# Audio **&** Speech **Tech**nology

**[1] Background**

---

# **S**ound/**A**udio/**V**oice/**S**peech **?**

## **S**ound
Mechanical wave that is an oscillation of pressure transmitted through a solid, liquid, or gas, composed of frequencies within the range of hearing.

## **A**udio
Audible sound coming from a recording, transmission or electronic device.
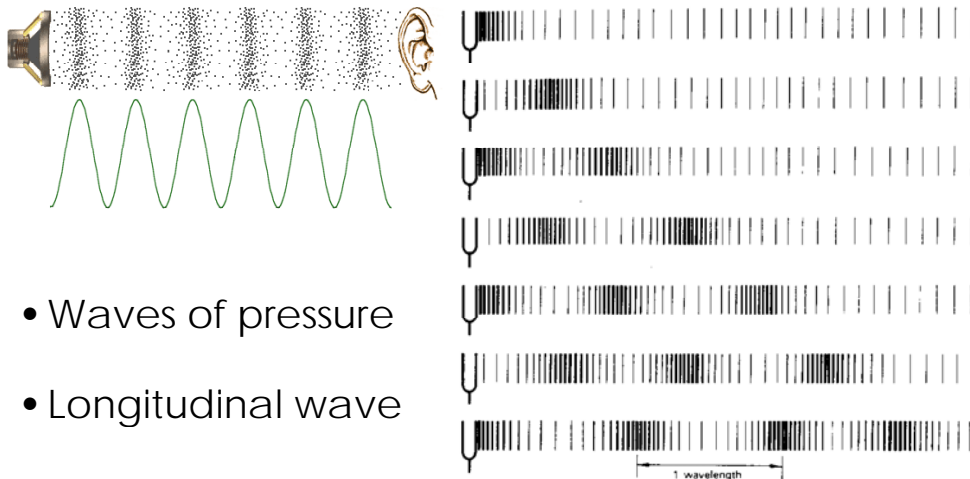
## **V**oice
The sound produced by the vocal organs of a vertebrate, especially a human.

## **S**peech
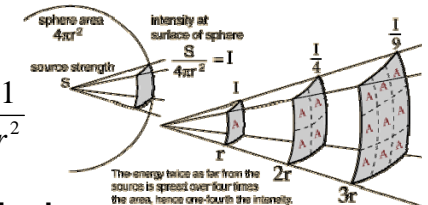Vocal communication which is an expression thoughts, feelings, and ideas orally.

---

# Physics of **Sound**



- Waves of pressure
- Longitudinal wave

http://www.mediacollege.com/audio/images/loudspeaker-waveform.gif

http://personal.cityu.edu.hk/~bsapplec/Fire/Image222.gif

---

# **Sound** Intensity

**S**ound Intensity

$$I = \frac{Sound\ Power}{Area} = \frac{Sound\ Power}{4\pi r^2} \propto \frac{1}{r^2}$$



**S**ound Intensity Level in Decibels

$$I_{dB} = 10\log_{10}\left[\frac{I}{I_0}\right]$$

**S**tandard threshold of hearing intensity

$$I_0 = 10^{-12}\ watt/m^2$$

http://toonz.ca/bose/wiki/images/1/1e/IntensitySurfaceSphere.gif
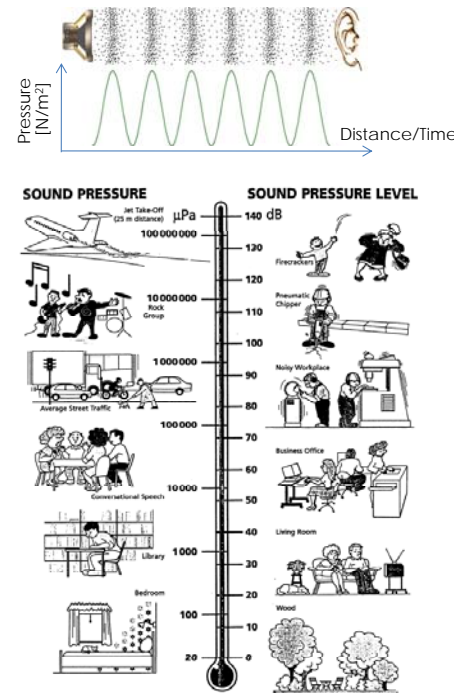
# Sound Pressure

