# Audio Signal Processing

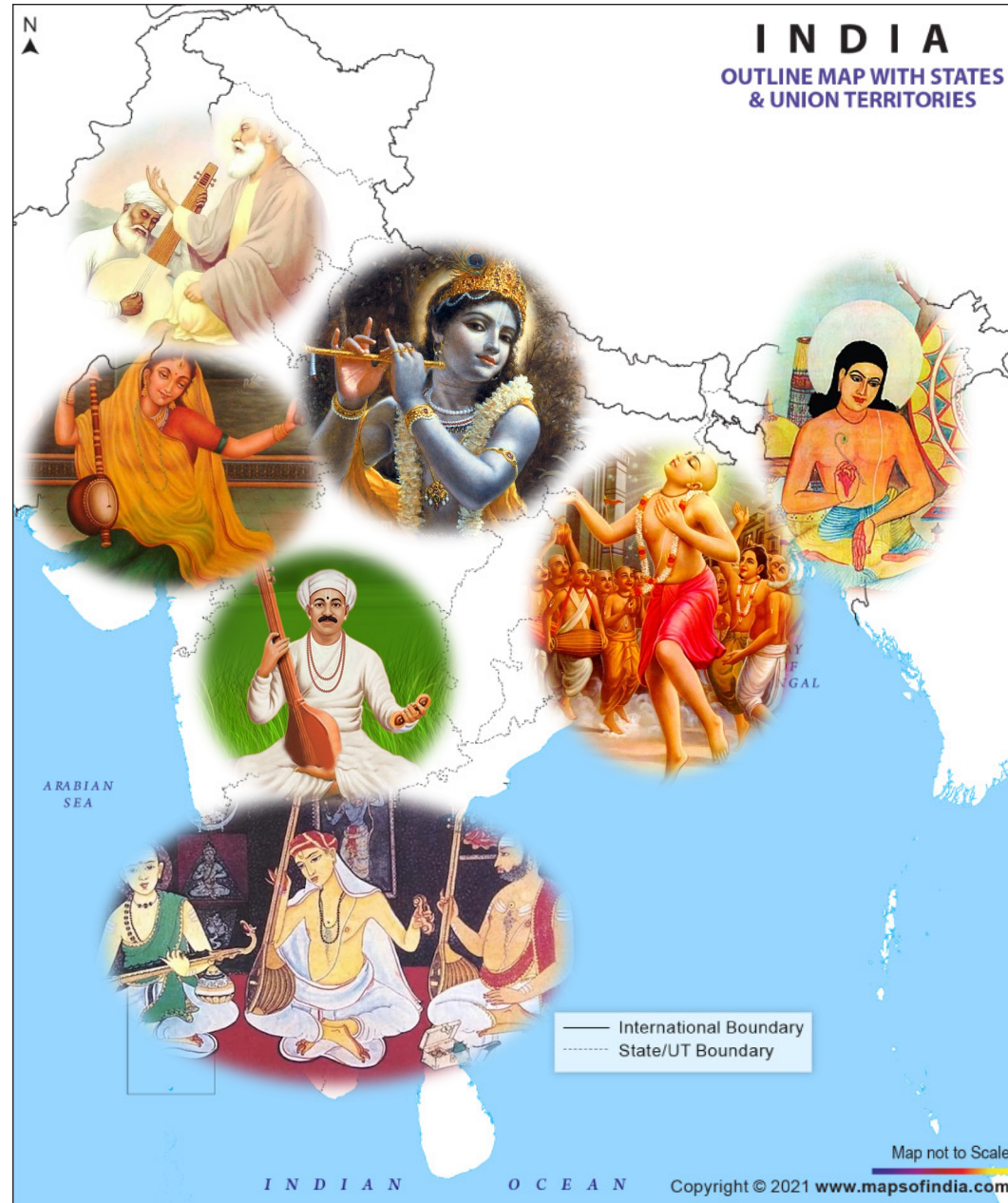## Vipul Arora
## Department of EE, IIT Kanpur

Vipul Arora
Department of EE, IIT Kanpur

MADHAV *lab*
Machine Analysis of Data
for Human Audition and Vision

# WiSSAP Class Rules

1. Ask questions
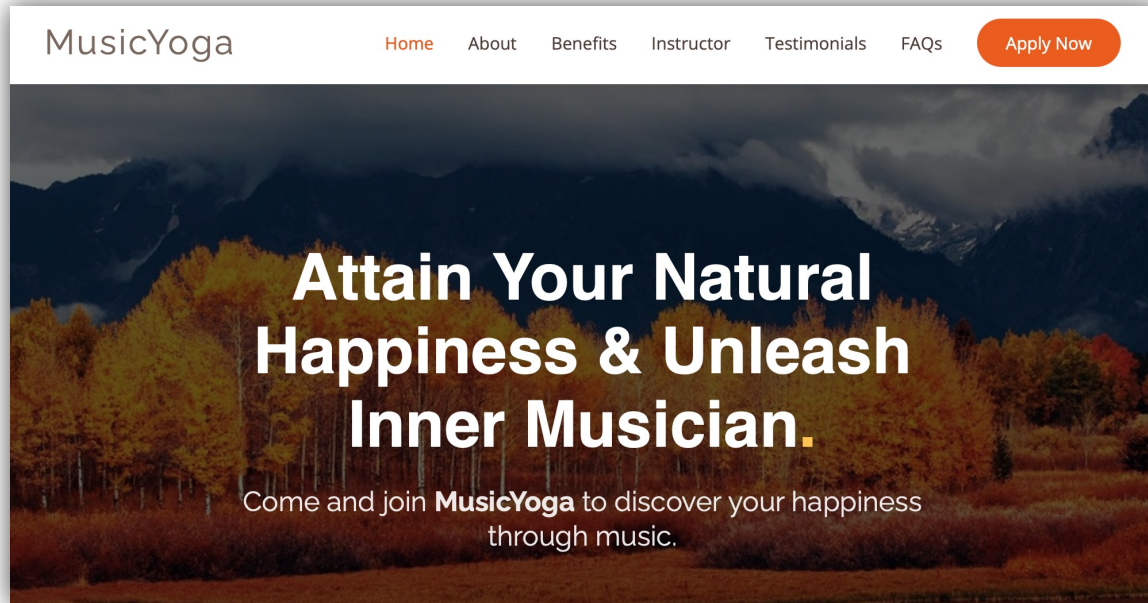2. Ask questions
3. Ask questions
4. Ask questions
5. …

Source: internet

# Music in Society









Source: internet

# Music and Health
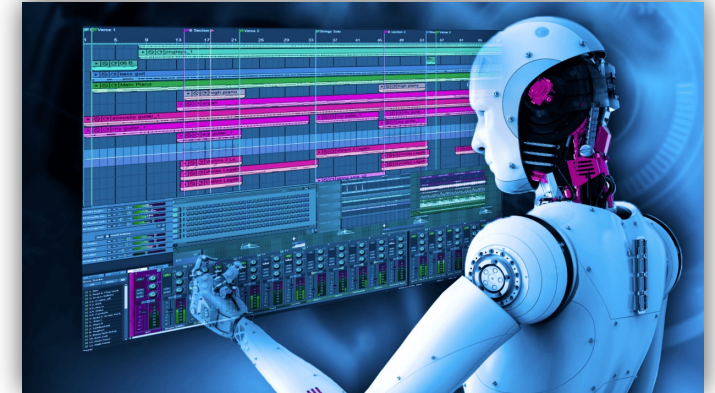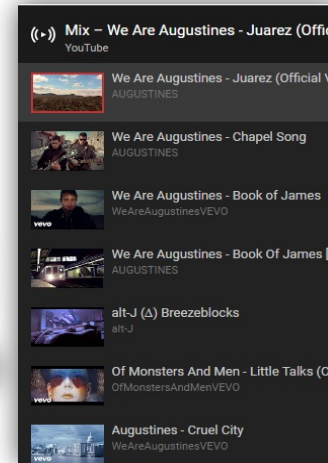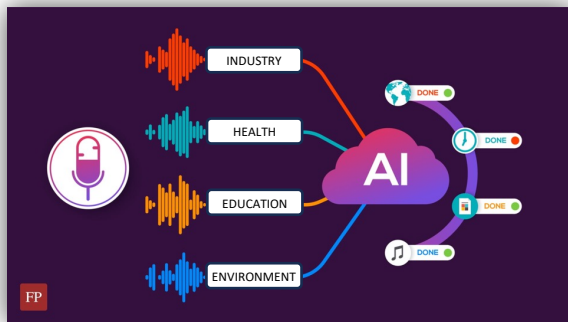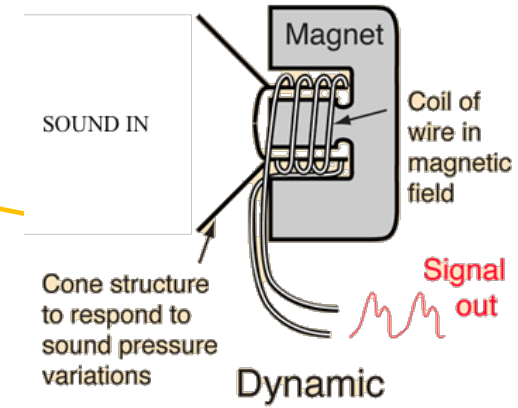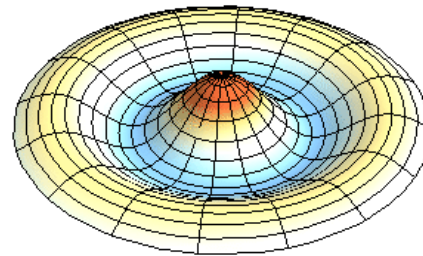
# Music and Development

