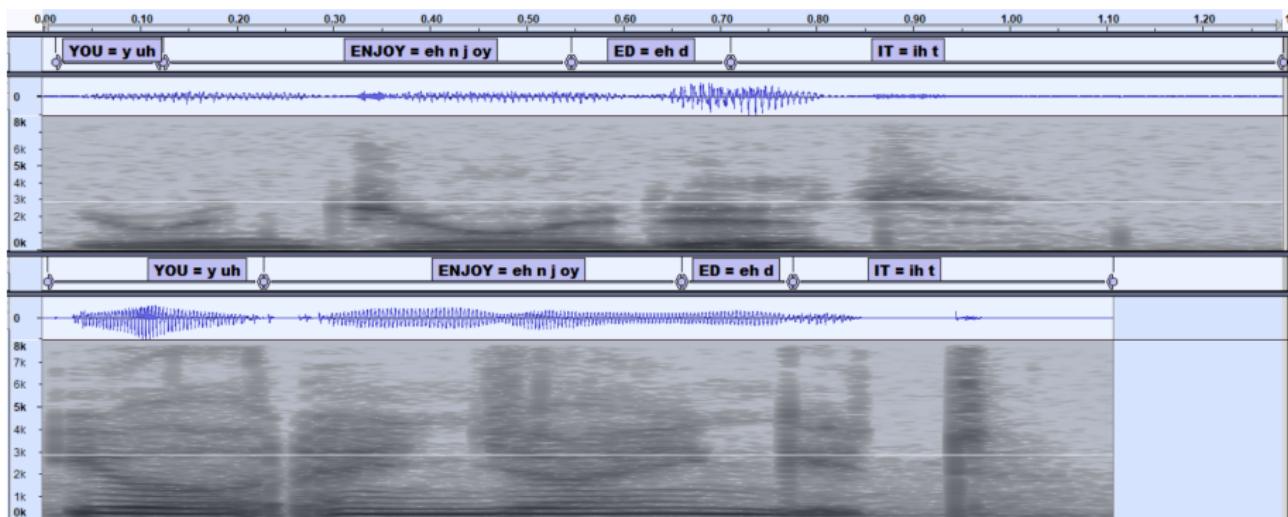
## For the exemplary speech



└ Short-time Analysis

  └ Short-time Fourier Transform (STFT)

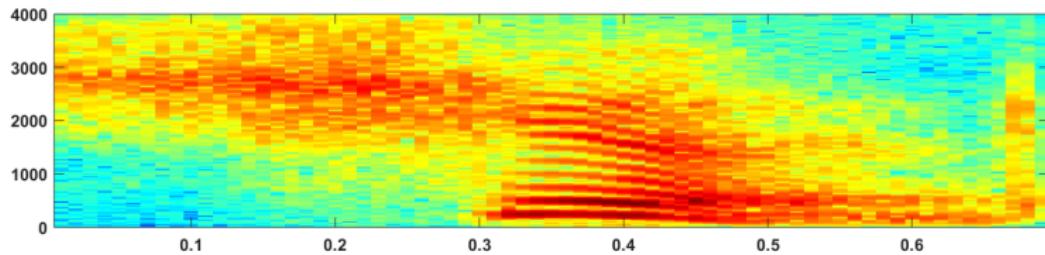
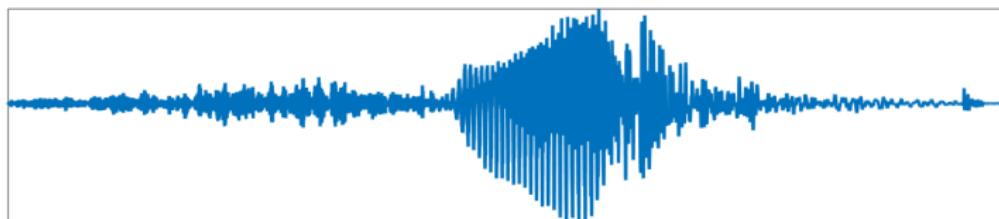
## For the exemplary speech