Sound Pressure

$$p = \text{rms of pressure wave} \propto \frac{1}{r}$$

Sound Pressure Level (SPL)
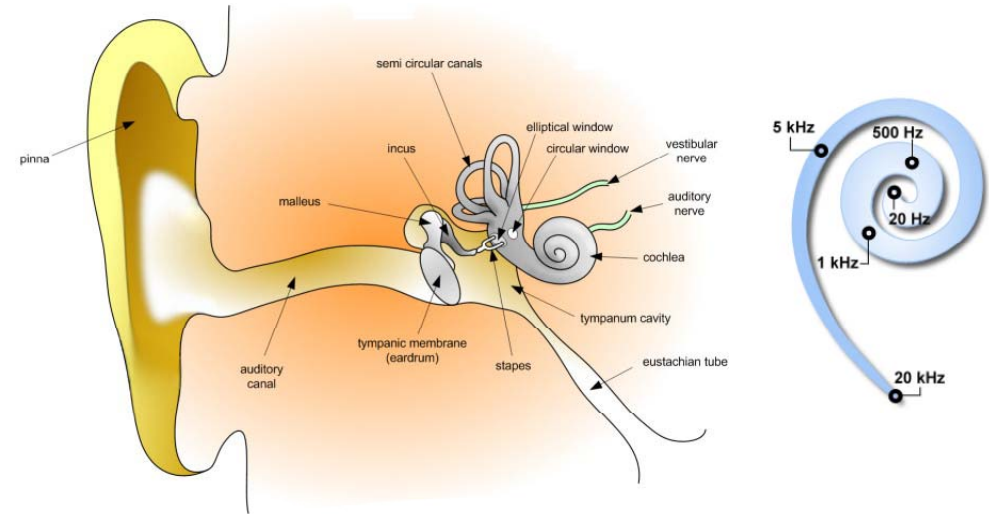
$$P_{dB} = 20\log_{10}\left[\frac{p}{p_0}\right]$$

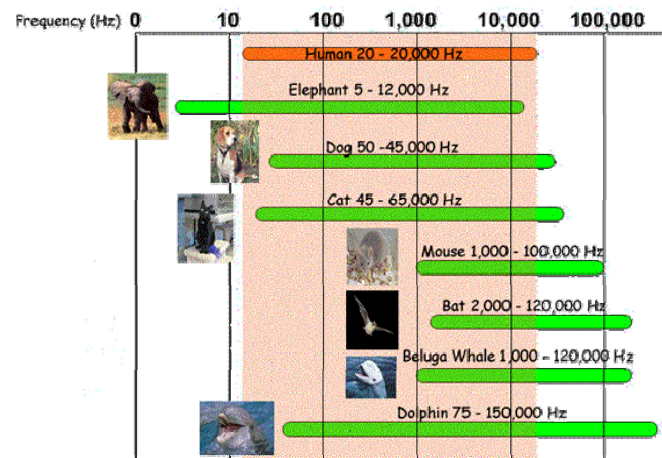Standard Threshold of Hearing Pressure
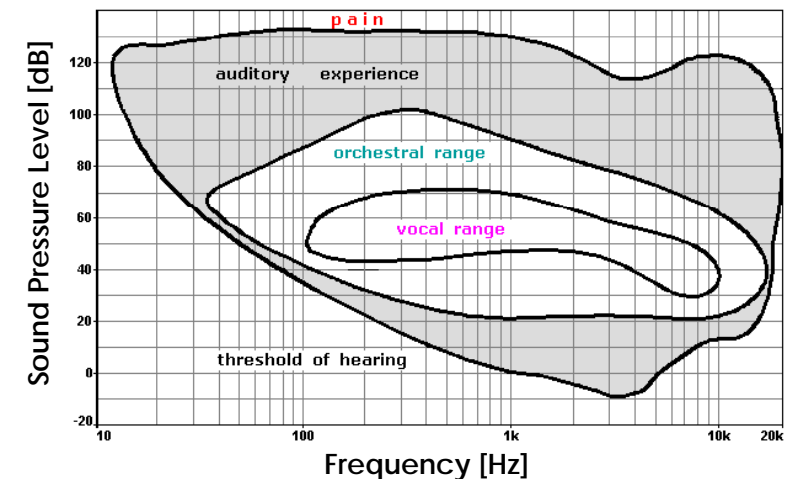
$$p_0 = 2\times10^{-5}\ N/m^2$$

Pressure [N/m²] — Distance/Time

SOUND PRESSURE — SOUND PRESSURE LEVEL

*Brüel and Kjær*

---

# Human Hearing

http://hendrix2.uoregon.edu/~dlivelyb/phys152/images/HumanEar.jpg
http://cnx.org/content/m43048/latest/Picture%202.png