# Audio Signal Processing and AI

# Listen



This Photo by Unknown Author is licensed under CC BY-NC-ND



This Photo by Unknown Author is licensed under CC BY-NC

# Digital Audio

SOUND IN

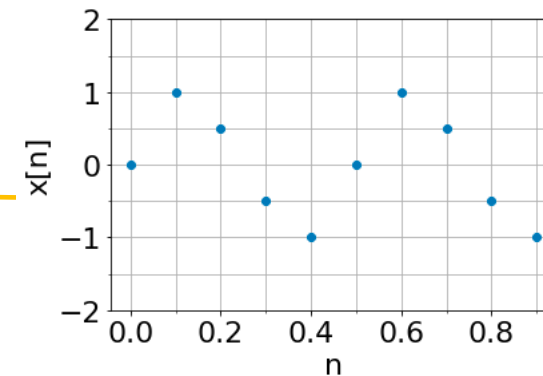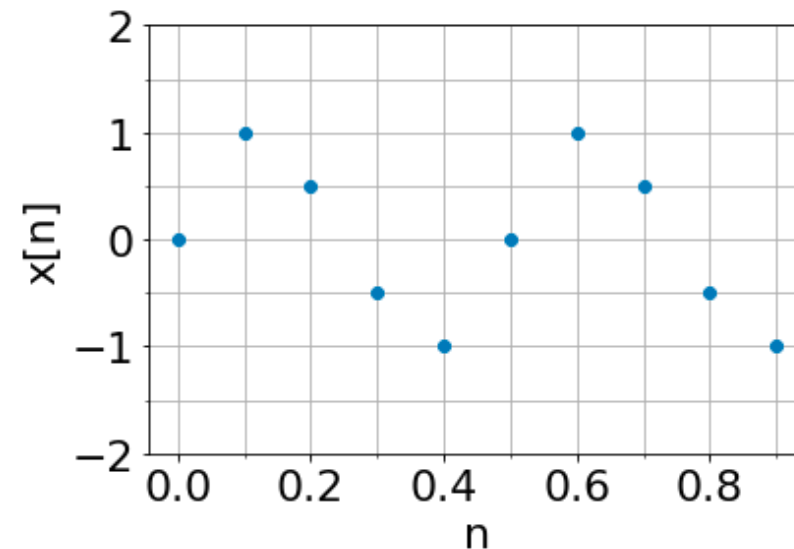Magnet

Coil of wire in magnetic field

Cone structure to respond to sound pressure variations

**Signal out**

Dynamic

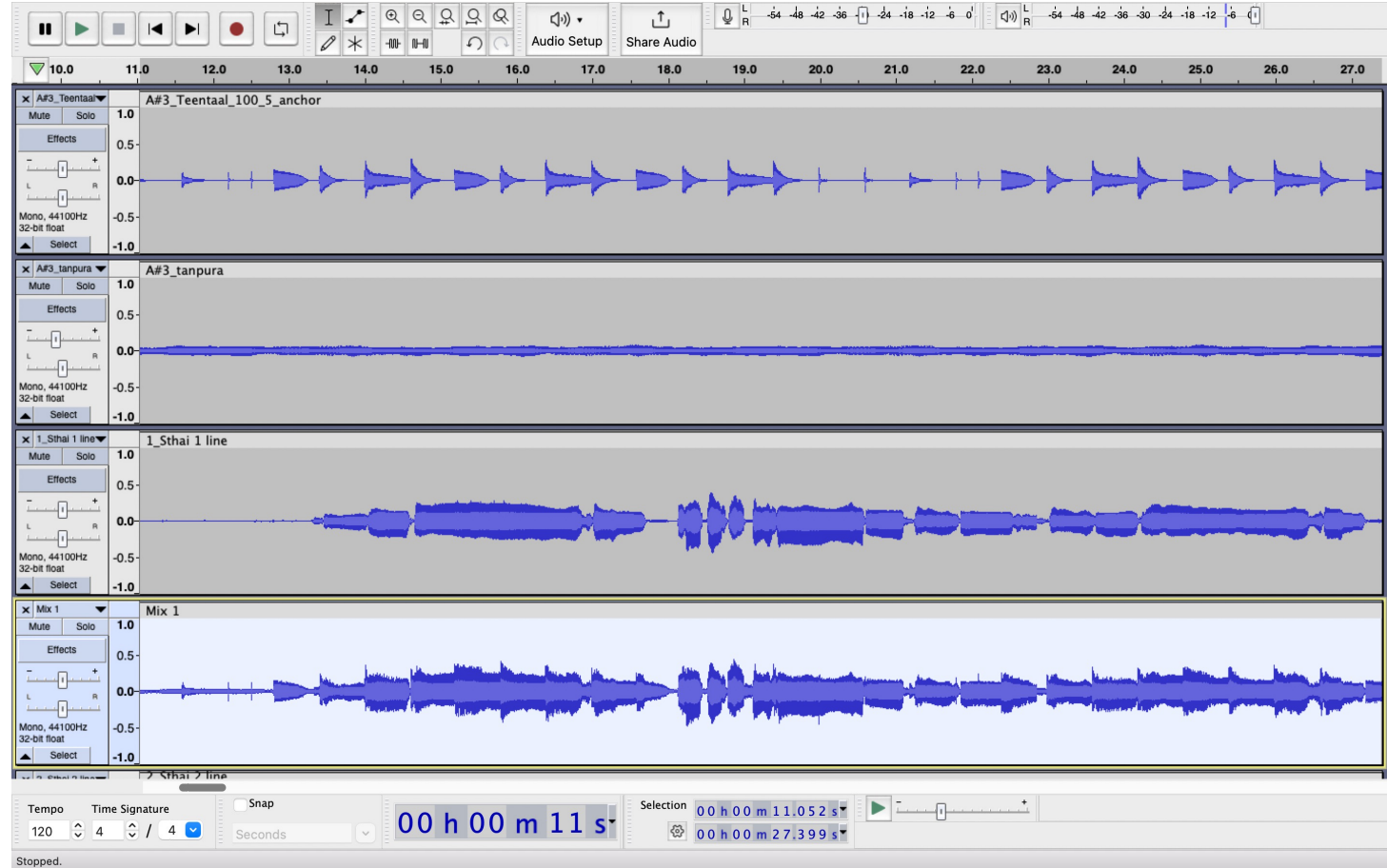http://hyperphysics.phy-astr.gsu.edu/hbase/Audio/mic.html

x[n]

# Sampling

- Nyquist Sampling theorem
- Humans can hear in the range 20Hz to 20kHz
- Popular: 44.1kHz for CD recordings

# Quantization

- Converting $x \in \mathbb{R}$ to a digital number
- Q bits per sample => $2^Q$ possible integer values per sample
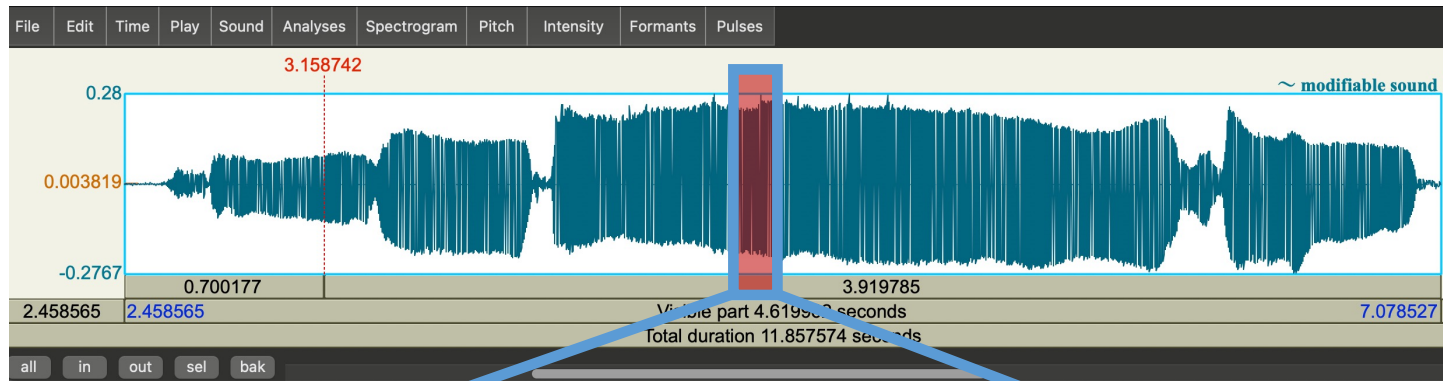- Popular: 16 bits per sample for CD recordings