└ Short-time Analysis

  └ Short-time Fourier Transform (STFT)

## For an exemplary speech



1 Introduction

2 Short-time Analysis

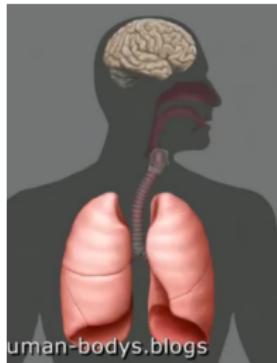
- Short-time Fourier Transform (STFT)

3 Production based knowledge

- Speech Excitation
- Vocal tract

4 Conclusion

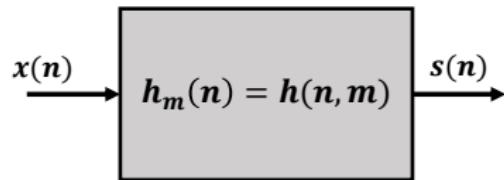
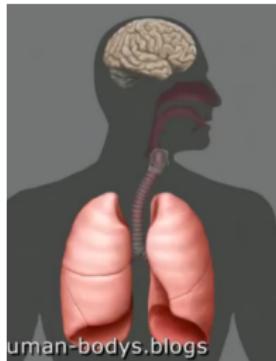
## Speech production approximation



uman-bodys.blogs



## Speech production approximation



$$s(n) = \sum_{m=-\infty}^{\infty} h(n, m)x(n - m)$$

└ Production based knowledge

  └ Speech Excitation

## 1 Introduction

## 2 Short-time Analysis

  ■ Short-time Fourier Transform (STFT)

## 3 Production based knowledge

  ■ Speech Excitation

  ■ Vocal tract

## 4 Conclusion

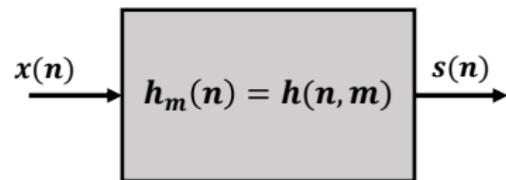
└ Production based knowledge

  └ Speech Excitation

## Vocal chord movements

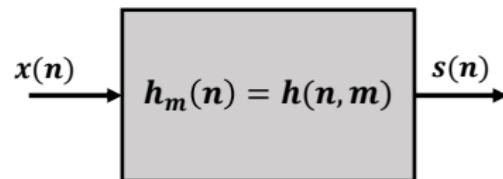


## VuV based modelling

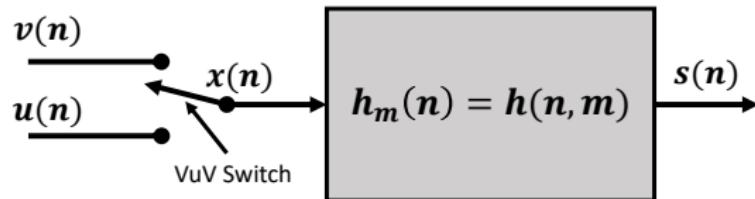


$$s(n) = \sum_{m=-\infty}^{\infty} h(n, m)x(n - m)$$

## VuV based modelling



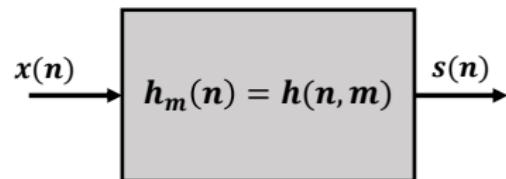
$$s(n) = \sum_{m=-\infty}^{\infty} h(n, m)x(n - m)$$