---

# Human Hearing

http://www.cyberphysics.co.uk/graphics/diagrams/hearing_range.gif

---

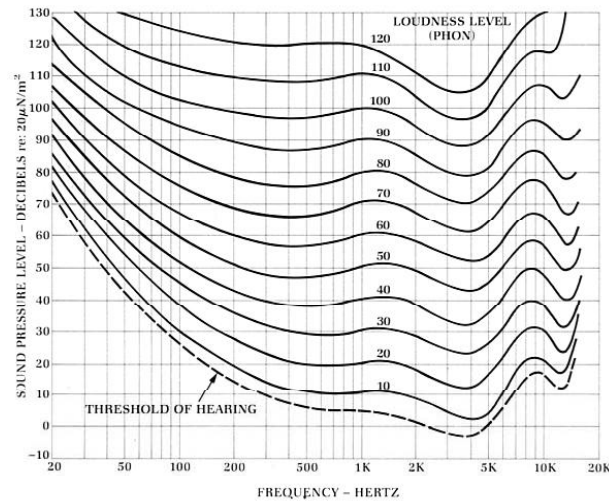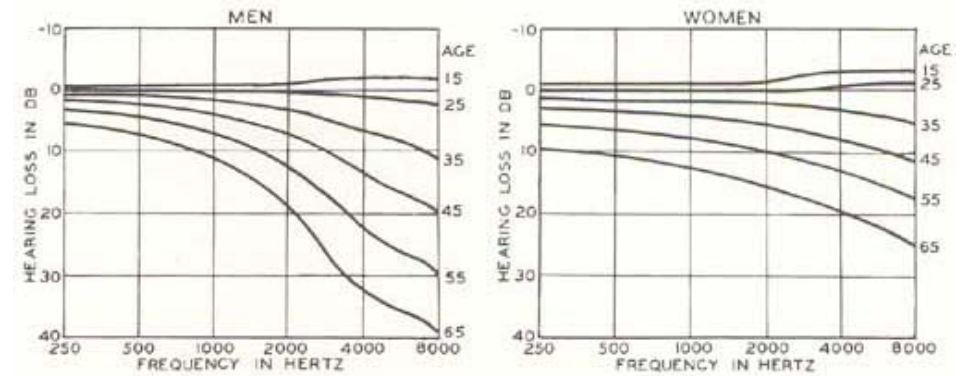# Human Hearing

*http://sound.westhost.com/articles/fadb.htm*

# Sound Loudness

- Subjective Term Based on Human Perception
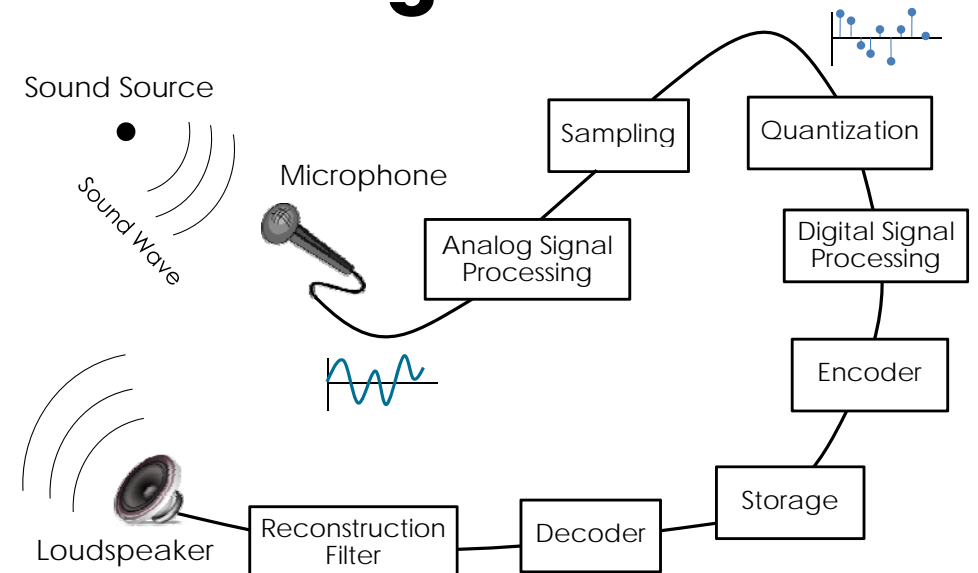
- Same Intensity

$$\neq$$

Same Loudness

# Hearing Age

# Auditory Masking



**Used in Lossy Compression**

# Audio Signal I/O

Sound Source

Microphone

Sound Wave

Sampling

Quantization

Analog Signal Processing

Digital Signal Processing

Encoder

Storage

Loudspeaker

Reconstruction Filter

Decoder

# Micro**phone**

**T**ransducer : Sound $\Rightarrow$ Electrical signal.

- **C**ondenser microphone
- **D**ynamic microphone
- **P**iezoelectric microphone
- **F**iber optic microphone
- **L**aser microphone
- **M**EMS microphone
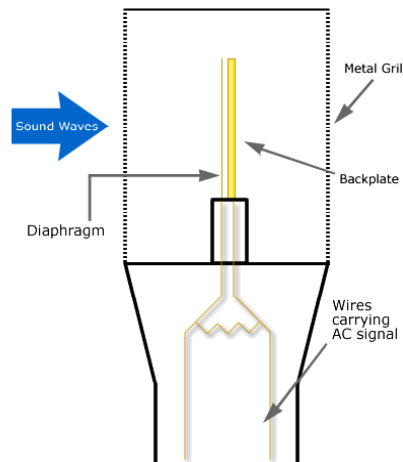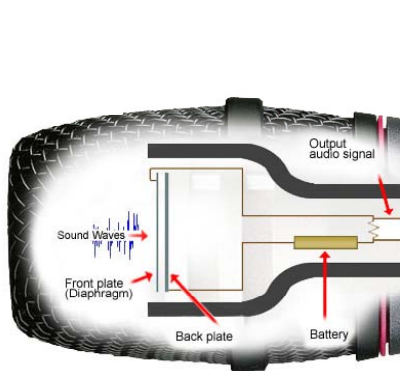  (MicroElectrical-Mechanical System)