# Waveforms



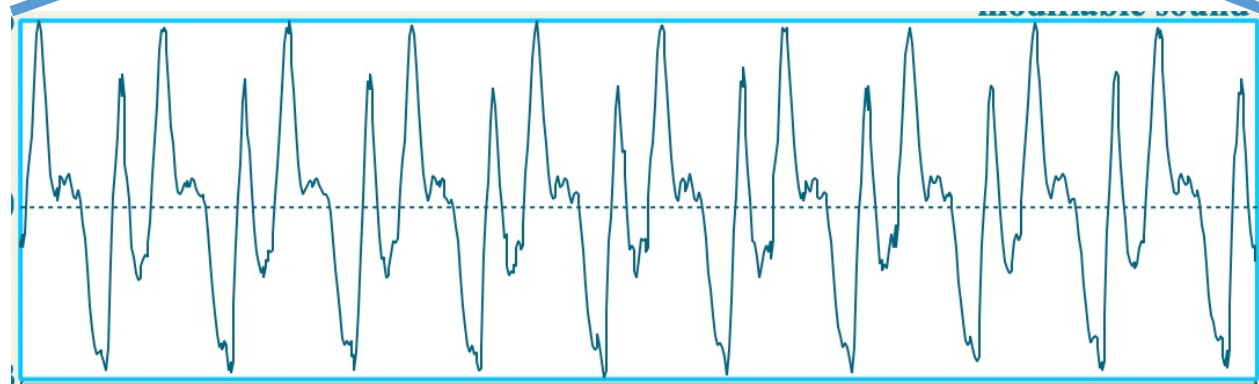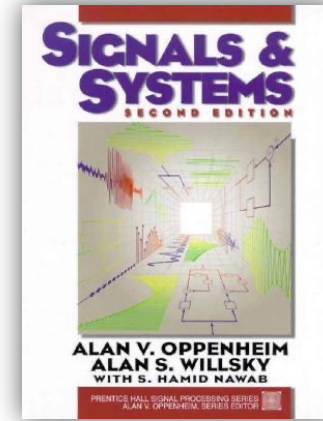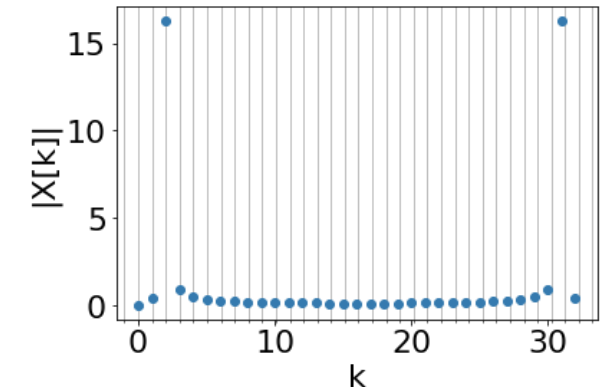Audacity

# Processing … we need mathematical model



Praat

# Fourier Transform

$$X[k] = \sum_{n=0}^{N-1} x[n] e^{-j\frac{2\pi}{N}kn} \; ; n = 0, \ldots, N-1$$

x[n] = sin(2$\pi$ * 2/32 * n)

x[n] = 0.5 * sin(2$\pi$ * 2/32 * n)
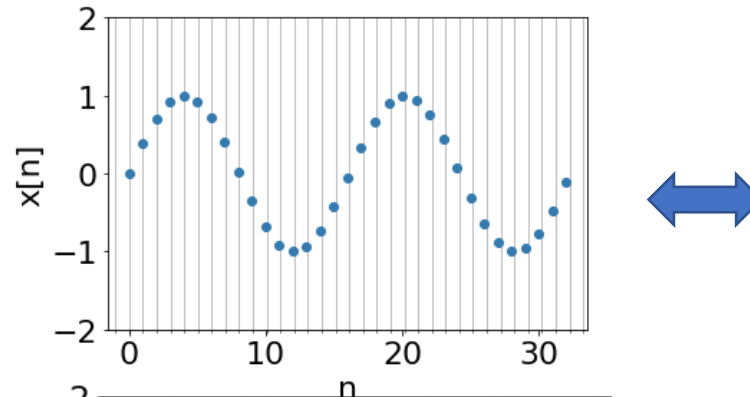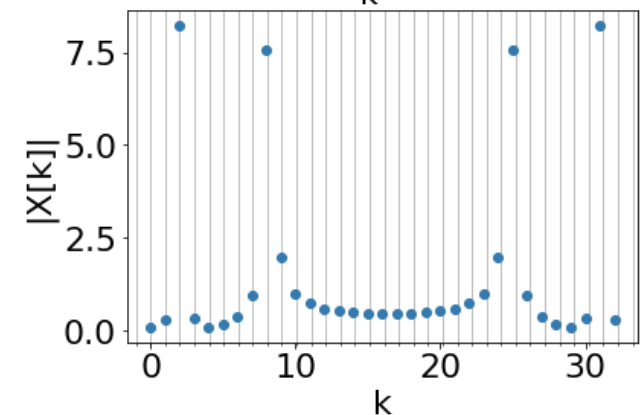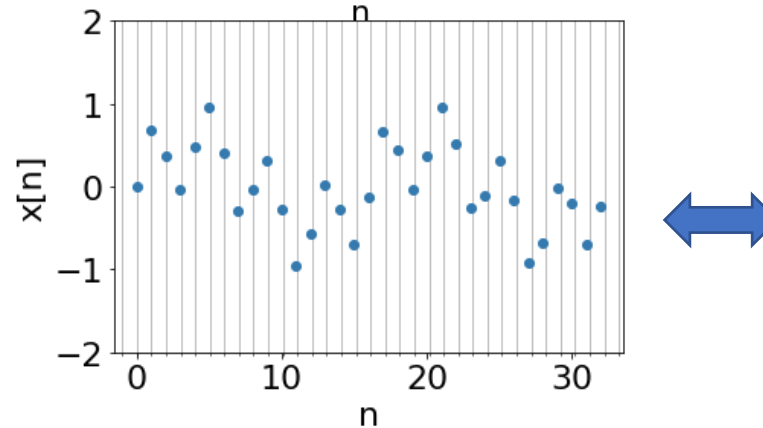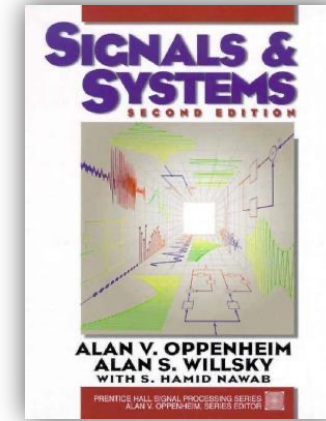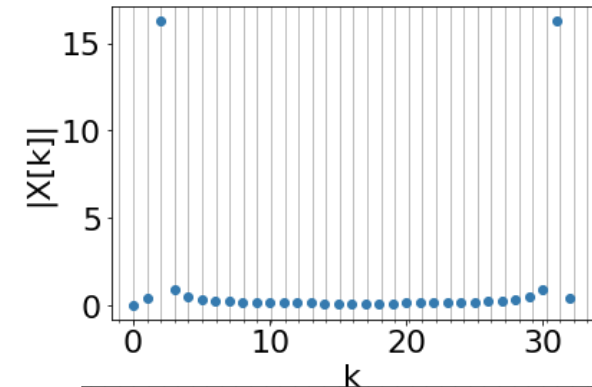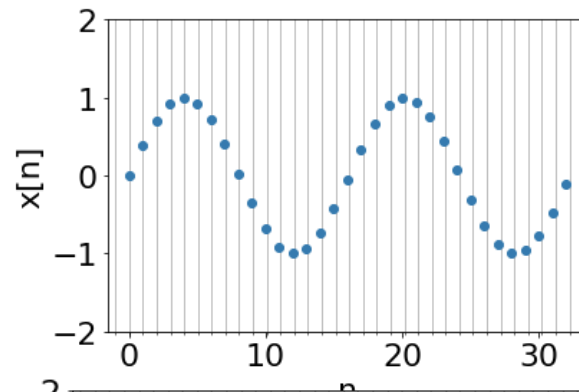+ 0.5 * sin(2$\pi$ * 8/32 * n)

# Fourier Transform
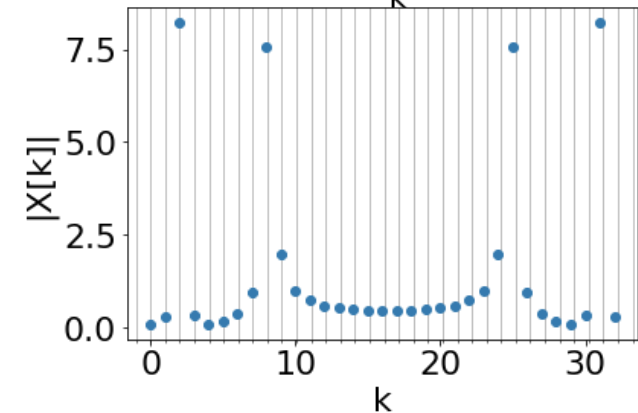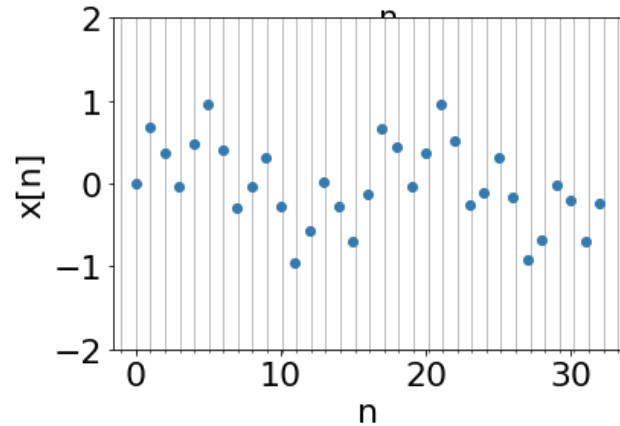
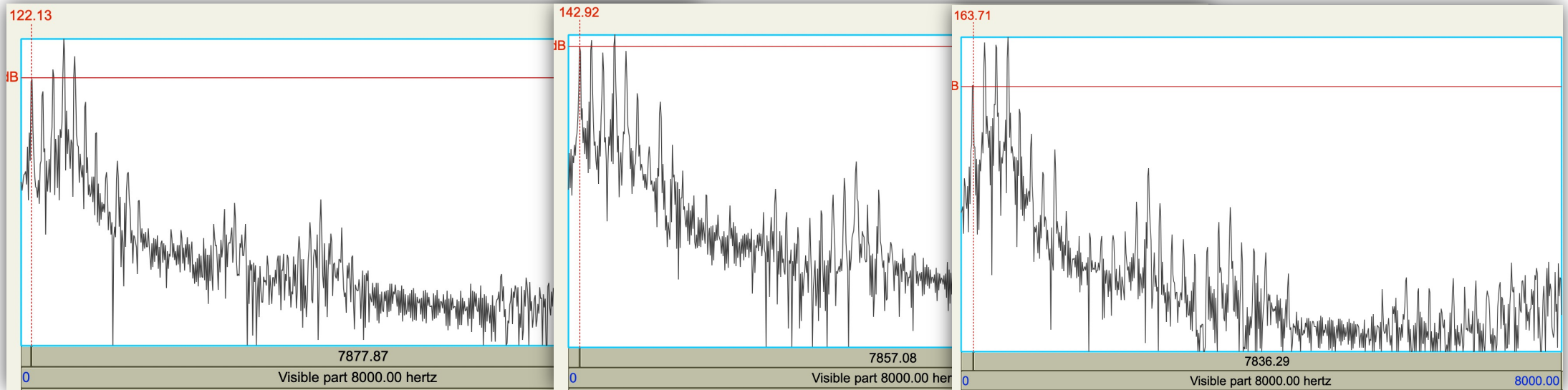$$X[k] = \sum_{n=0}^{N-1} x[n]e^{-j\frac{2\pi}{N}kn} \; ; n = 0, \ldots, N-1$$



