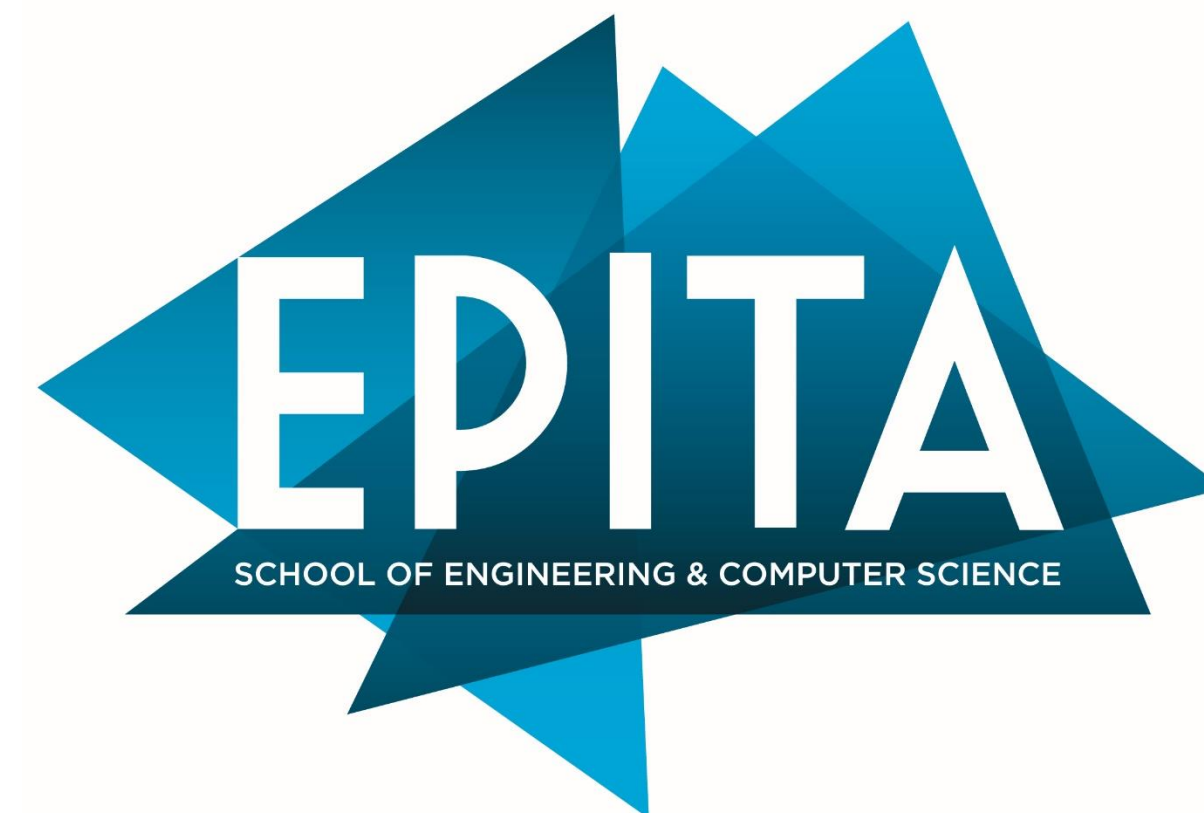




Arab Republic of Egypt
Ministry of Communications
and Information Technology

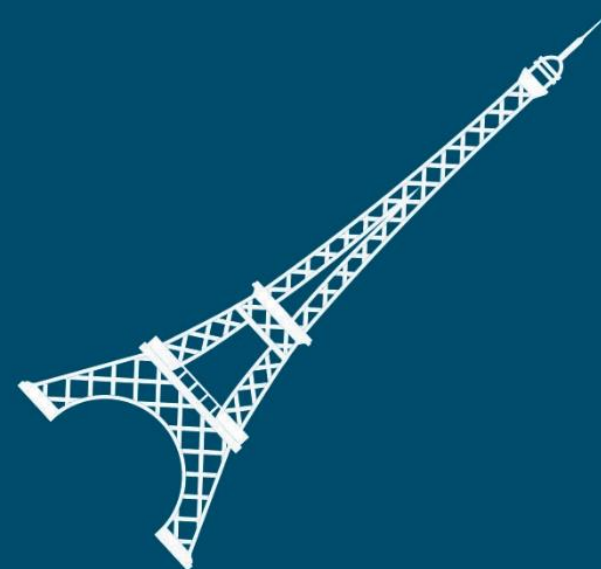


Information
Technology
Institute



A.I. IN AUDIO & SIGNAL PROCESSING

Session 1: Signal, Audio, Speech encoding



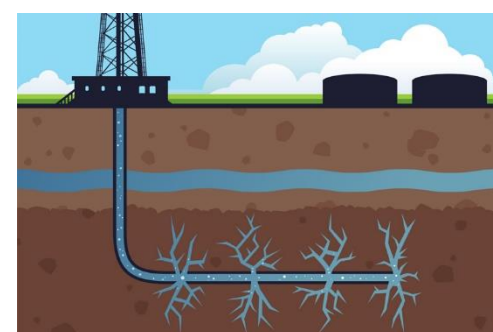
INTRODUCTION



Gaël Laqueille



CentraleSupélec
(Supélec)



Schlumberger



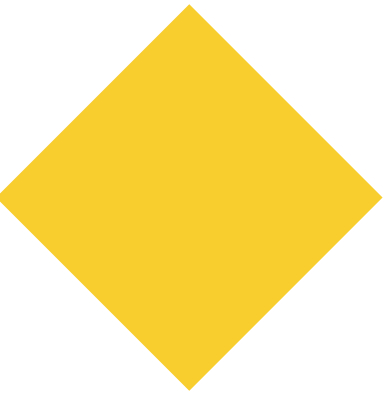
Parrot

move wireless



faurecia
inspiring mobility

COURSE STRUCTURE



Quick Summary

Audio processing for AI

- Signal, audio, speech encoding (4h)
- Deep learning for audio processing (4h)

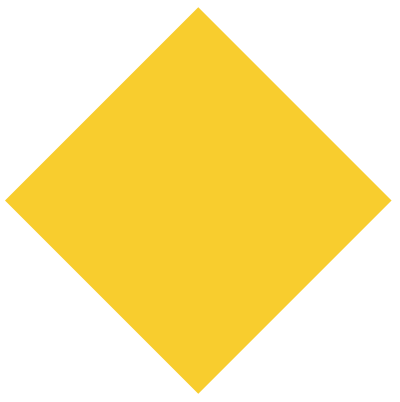
Automata for language modelling

- HMM for speech processing (4h)
- Automata and transducer (4h)

Towards speaking with an AI-bot

- Speech synthesis (4h)
- Automatic speech recognition (4h)
- Speaker and emotion recognition (4h)