# Dynamic
## Micro**phone**

- Electromagnetic Induction
- Inexpensive
- Robust & Resistant to Moisture
- High Gain
- Attenuate in High Frequencies



Microphone

Diaphragm

Sound Waves

Coil

Cone

Signal Out

Magnet

Wires carrying AC signal

Sound Waves

Magnet

Coil of wire (attached to diaphragm)

Diaphragm

# Condenser
## Micro**phone**

- Capacitive Microphone
- Wider & Flatter Frequency Response
- More Expensive
- Fragile
- Need Power



Sound Waves

Output audio signal

Front plate (Diaphragm)

Back plate

Battery

Metal Grill

Sound Waves

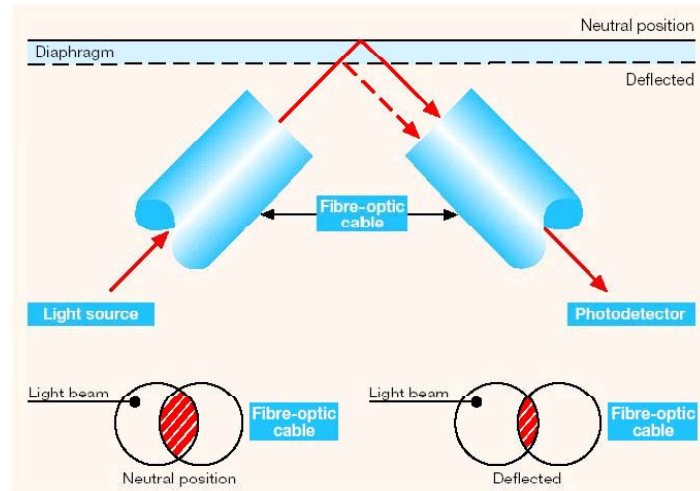Backplate

Diaphragm

Wires carrying AC signal

# Piezoelectric
## Micro**phone**

- Piezoelectricity : Pressure $\Rightarrow$ Voltage
- Contact Microphone
- Durable
- High Pressure/Temp Environments



sound pressure

diaphragm

piezoelectric crystal
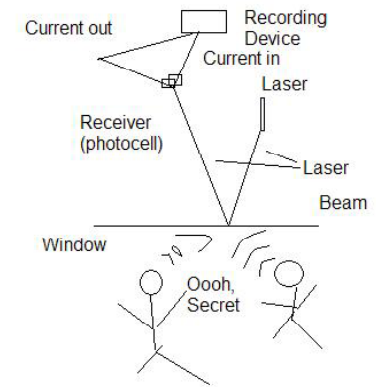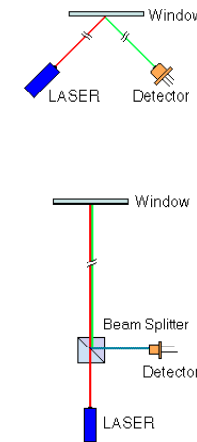
metal

# Fiber Optic
## Microphone

- Light Intensity
- High Dynamic and Frequency Range
- Robust, resistant to changes in heat and moisture, do not influence by any electrical, magnetic or radioactive fields



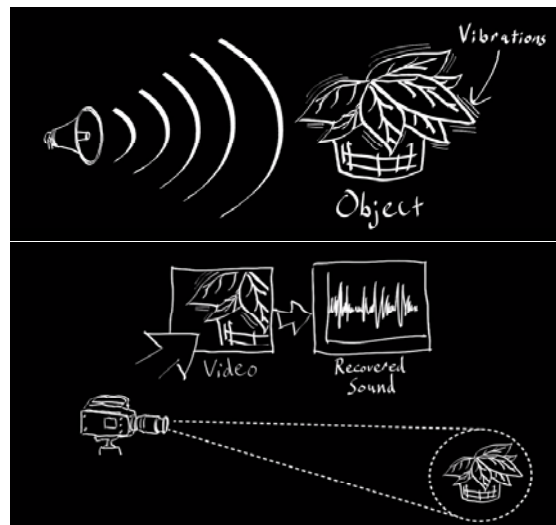http://www.prosoundweb.com/images/uploads/SennDiagram.jpg