x[n] = sin(2$\pi$ * 2/32 * n)

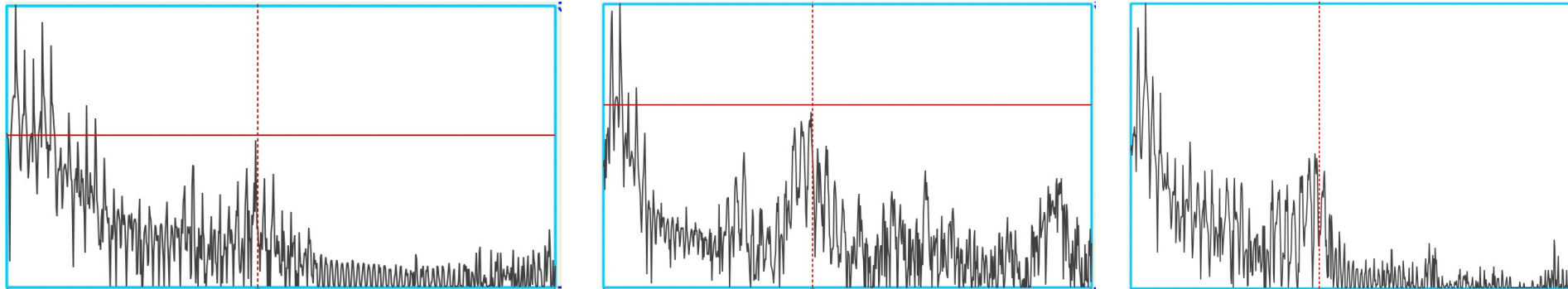x[n] = 0.5 * sin(2$\pi$ * 2/32 * n)
+ 0.5 * sin(2$\pi$ * 8/32 * n)

# Varying the Pitch

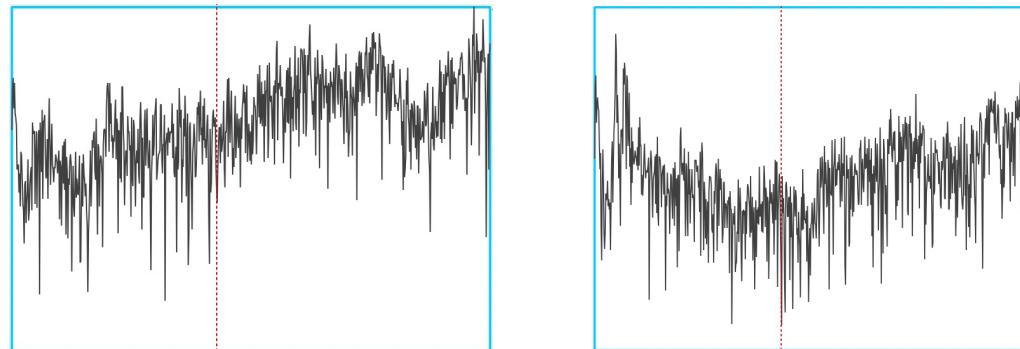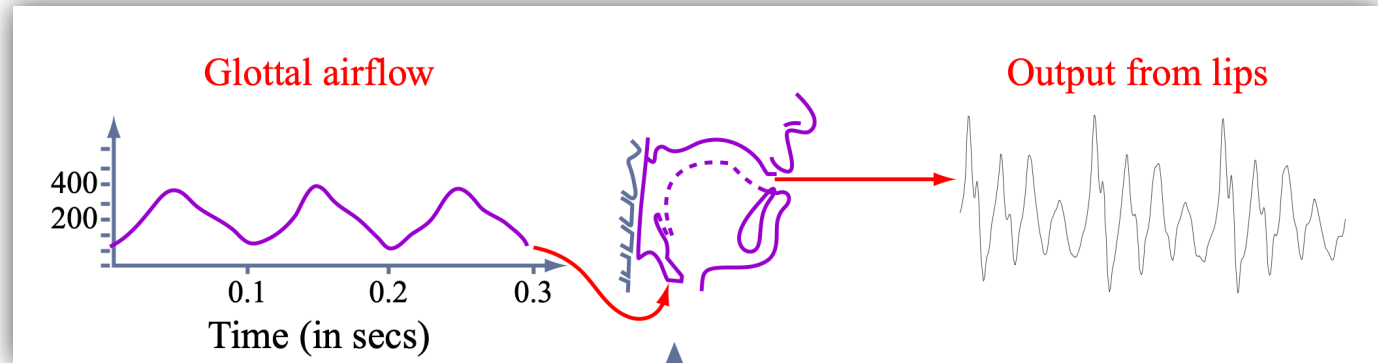# Spectra of speech

- Sounds with periodic waveforms: /a/, /i/, /m/



- Sounds with aperiodic waveforms: /s/, /f/


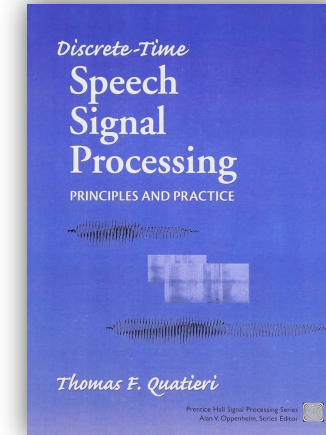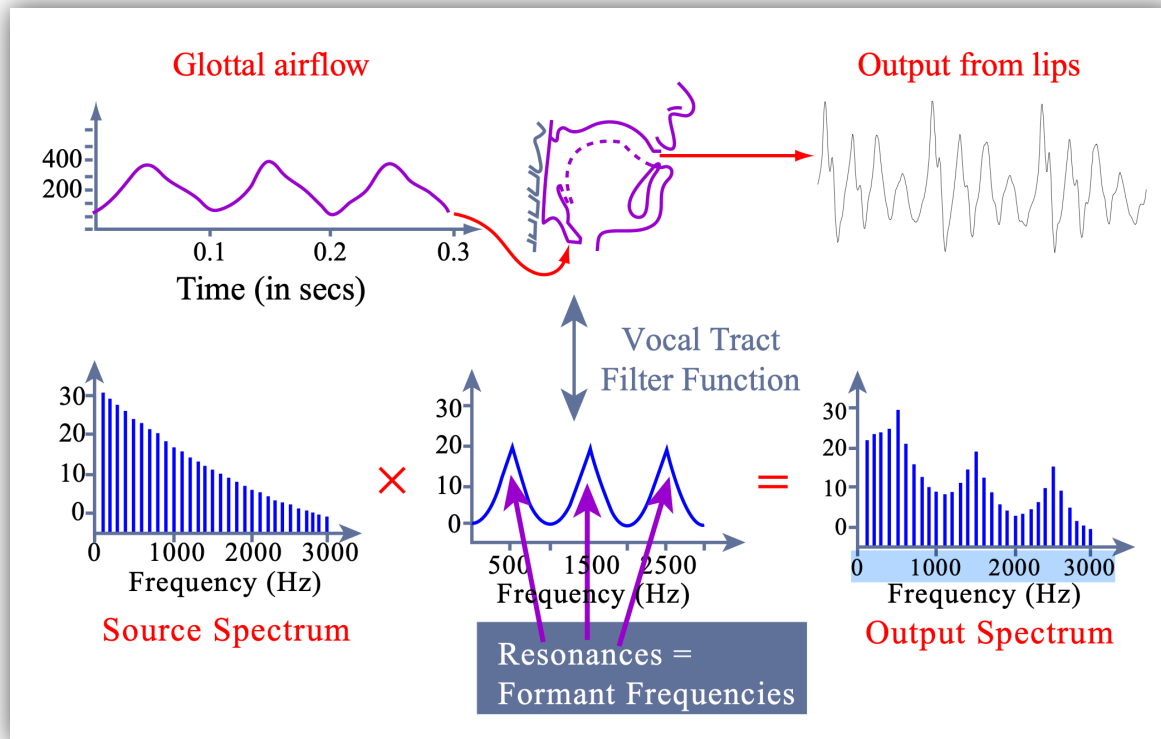
Praat