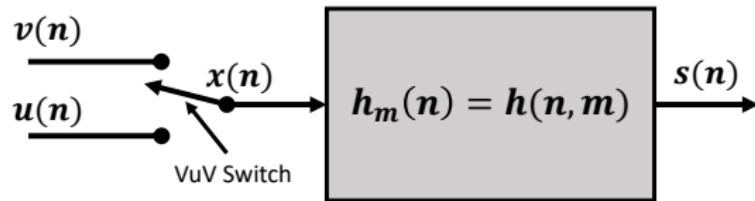
└ Production based knowledge

└ Speech Excitation

## VuV based modelling



$$s(n) = \sum_{m=-\infty}^{\infty} h(n, m)x(n - m)$$



$$s_v(n) = \sum_{m=-\infty}^{\infty} h(n, m)v(n - m); \quad s_u(n) = \sum_{m=-\infty}^{\infty} h(n, m)u(n - m)$$

- └ Production based knowledge

- └ Speech Excitation

## VuV detection

└ Production based knowledge

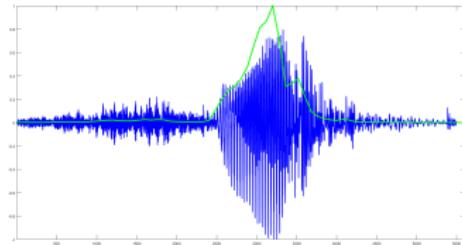
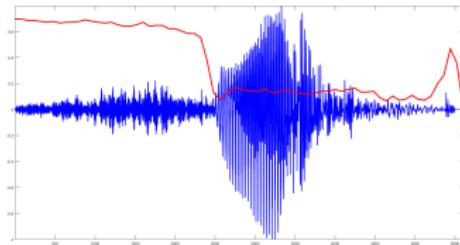
  └ Speech Excitation

## VuV detection

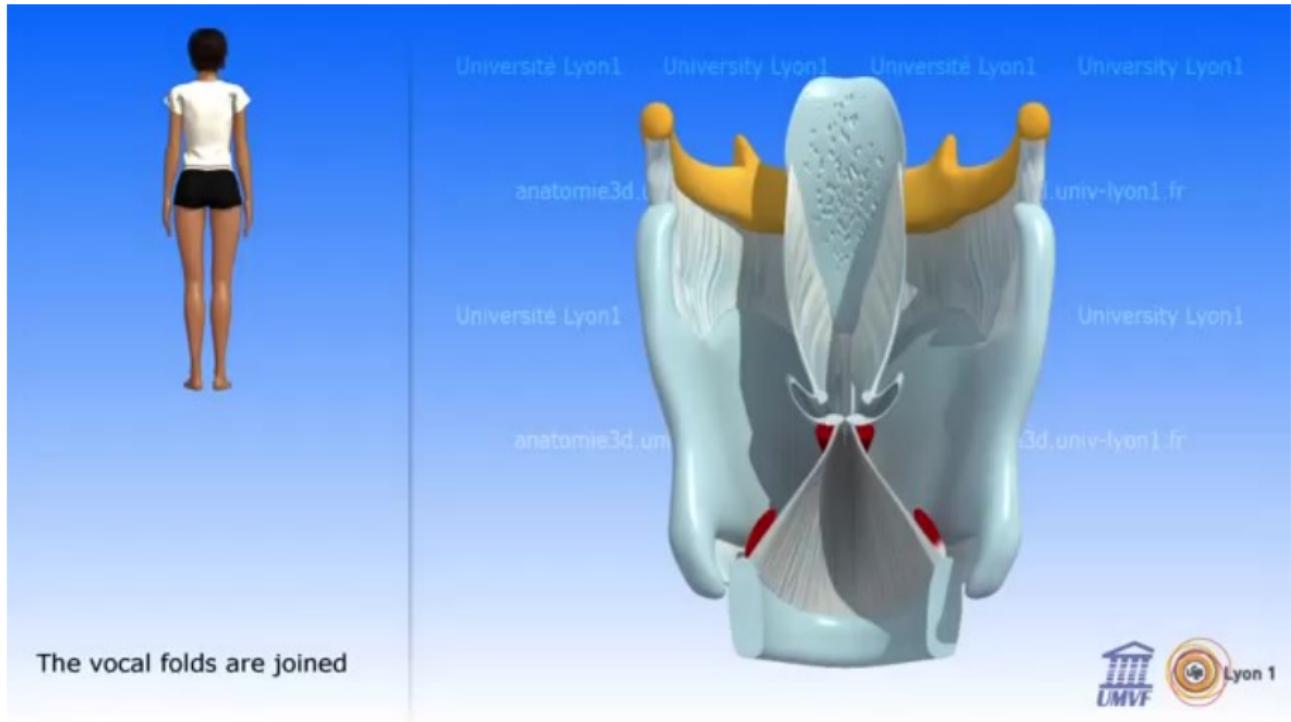
- Short-time zero crossing rate
- Short-time energy