# Laser
## Microphone

- Light Intensity
- Pick up Sound at a Distance



http://williamson-labs.com/laser-mic.htm
http://hackedgadgets.com/2007/08/16/the-laser-listener/
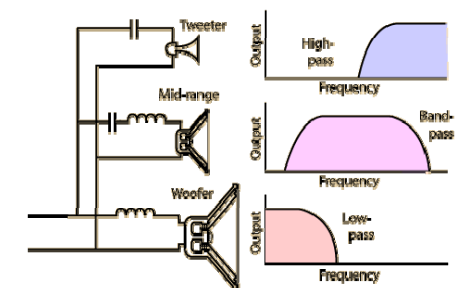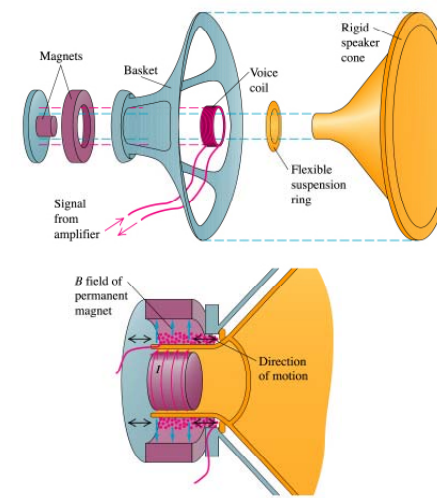
# Visual
## Microphone

- Using High Speed Video
- Image Processing
- Passive Recovery of Sound from Video

# Loudspeaker

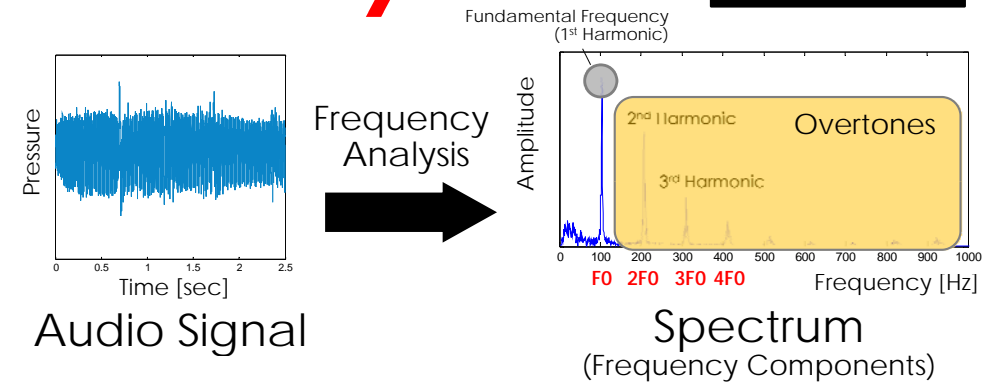## Transducer : Electrical signal $\Rightarrow$ Sound



Audio Crossover

http://wiki.vpa.mtu.edu/wiki/index.php/Speakers
http://hyperphysics.phy-astr.gsu.edu/hbase/audio/imgaud/cross6.gif

# Sampling

$s(t)$

$t$

Analog Audio Signal

$s[n]$

0 1 2 3 ⋯

$n$

Discrete Audio Data

To avoid Aliasing
Sampling Frequency ($F_S$) > 2 × Signal Max. Frequency

human can hear
~ 20 Hz-20 kHz

➡ CD : 44.1 kHz
Digital Audio Tape : 48 kHz
MP3 : 32 kHz, 44.1 kHz, 48 kHz, etc.

human speech
~ 5 Hz-4 kHz

➡ Telephone-quality Audios : 8 kHz

# Frequency Analysis

Harmonic Structure

Fundamental Frequency
(1st Harmonic)

Pressure

Time [sec]

Audio Signal

Frequency
Analysis

Amplitude

2nd Harmonic     Overtones

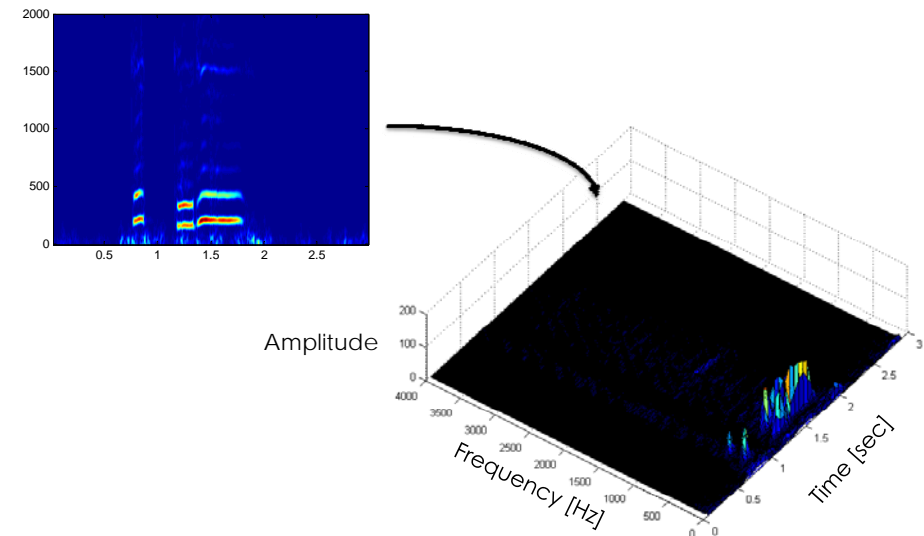3rd Harmonic

F0  2F0  3F0  4F0     Frequency [Hz]