SESSION 1: SIGNAL, AUDIO, SPEECH ENCODING



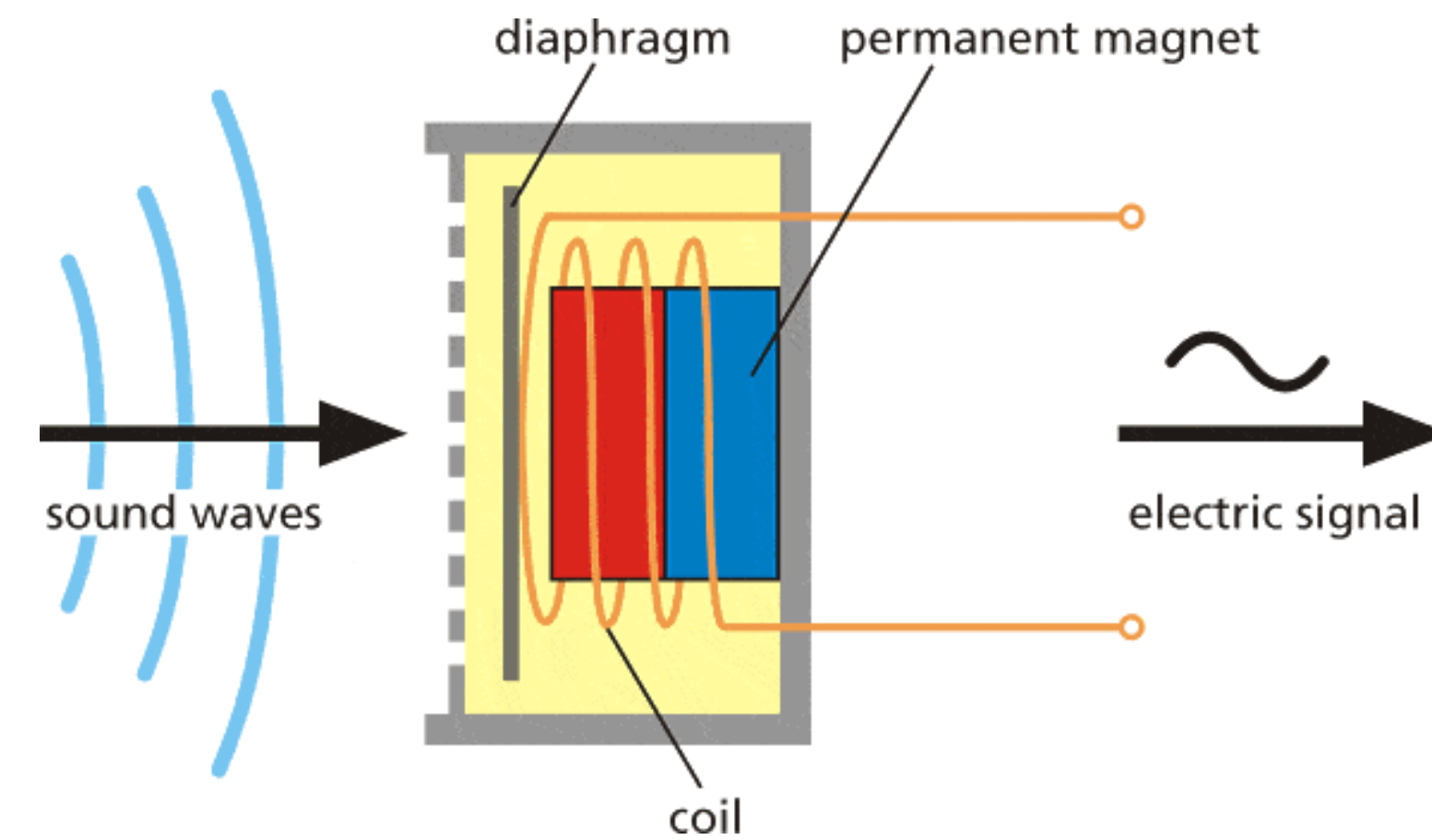
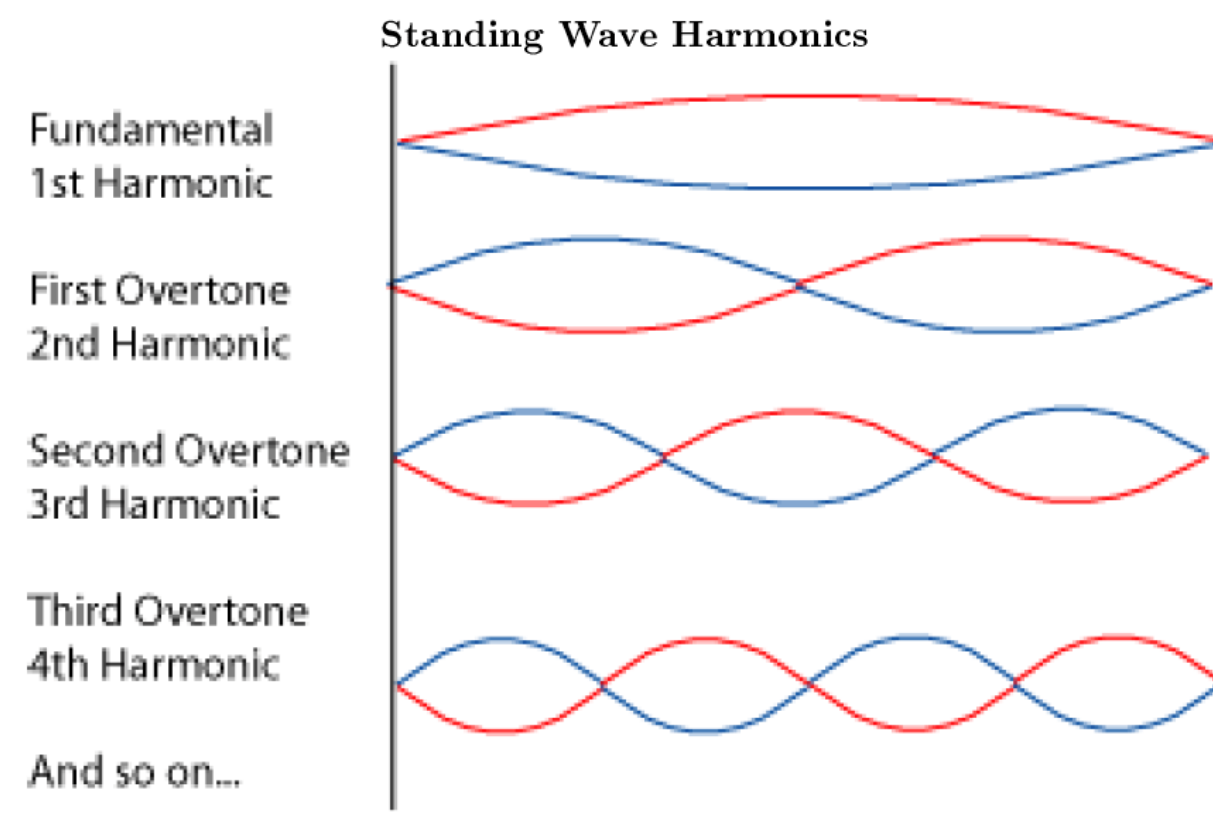
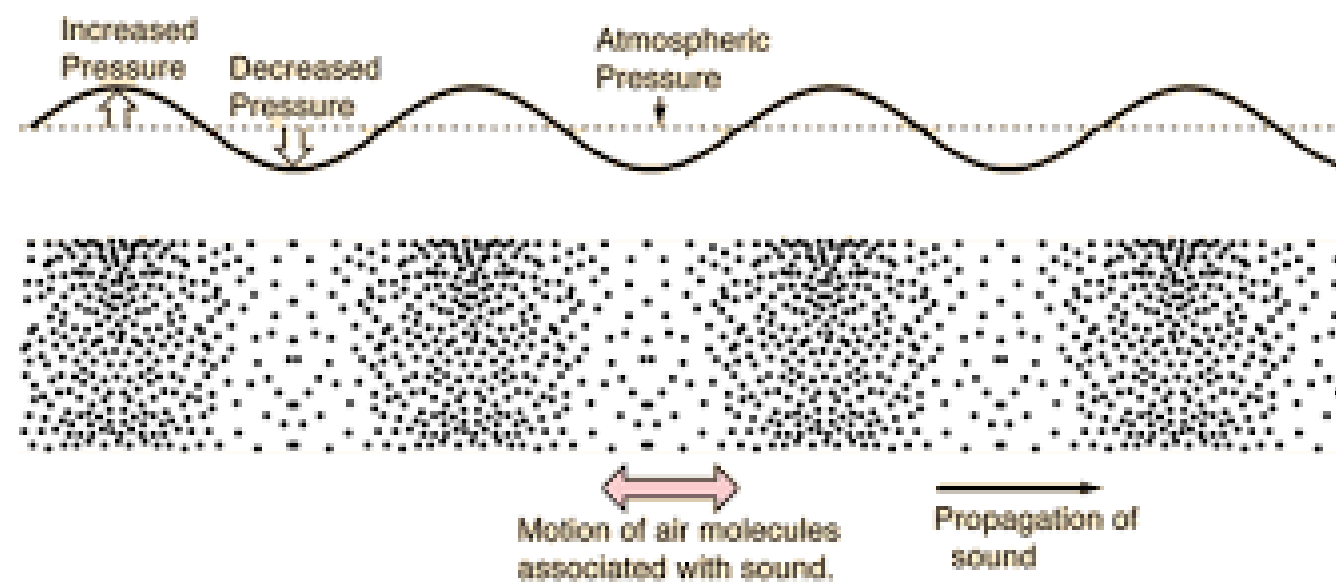
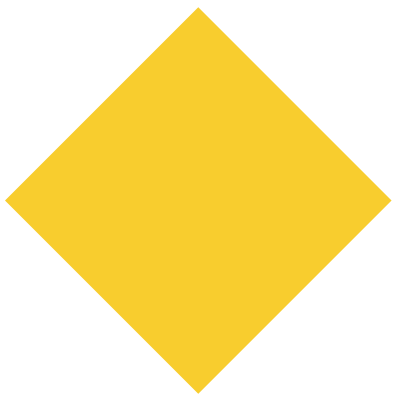
Quick Summary