## VuV detection

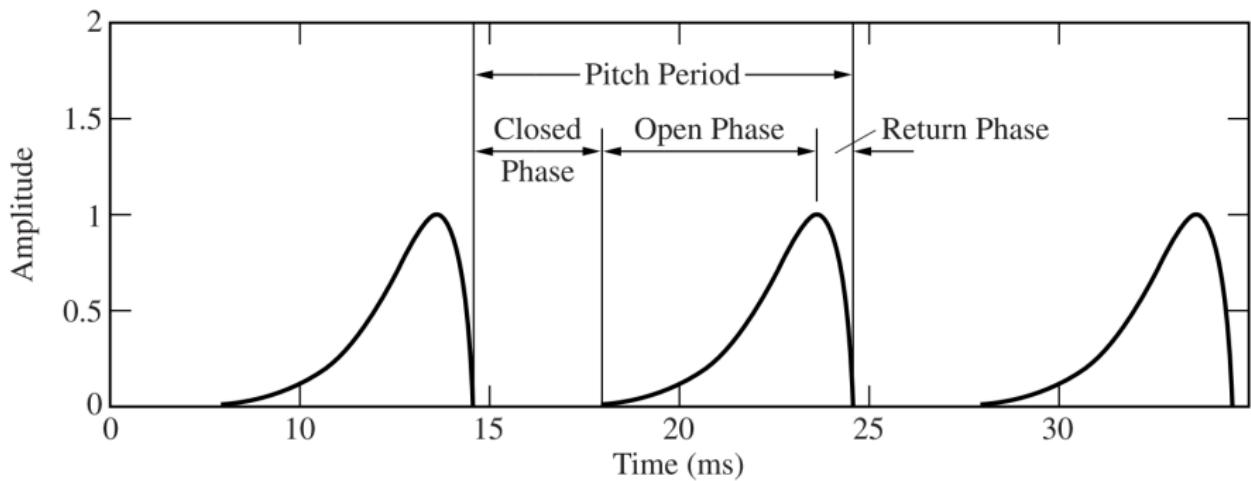
- Short-time zero crossing rate
- Short-time energy



## Airflow movement at the vocal chords



## Glottal airflow velocity



└ Production based knowledge

  └ Vocal tract

## 1 Introduction

## 2 Short-time Analysis

  ■ Short-time Fourier Transform (STFT)