Spectrum
(Frequency Components)

- Continuous System ⇒ Fourier Transform
- Digital System ⇒ Discrete Fourier Transform [DFT]
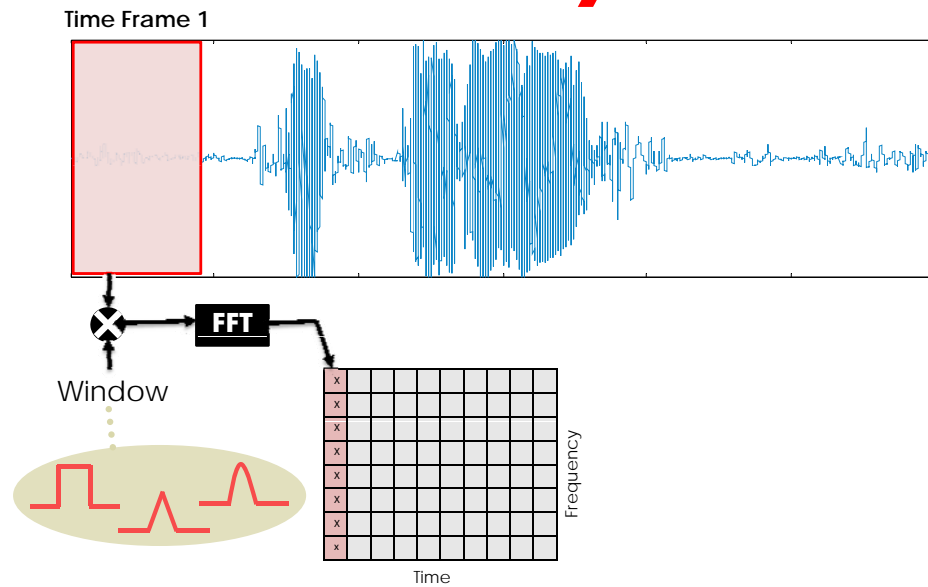- Fast Fourier Transform [FFT] = Fast Algo. for DFT

# Time-Frequency Analysis

Pressure

Time [sec]

Audio Signal
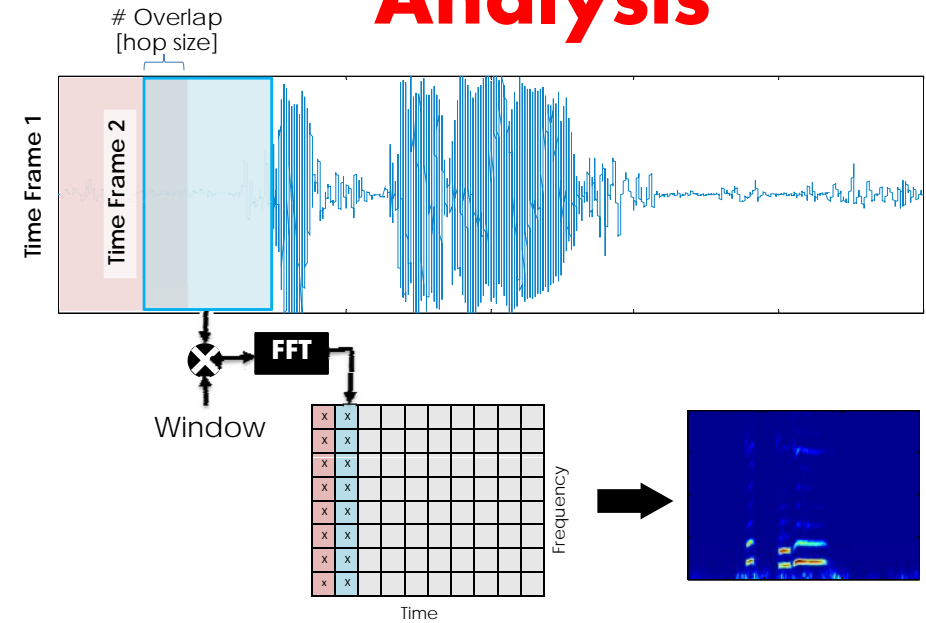
T-F
Analysis

Frequency [Hz]

Time [sec]

Spectrogram
(Frequency Components over Time)

- Short-Time Fourier Transform

# Time-Frequency Analysis

Amplitude

Frequency [Hz]

Time [sec]