# Filter Theory



$$y[n] = x[n] * h[n] = \sum_{k=-\infty}^{\infty} x[k]h[n-k]$$

$$Y[k] = X[k]H[k]$$

# Source Filter Model

- Linear Time-Invariant (LTI) filters

# How do I make real applications with this?
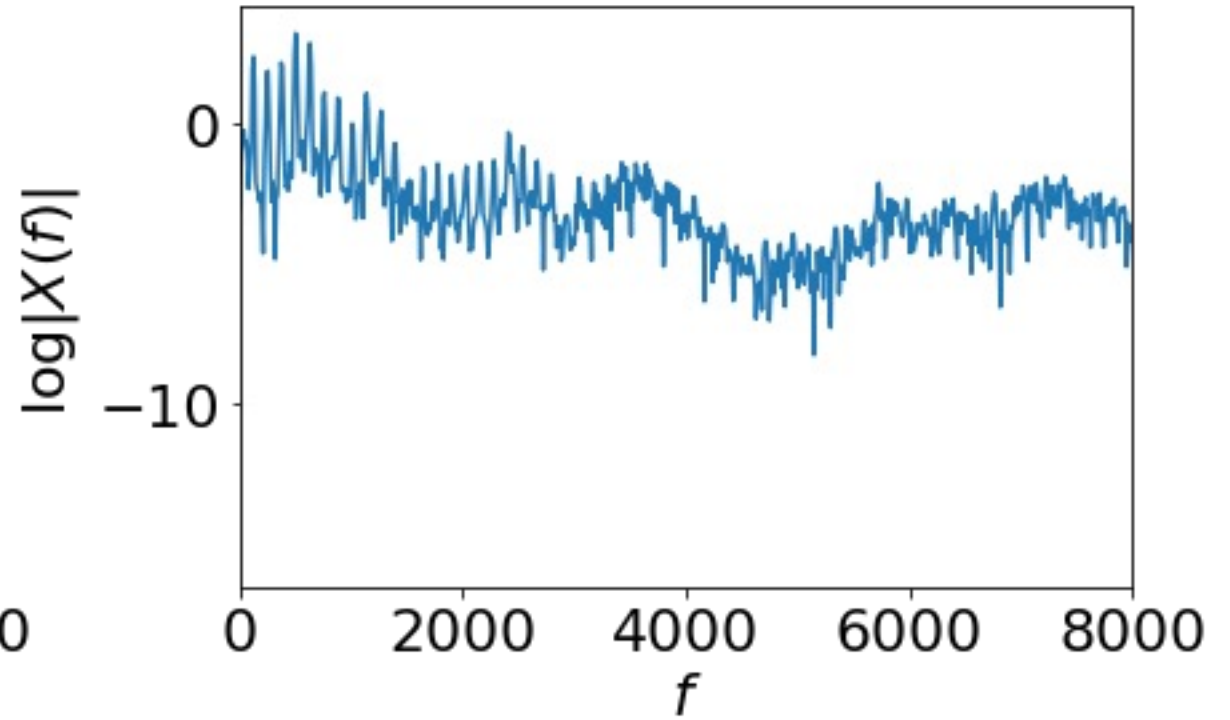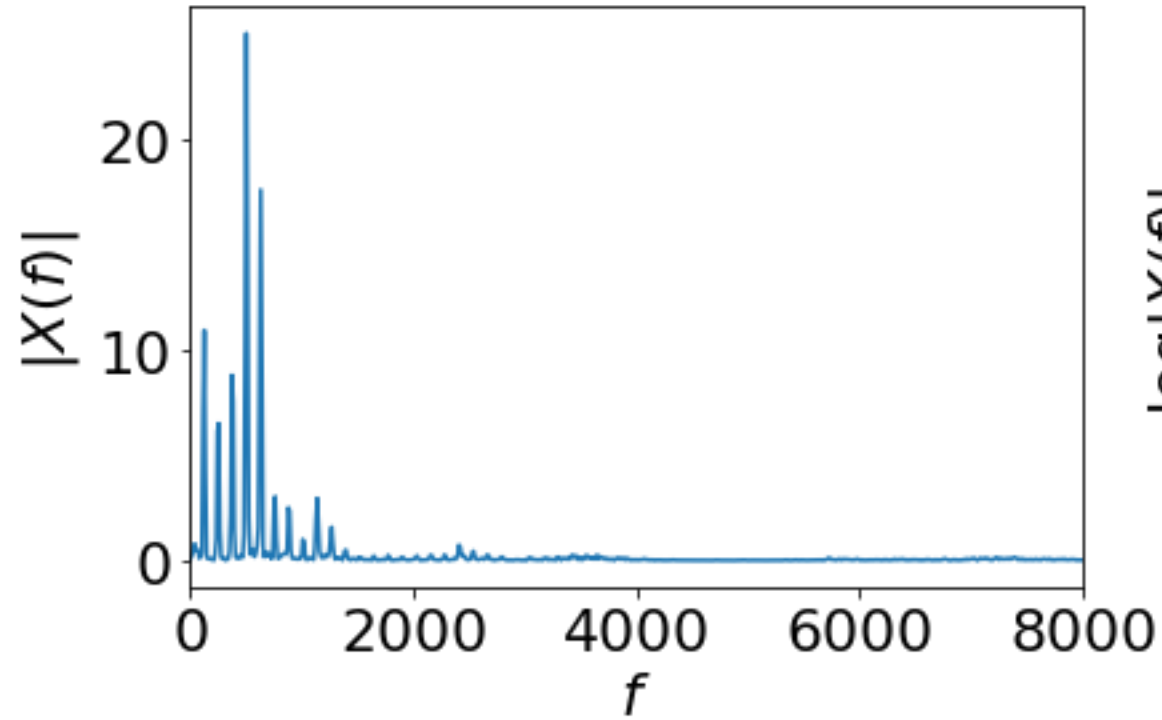
# Designing Representations

Representations should be

- minimal in size

- distinguishing for what we are interested in

- invariant to what we are not interested in


- Design the space so it may have uniform sensitivity (more in Audio Retrieval hands-on by Anup)

# Designing Representations

- for Pitch
  - Look at the peaks of spectrum
- for instrument/phoneme
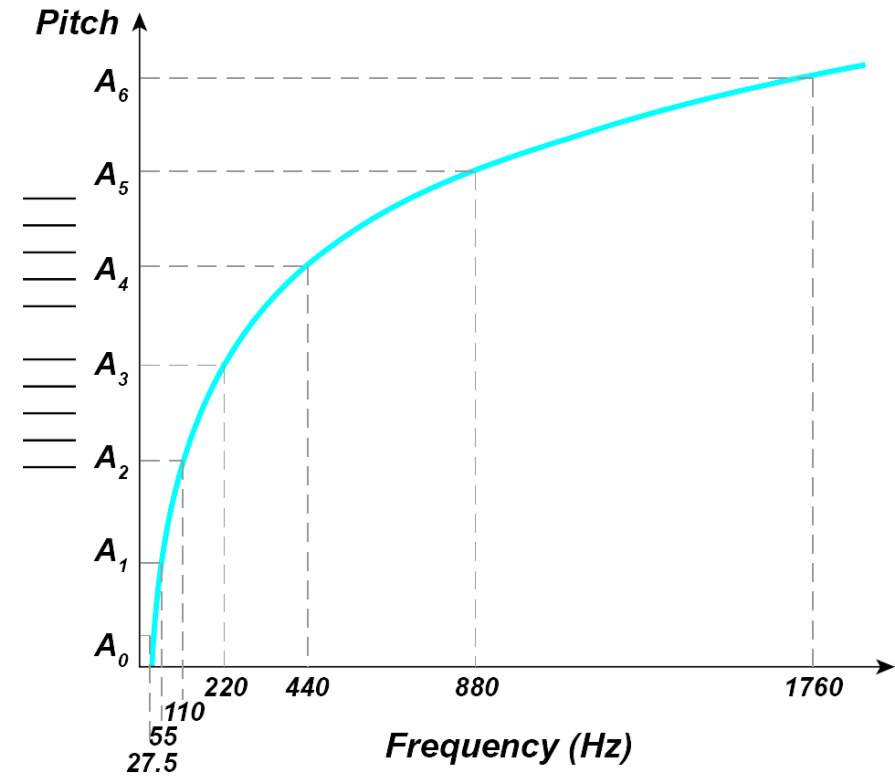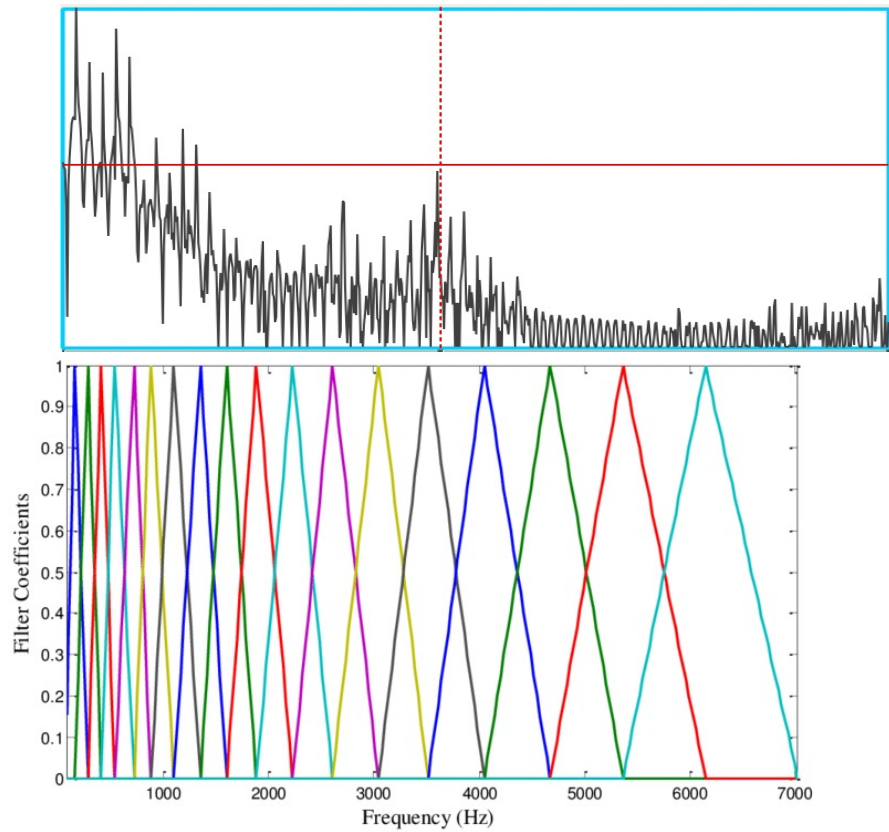  - Look at the spectral envelope