1. The physics of sound
2. Signal representation and purpose
3. Signal characteristics
4. Signal models
5. Classical ML approaches

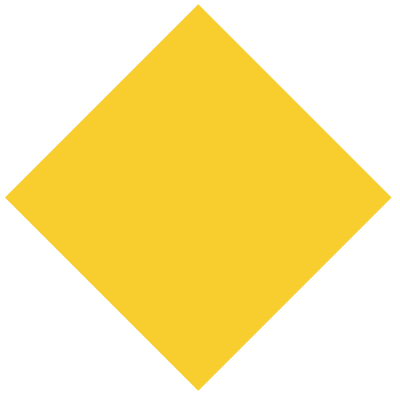
SIGNAL, AUDIO, SPEECH ENCODING.

The physics of sound

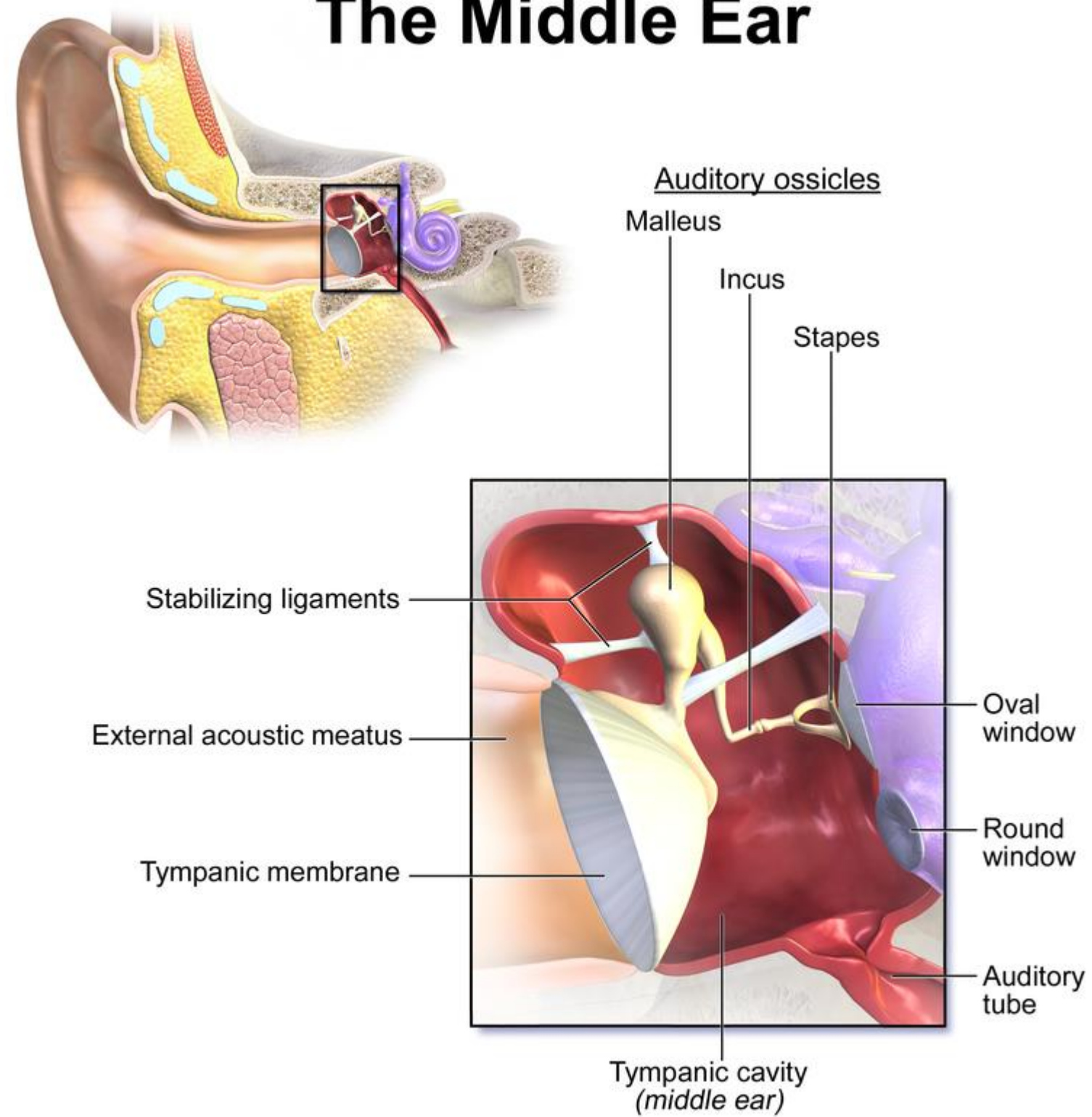
THE PHYSICS OF SOUND



THE PHYSICS OF SOUND



The Middle Ear



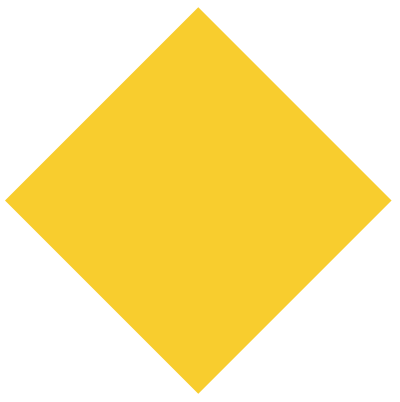
The Internal Ear



SIGNAL, AUDIO, SPEECH ENCODING.