# Time-Frequency Analysis

**Time Frame 1**



Window

FFT

Frequency

Time

# Time-Frequency Analysis

# Overlap [hop size]

Time Frame 1

Time Frame 2

Window

FFT

Frequency

Time

# Spectrogram

**White Noise** 🔊

**Speech** 🔊

Frequency

Time

Frequency

Time

# Spectrogram

**Pitched Sounds** 🔊

**Unpitched Sounds** 🔊

Frequency

Time

Frequency

Time

# Spectrogram

### 8-bit Digital Sound 🔊

Frequency | Time

### Saxophone 🔊

Frequency | Time

# Spectrogram

### Piano 🔊

Frequency | Time

### Guitar 🔊

Frequency | Time

# Spectrogram

**High Vibrato & Tremolo**

### Violin 🔊

Frequency | Time

# Spectrogram

### Thai Flute 🔊

Frequency | Time

# Spectrogram

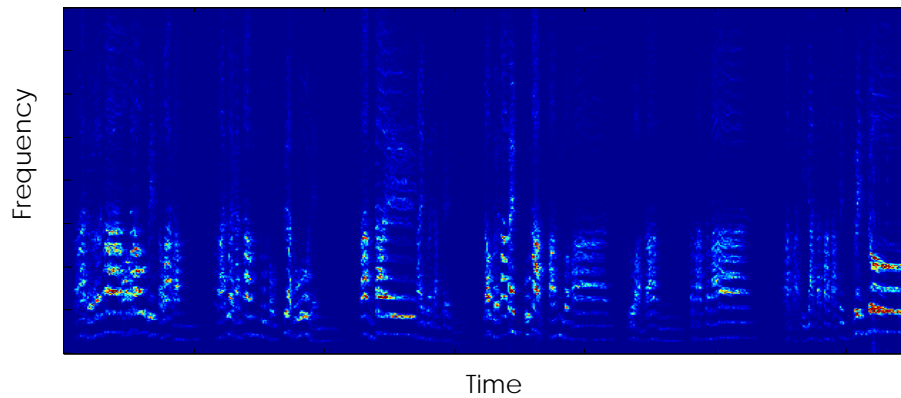## Hand Drum 🔊

# Spectrogram

## Guitar + Tapping 🔊

# Spectrogram

## Sing-along 🔊

# Spectrogram

## Polyphonic 🔊

# Frequency Selective Filter

**Noisy Signal** 🔊



Low Frequency Noise

# Frequency Selective Filter

**High-pass Filter**

# Frequency Selective Filter

**Restored Signal** 🔊



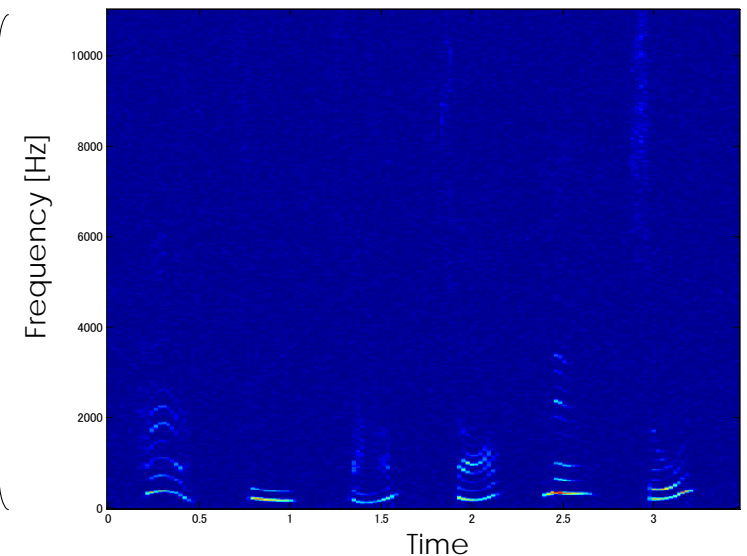Low frequency noise was removed
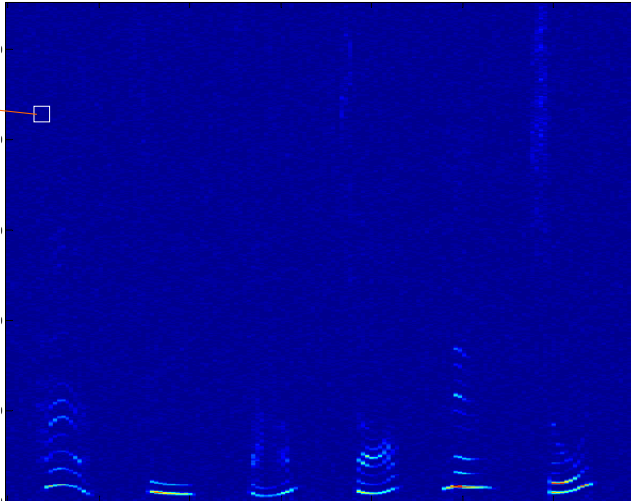
# Thresholding