# Amplitude



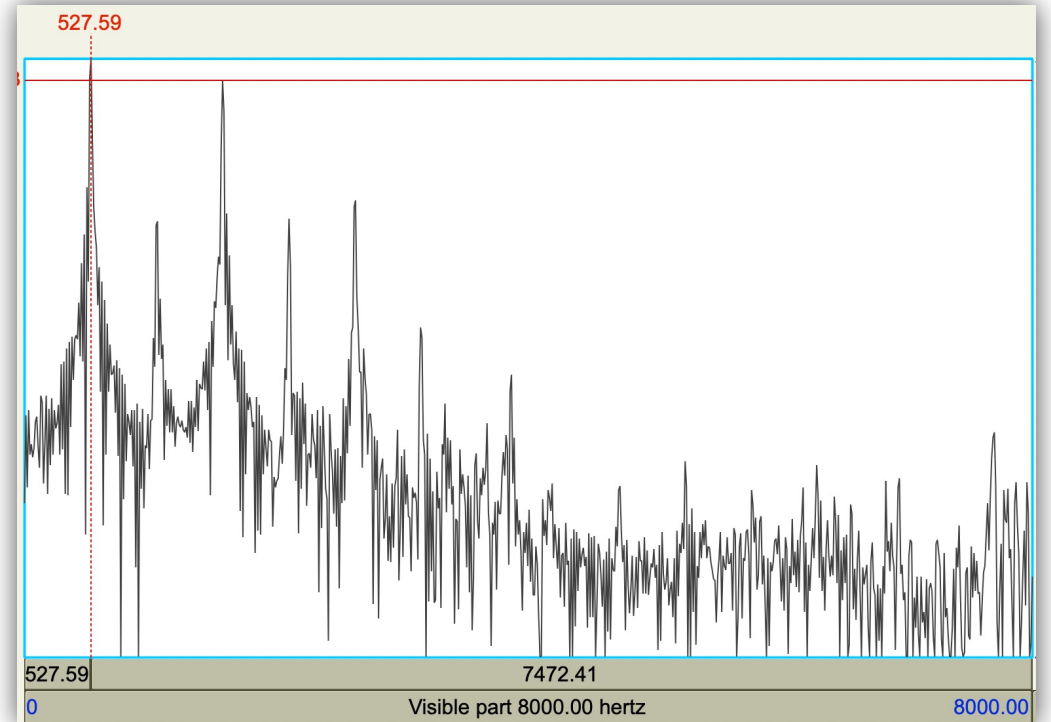$$A_{dB} = 20 \log_{10} A$$

# Frequency

- $\tilde{f} \propto \log f$

# Spectral Envelope

- $|X[\tilde{f}]|_{dB}$
- Mel-frequency, dB amplitude
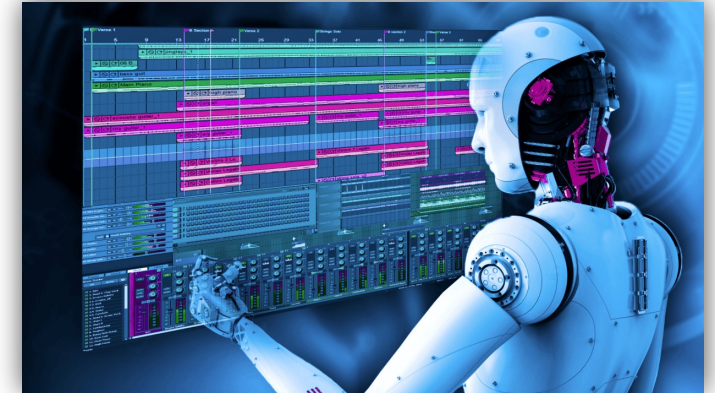- Take low frequency components of Fourier transform (DCT) of $|X[\tilde{f}]|_{dB}$
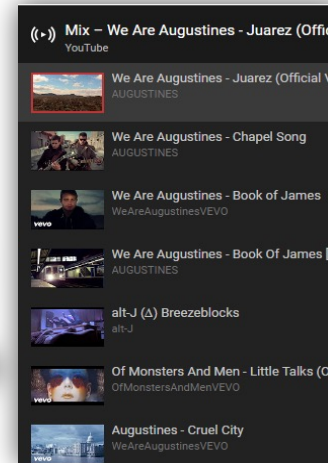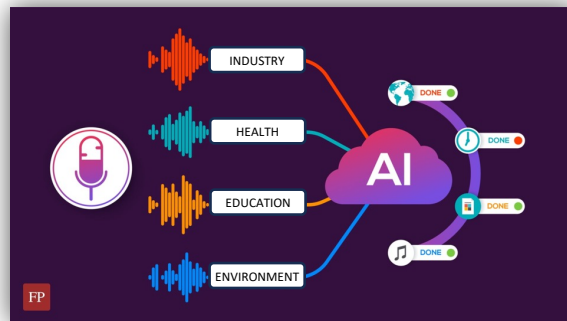
# You are ready!!!