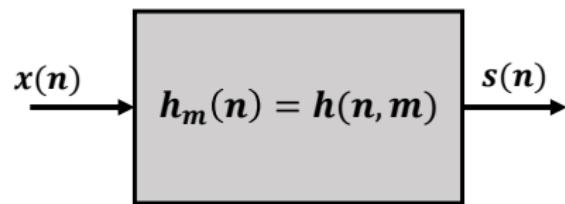
## 3 Production based knowledge

  ■ Speech Excitation

  ■ Vocal tract

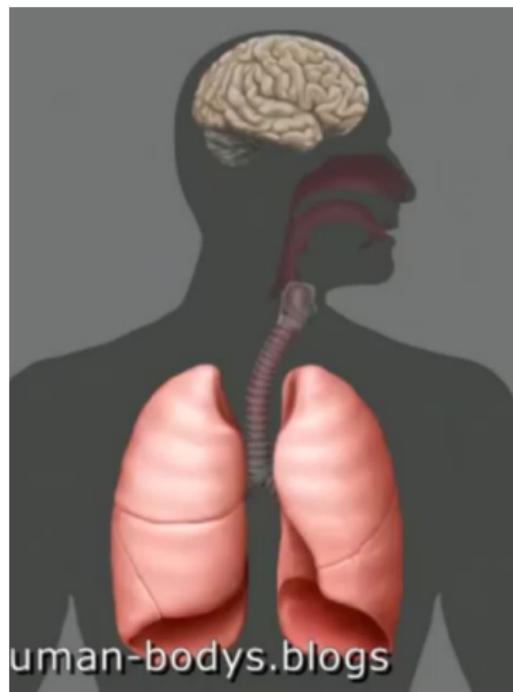
## 4 Conclusion

## Vocal-tract modelling

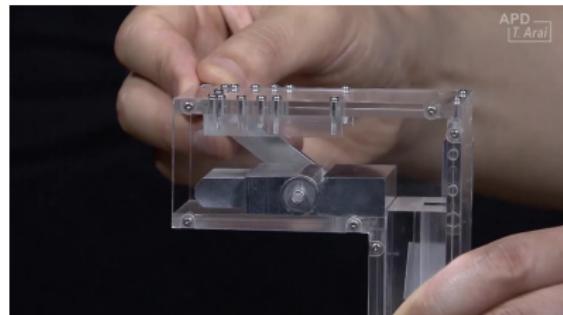
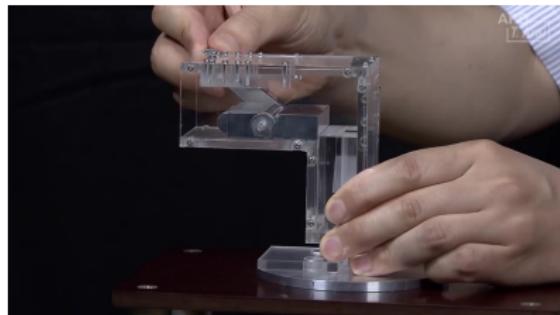


$$s(n) = \sum_{m=-\infty}^{\infty} h(n, m)x(n - m)$$

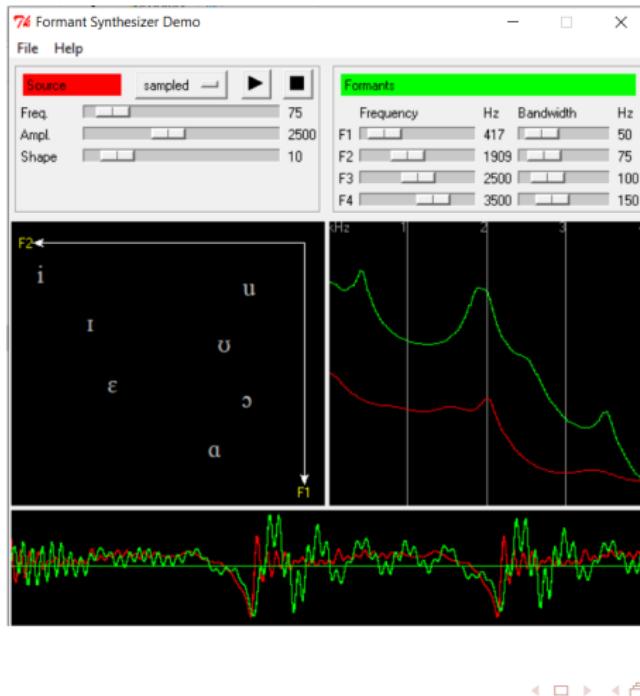
## Vocal-tract approximation - cascade of tubes



## Illustrations of approximation



# Vocal tract modelling



## 1 Introduction

## 2 Short-time Analysis

- Short-time Fourier Transform (STFT)

## 3 Production based knowledge

- Speech Excitation
- Vocal tract

## 4 Conclusion

# Conclusion

- What is nature of the speech signal?
- What is short-time analysis and frames?
- What are relative lengths of frames and segments?
- What are the voiced and unvoiced sounds?
- What is pitch?
- What are the formant frequencies?
- How the sine signal characterize?
- How the STFT (spectrogram) is useful over DFT?

Thank you