Signal representation &
purpose

REPRESENTATION & PURPOSE



Different types of audio content

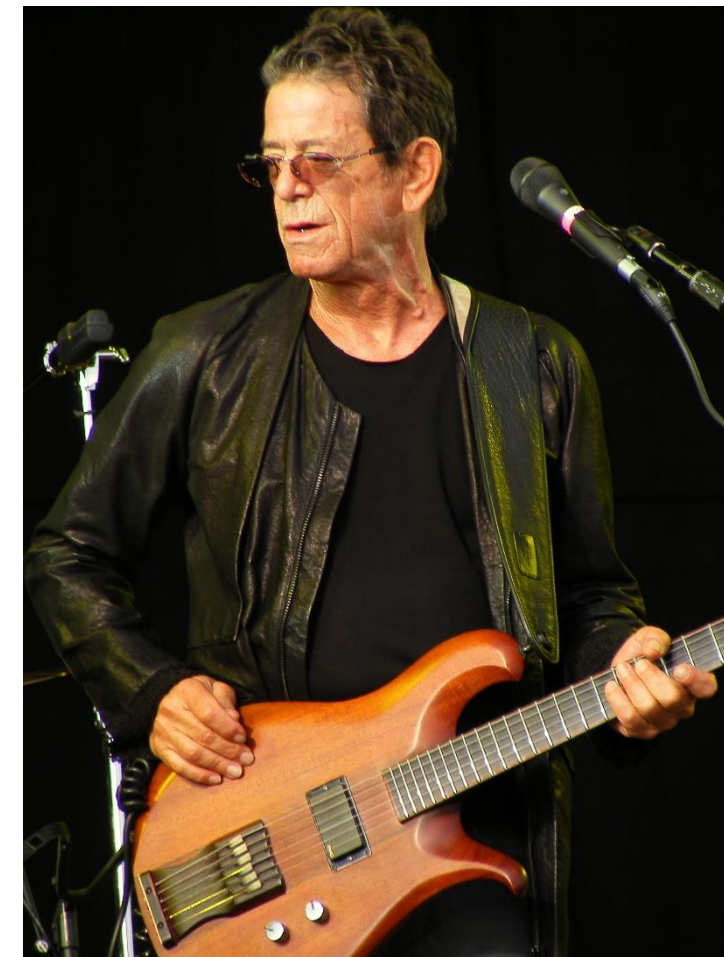
Noise



Multisource polyphonic with unstructure sound sources



Music



Multi timbre polyphonic with structured sound sources



Speech

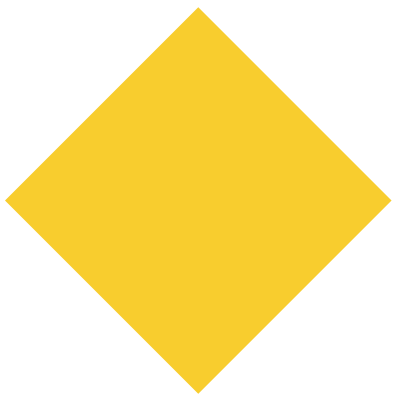


Monophonic with structured sound source



Extrapolation to other signals content: seismic, EEG,...

REPRESENTATION & PURPOSE



Signal & Audio symbolic representations

Codification



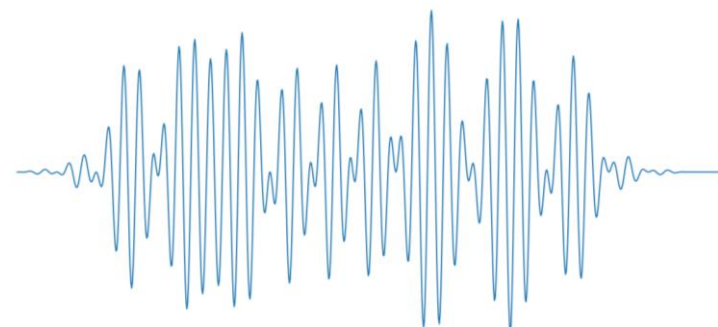
Production



Recording



Digitalization

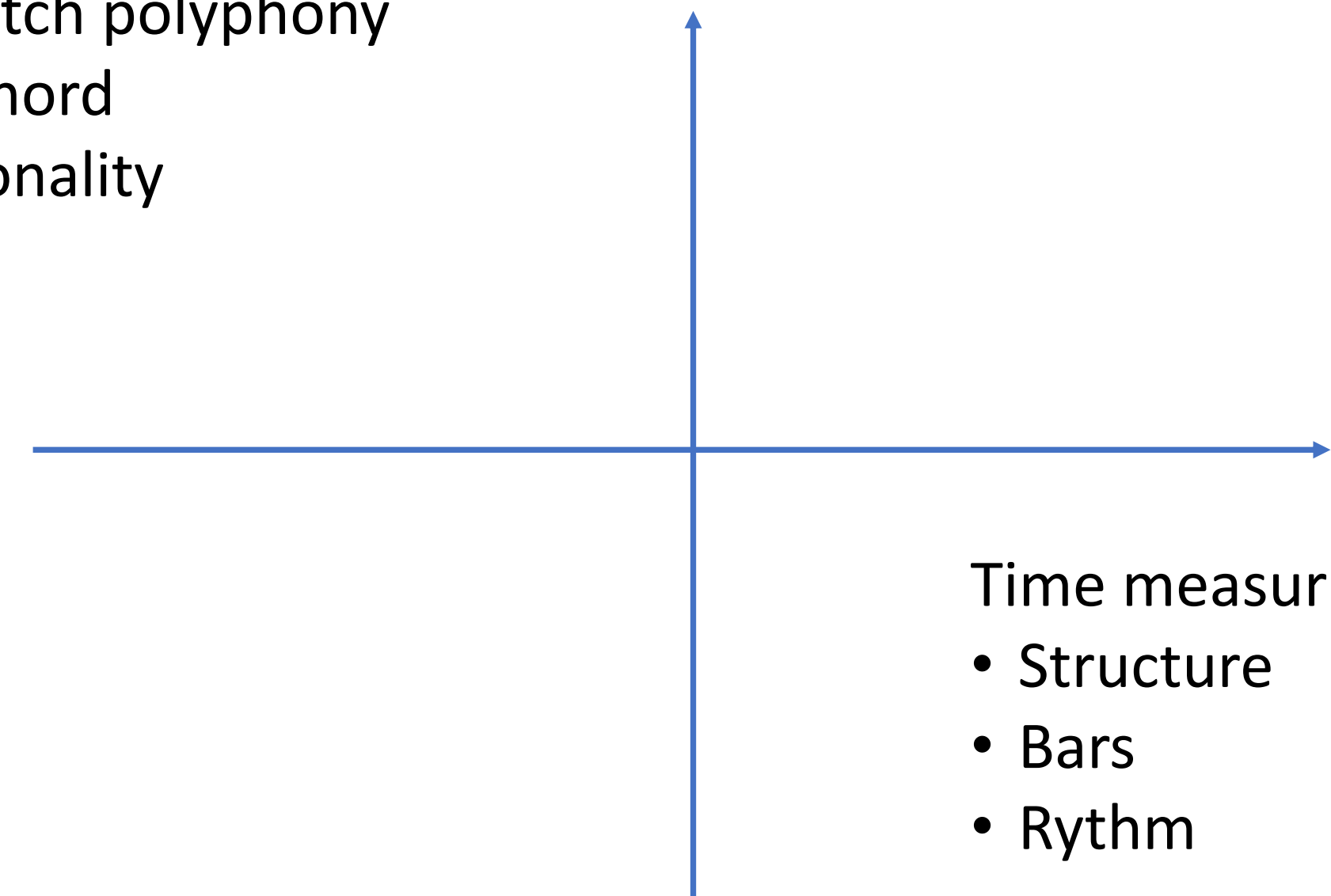


(Transcription)

Vertical and horizontal organization

Frequency measurement:

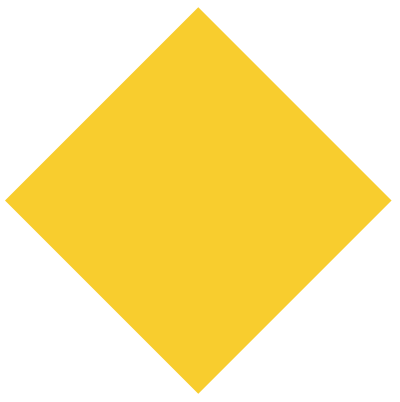
- Sources
- Timbres
- Pitch polyphony
- Chord
- Tonality



Time measurement:

- Structure
- Bars
- Rythm

REPRESENTATION & PURPOSE



Applications of AI in signal processing

Musical Information Retrieval (ISMIR)

- Automatic Music Transcription
- Auto-tagging
- Music Recommendation
- Source separation
- Audio/Music generation
- Style transfer
- Fingerprinting

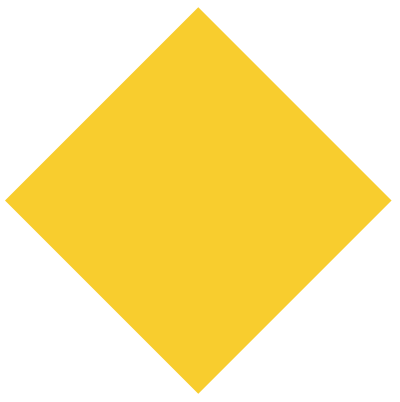
Spoken Content Retrieval

- Trigger word detection
- Segmentation
- Speech-to-text (ASR)
- Text-to-speech
- Emotion recognition
- Speaker recognition, diarization
- Denoising
- Active Noise Control

Other signal processing applications

- Time series prediction
(finance, insurance, healthcare,...)

REPRESENTATION & PURPOSE



Fourier series

$$\begin{aligned} f(t) &= \sum_{n=-\infty}^{+\infty} k_n \sin(2\pi n f_0 t + \phi_n) \\ &= \sum_{n=-\infty}^{+\infty} A_n \cos(2\pi n f_0 t) + B_n \sin(2\pi n f_0 t) \\ &= \sum_{n=-\infty}^{+\infty} c_n e^{j2\pi n f_0 t} \end{aligned}$$

Discrete Fourier Transform

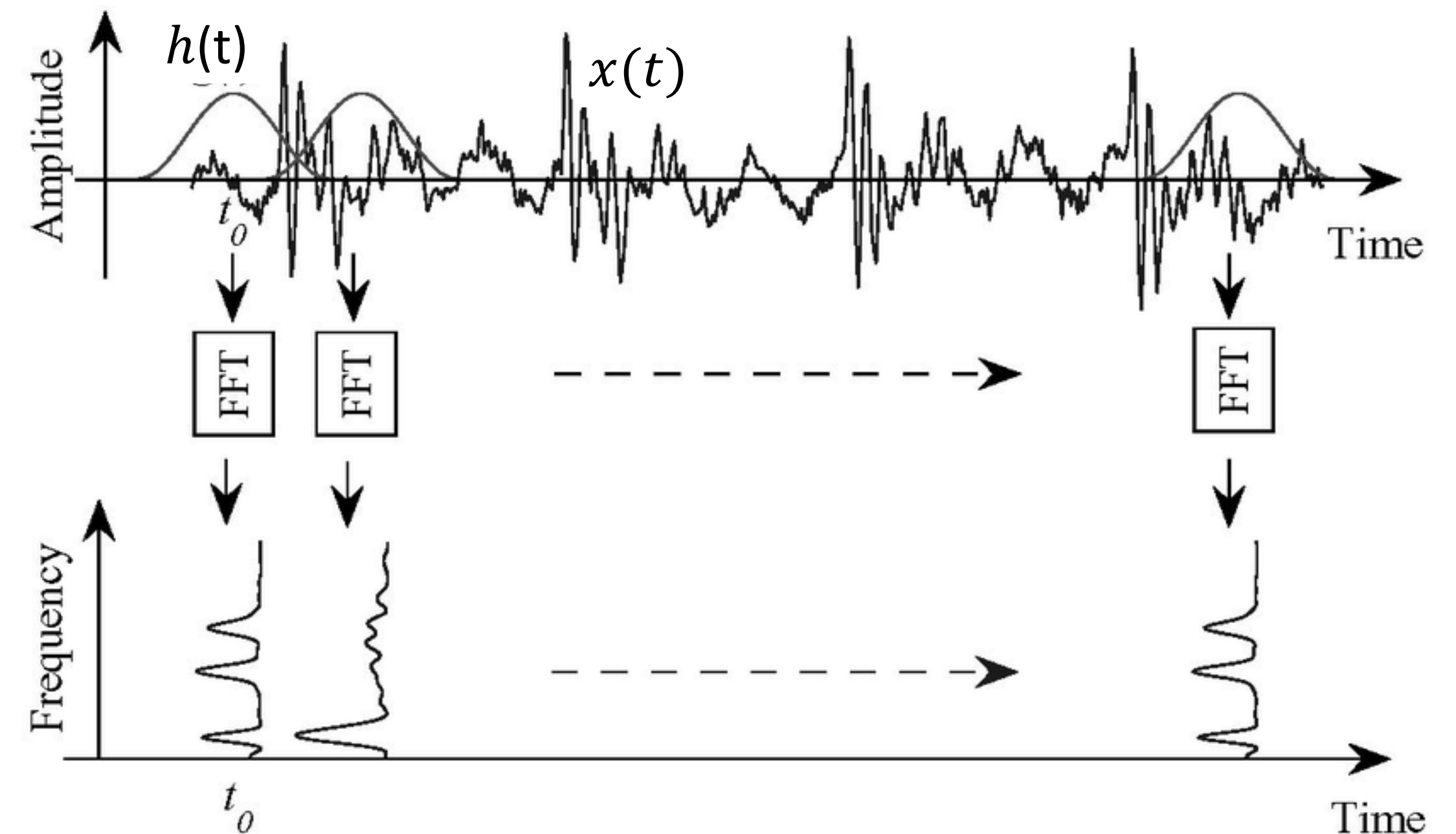
$$X(k) = \sum_{m=0}^{N-1} x(m) e^{-j2\pi \frac{k}{N} m}$$

Short-Time-Fourier-Transform (STFT)

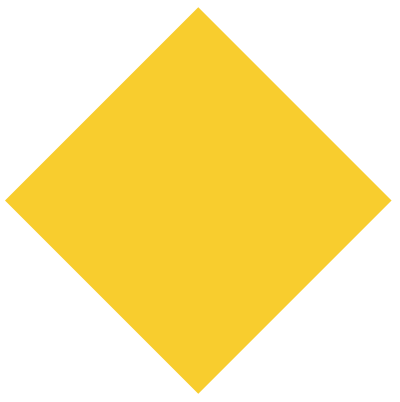
$$X(k, n) = \sum_{m=0}^{N-1} x(m) h(n - m) e^{-j2\pi \frac{k}{N} m}$$