# Not yet

- Dynamic behavior
- Time Series Analysis

# Short Time Fourier Transform



$$X[k, m] = \sum_{n=0}^{N-1} x[n]w[n - mH]e^{-j\frac{2\pi}{N}kn}; \quad k = 0, 1, ..., N-1$$
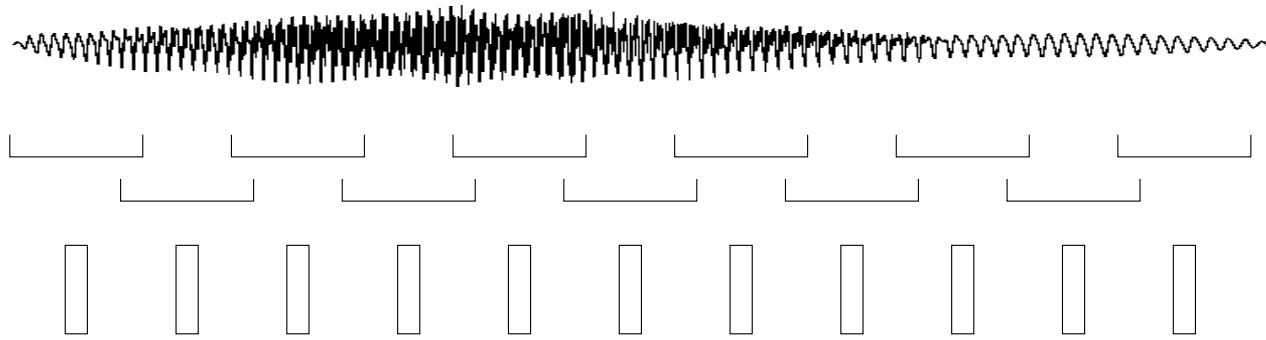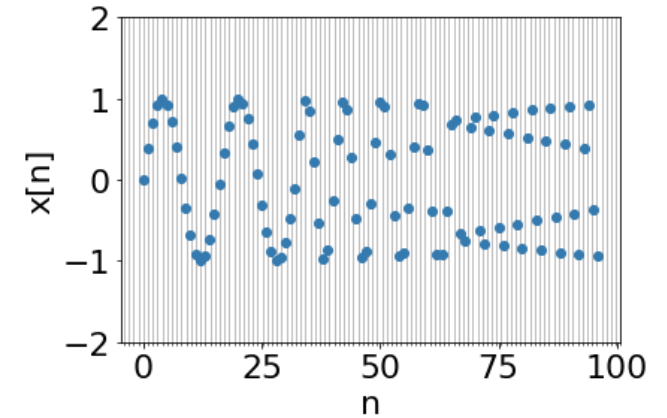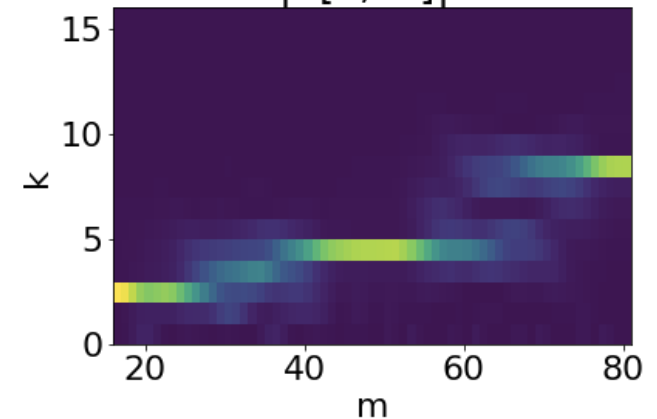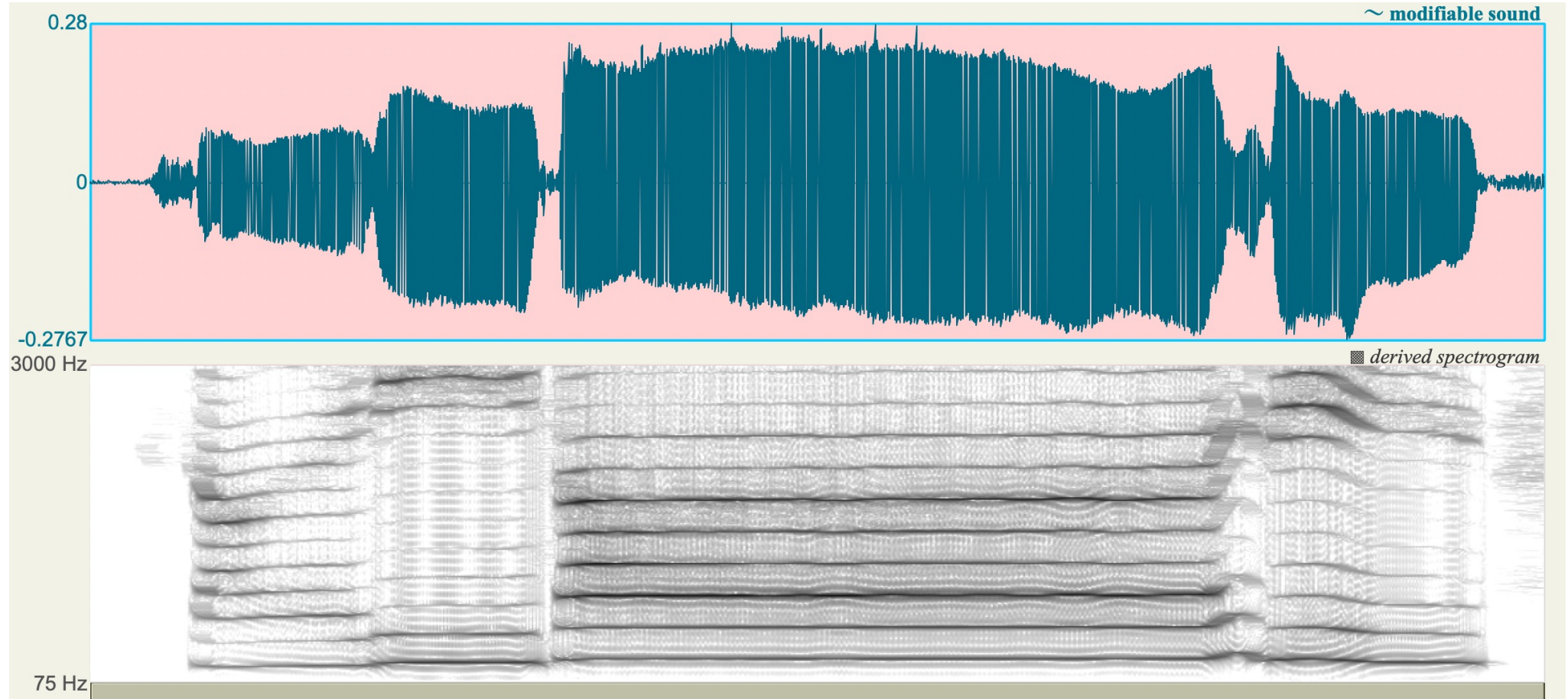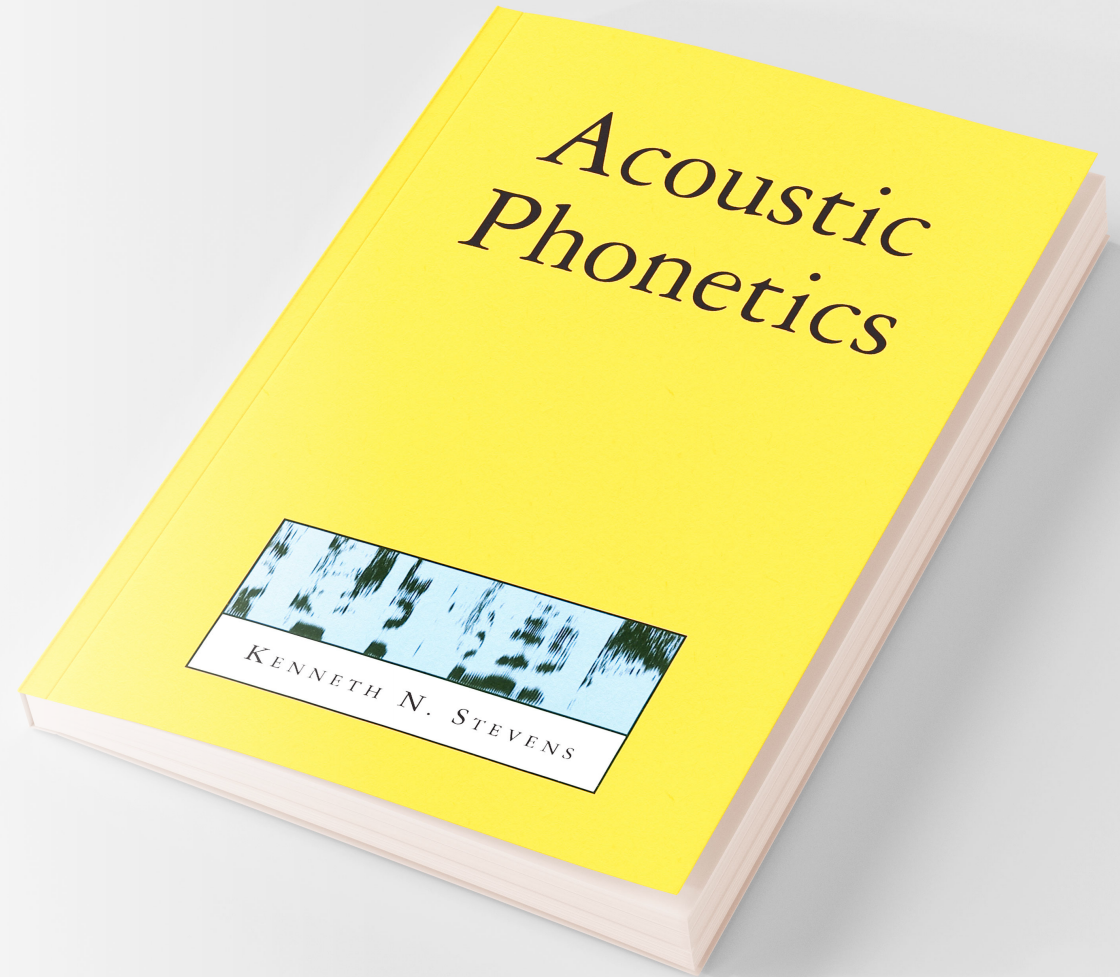
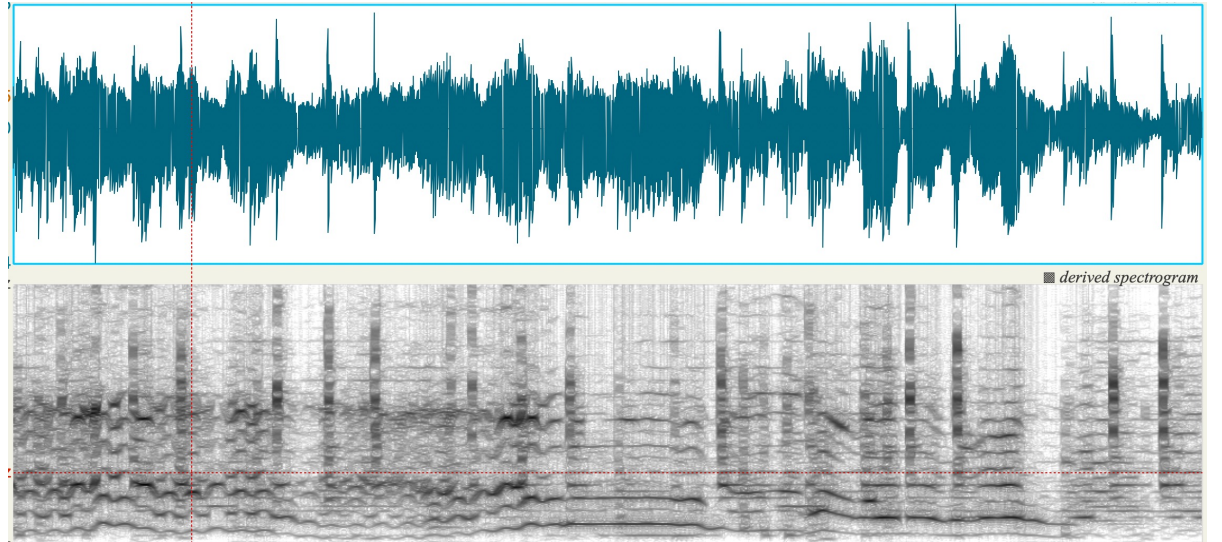# Short Time Fourier Transform
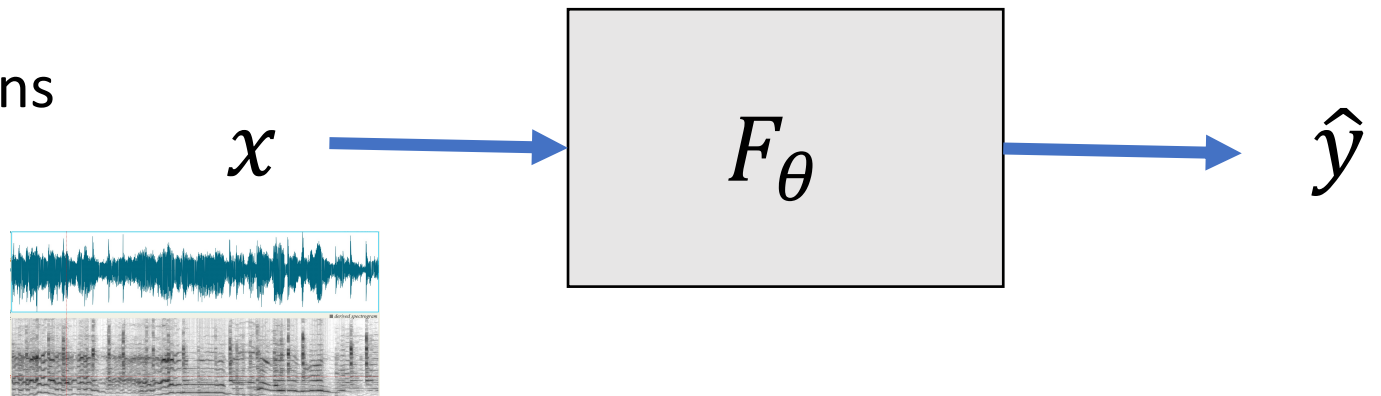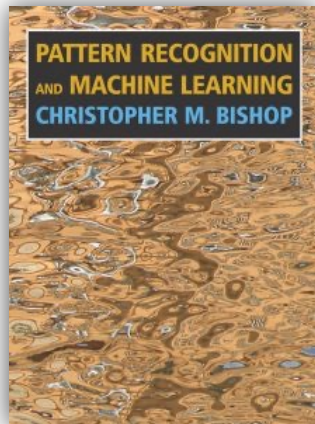
# It is possible