**Noisy Signal** 🔊
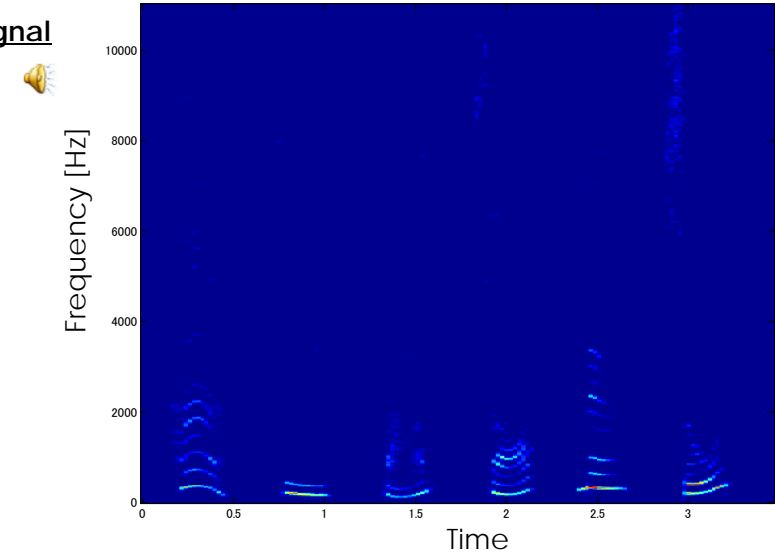
Additive White Gaussian Noise

# Thresholding

Find the power of signal $\sum |s|^2$ in each block

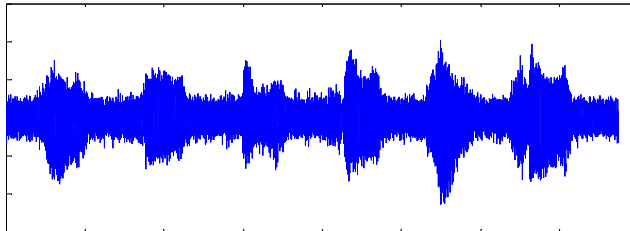If the power is less than threshold, the block is truncated to zero.
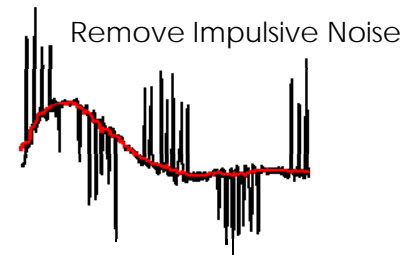
# Thresholding

**Restored Signal**

Frequency [Hz] vs Time

# Thresholding

**Noisy Signal**

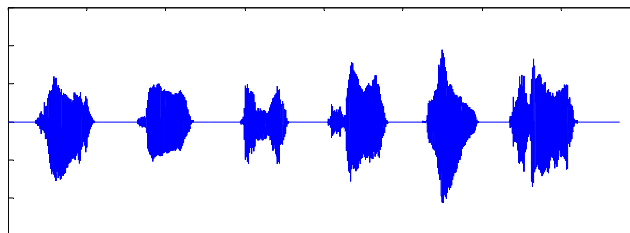**Restored Signal**

Remove Impulsive Noise
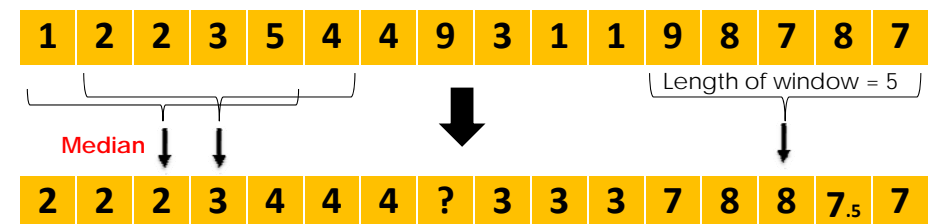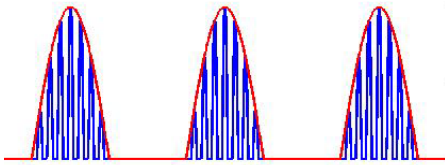
# Median Filter

- Computes the **median** of an array over sliding window with a given size.
- Can be used for removing an impulsive noises while preserving edges of the signal.

| 1 | 2 | 2 | 3 | 5 | 4 | 4 | 9 | 3 | 1 | 1 | 9 | 8 | 7 | 8 | 7 |

Length of window = 5

Median

| 2 | 2 | 2 | 3 | 4 | 4 | 4 | ? | 3 | 3 | 3 | 7 | 8 | 8 | 7.5 | 7 |

# Maximum Filter

## [Dilation Filter]

Estimate Signal Envelop

- Computes the **maxima** of an array over sliding window with a given size.
- Can be used for finding the envelop.

| 1 | 2 | 3 | 5 | 8 | 9 | 8 | 5 | 6 | 4 | 2 | 5 | 7 | 4 | 2 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

Length of window = 5

Max

| 3 | 5 | 8 | 9 | 9 | 9 | 9 | 9 | 8 | ? | 7 | 7 | 7 | 7 | 7 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|