h : windowing

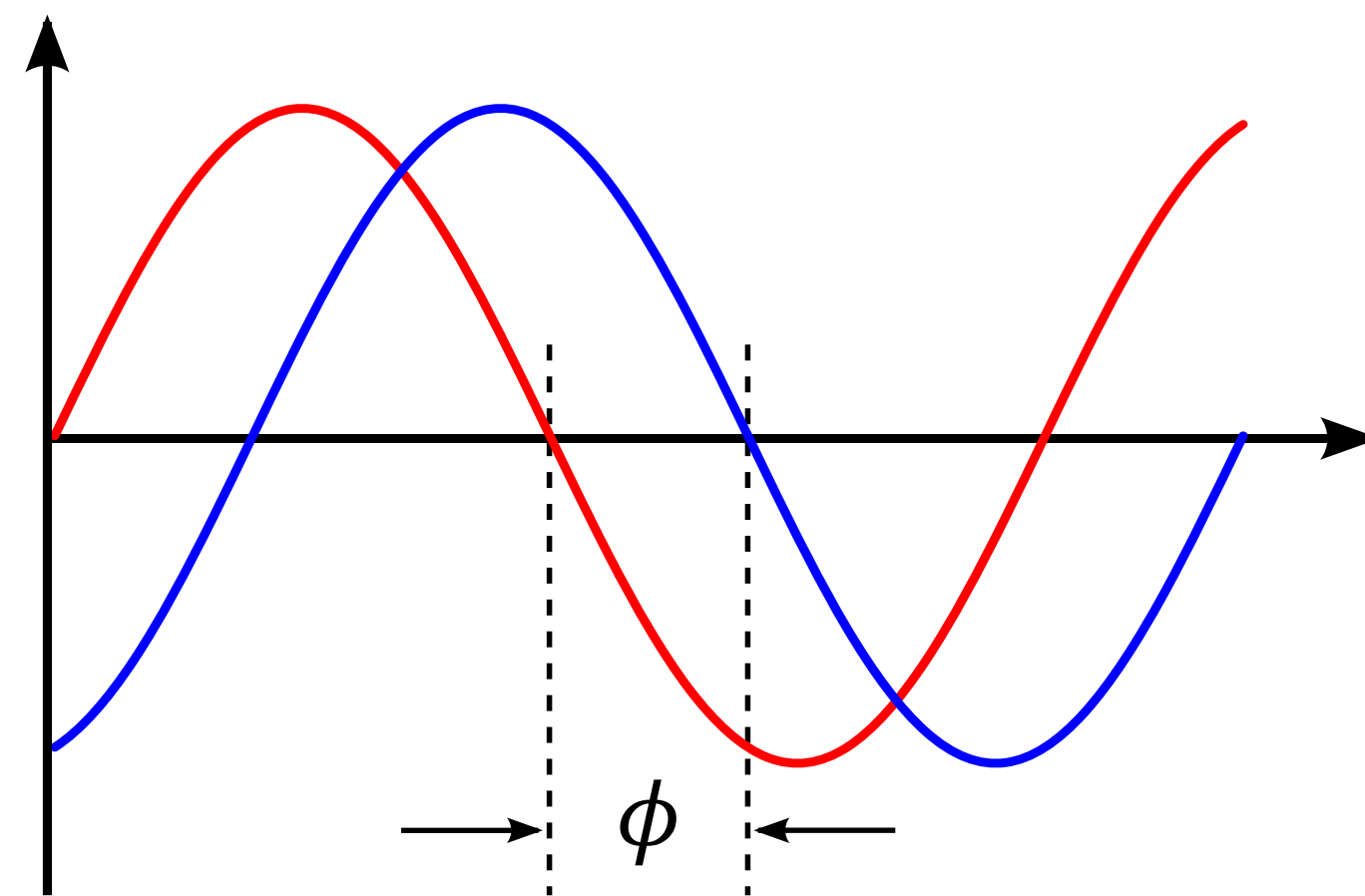
$k_n, A_n, B_n, c_n, x(m), h(m) \in \mathbb{R}$