But only in ideal situations

# Real-world Variations

- Context (co-articulation)
- Running speech
- Speakers, instruments
- Languages
- Recording equipment
- Acoustic conditions



derived spectrogram
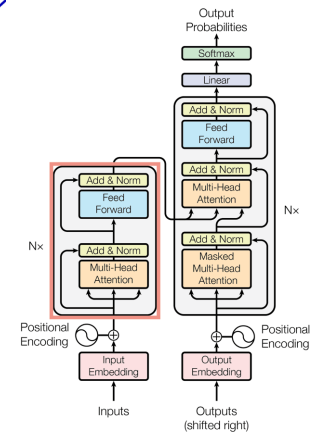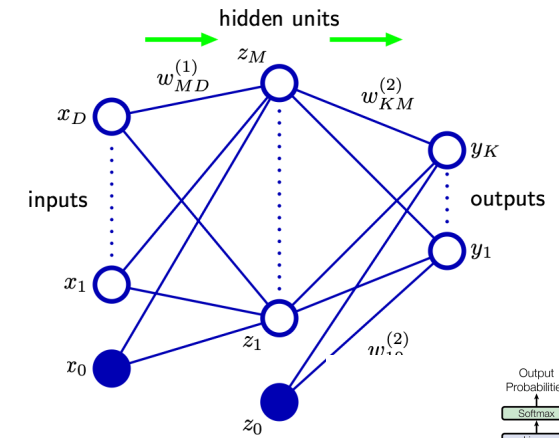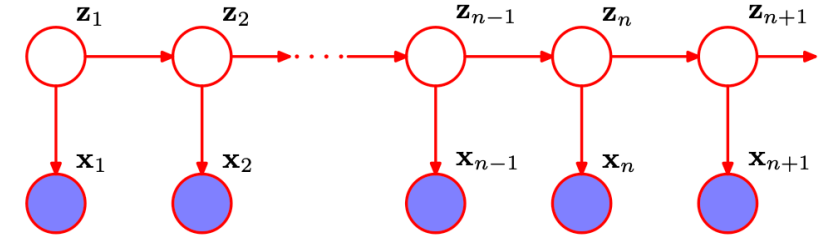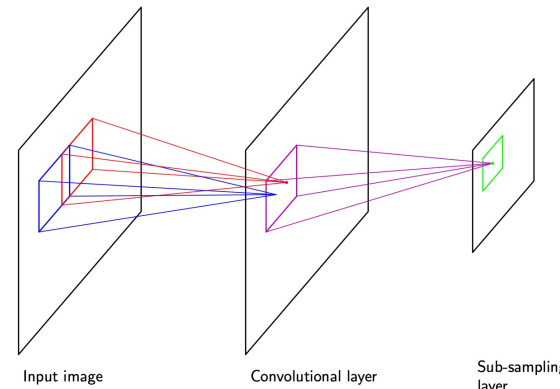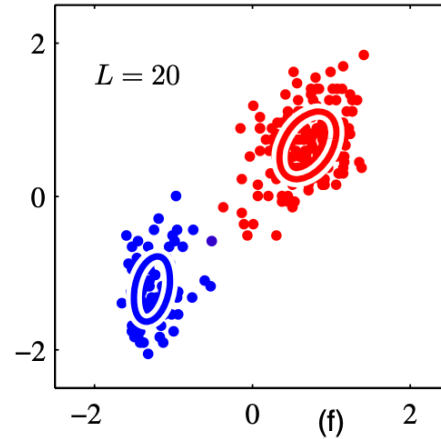
# Machine Learning

- Parametric models to learn the feature transformations
- Learn the mappings from
  - speech to text
  - audio to audio
  - audio to labels/classes
  - audio to recommendations

$$x \longrightarrow \boxed{F_\theta} \longrightarrow \hat{y}$$

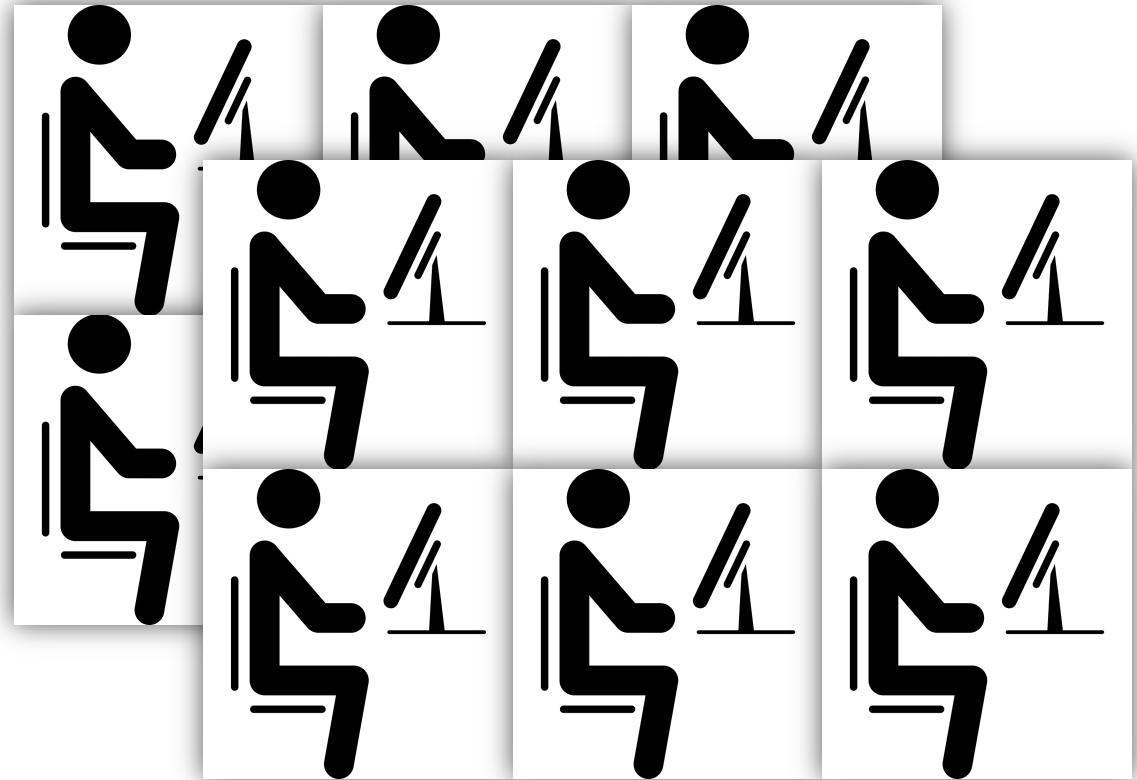$$\theta = \mathrm{argmin}_\theta \; \mathcal{L}(y, \hat{y}; \theta)$$

# Supervised Learning

- Gaussian Mixture Model
- Hidden Markov Model
- Multi-Layer Perceptron
- Support Vector Machine
- Convolutional Neural Network
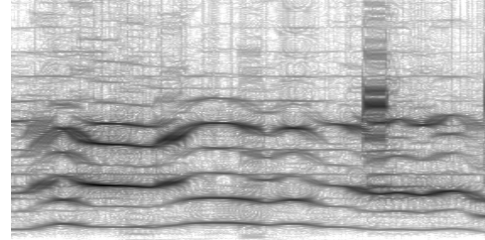- Recurrent Neural Network
- Transformers

Source: PRML Bishop and
https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf

# Bottleneck $(x, y)$

# Self-supervised Learning



Use $F_\theta$ instead of hand designed representations!

$$x \rightarrow \boxed{F_\theta \,\big|\, C_\phi} \rightarrow \hat{y}$$

$y =$