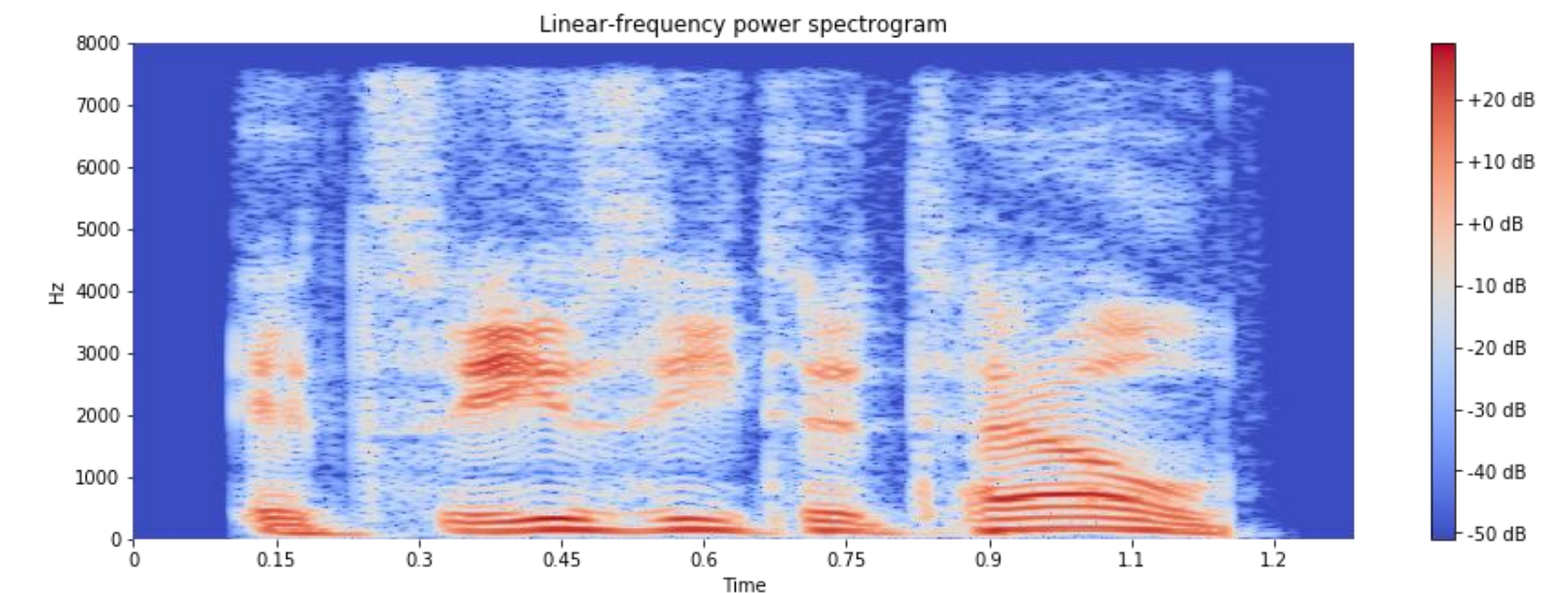
REPRESENTATION & PURPOSE



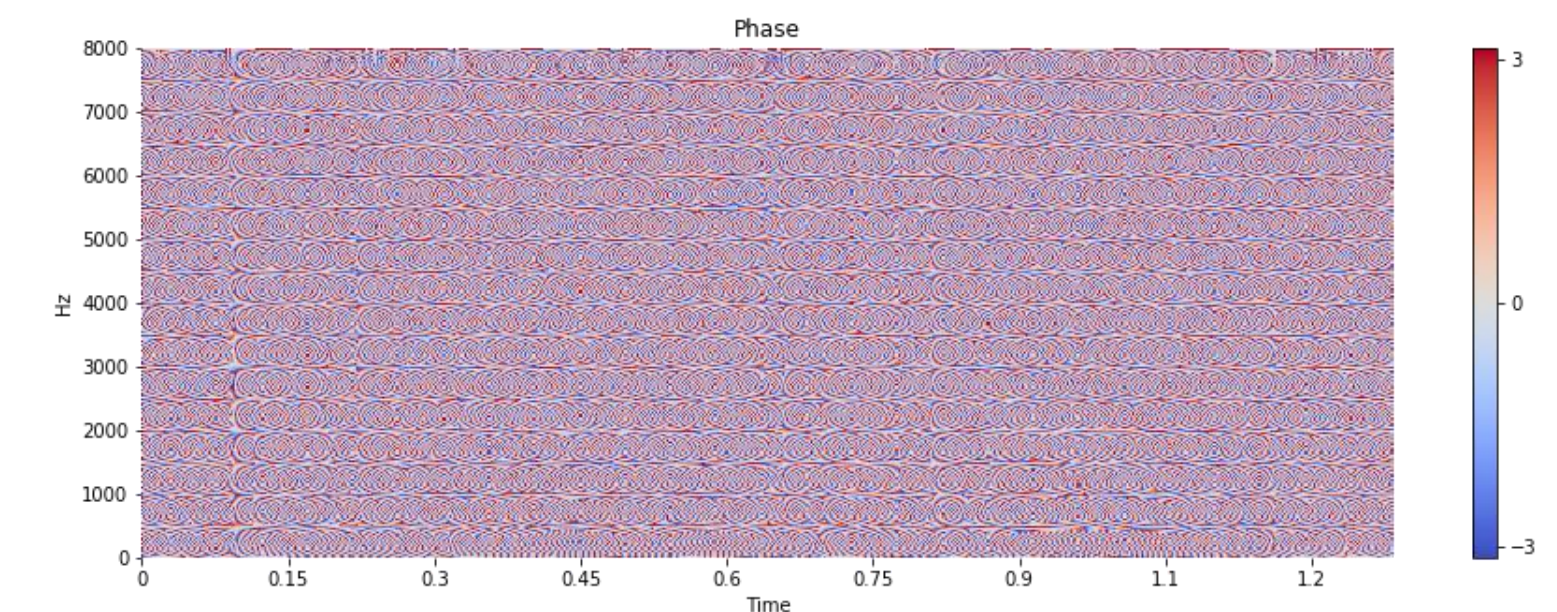
- $X(k) = |X(k)|e^{j\phi_X(k)}$
- Amplitude $|X(k)|$
- Phase $\phi_X(k)$, temporal localization of information
(Phase is necessary to reconstruct temporal signal)
- Instantaneous frequency



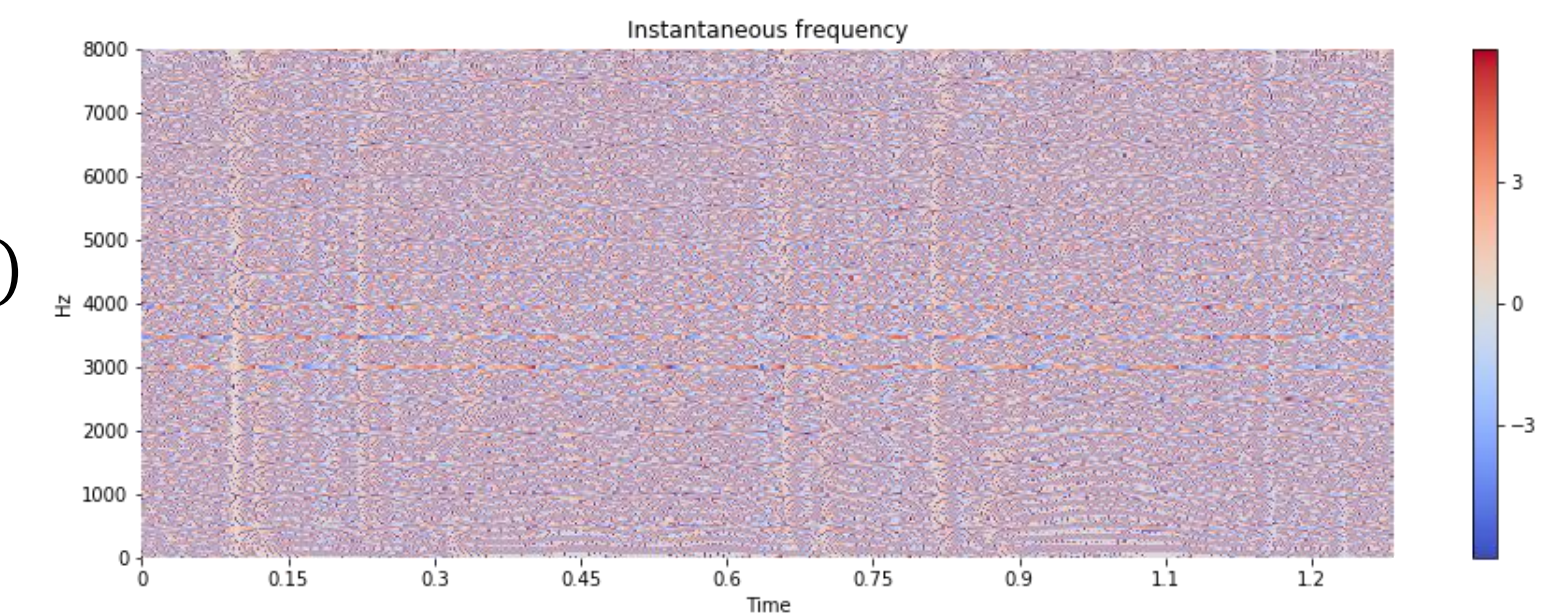
$$10 \log_{10} |X(k, n)|^2$$



$$\phi_X(k, n)$$

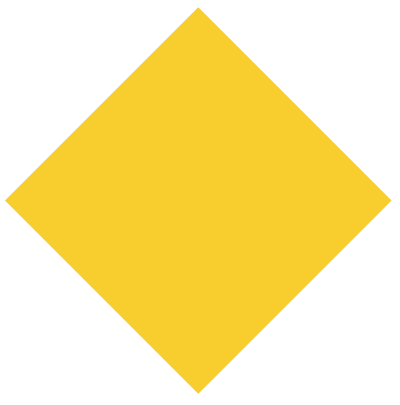


$$\phi_X(k, n) - \phi_X(k, n - 1)$$

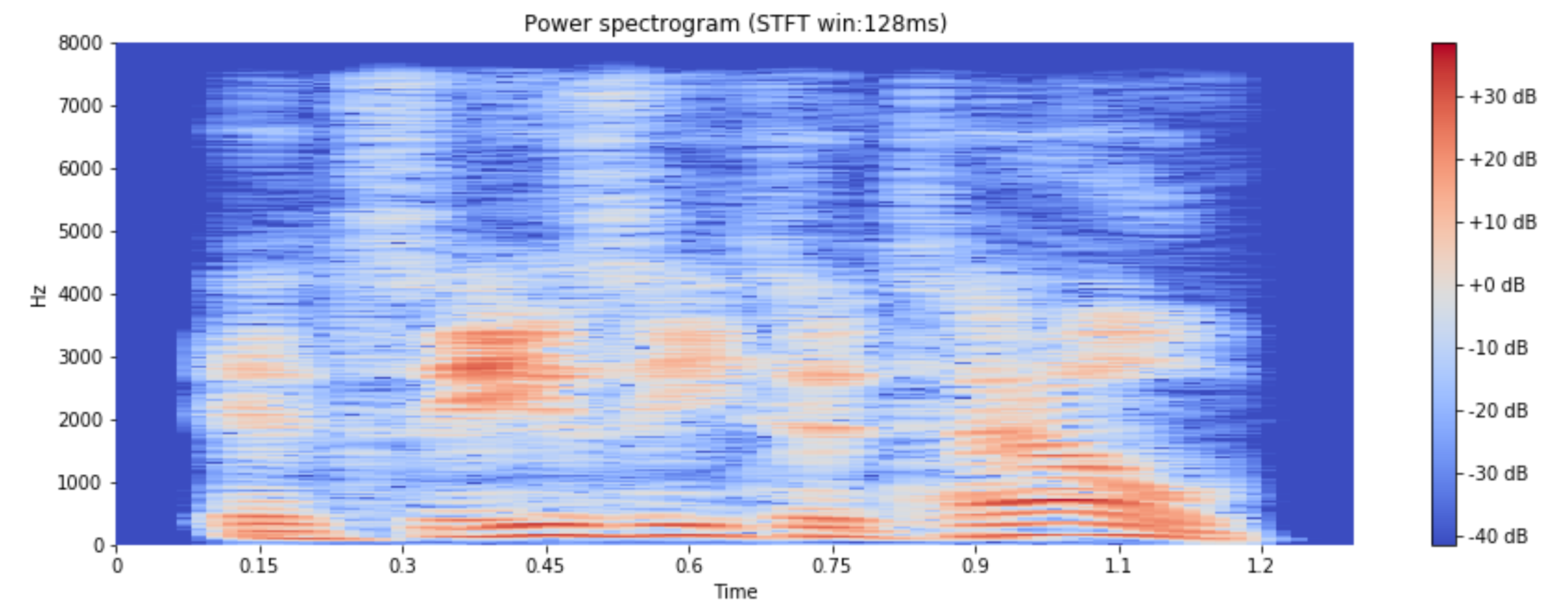
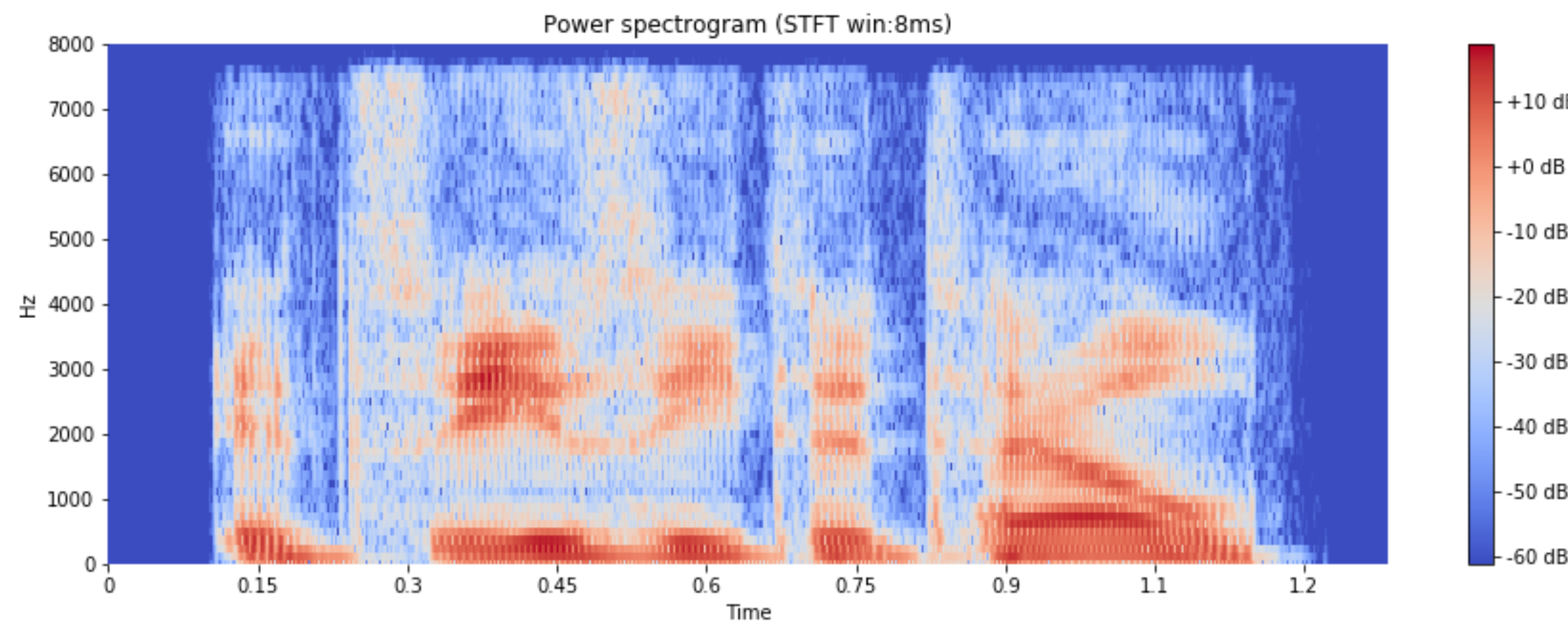


k : frequency bins
 n : time windows index

REPRESENTATION & PURPOSE



Remark on trade-off temporal/frequential resolution



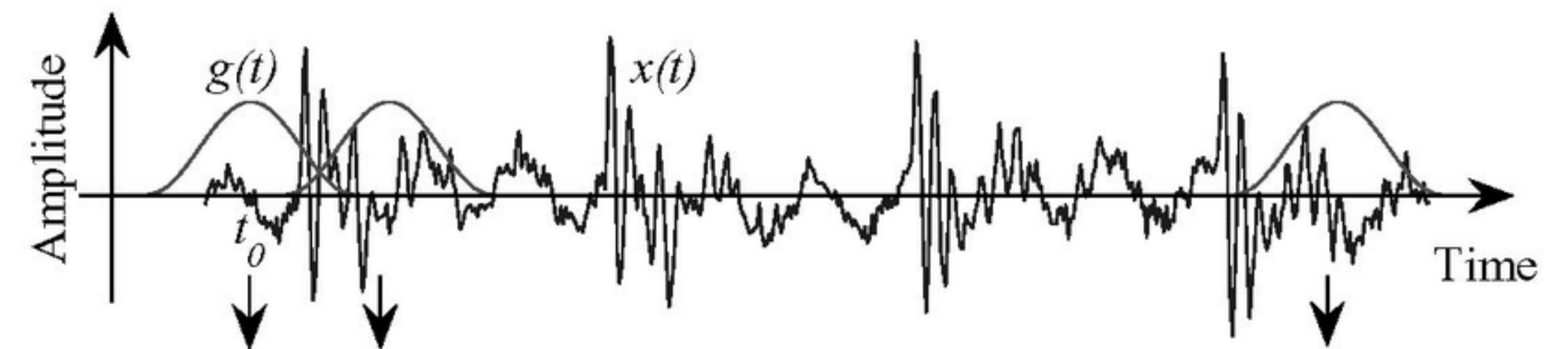
Remark on sound perception

- log perception of sound amplitude
→ magnitude spectrogram usually in dB

$$10 \log_{10} |X(k, n)|^2$$

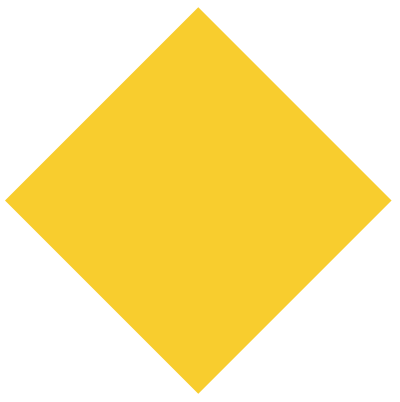
- log perception of sound frequency
→ magnitude spectrogram often in mel scale

$$10 \log_{10} |X(k|_{mel}, n)|^2$$



with mel scale: $m = 2595 \log_{10} \left(1 + \frac{f}{700} \right)$

REPRESENTATION & PURPOSE



Wavelet transform (WT)

(not used in practice, but concept is used to define CQT, see next slide)

- Continuous wavelet basis definition

$$\forall t \in \mathbb{R}, \psi_{s,\tau}(t) = \frac{1}{\sqrt{s}} \psi\left(\frac{t-\tau}{s}\right)$$

ψ is the “mother” wavelet

$\psi_{s,\tau}$ are the “child” wavelets that form a basis
 s and τ are the dilation and translation variables

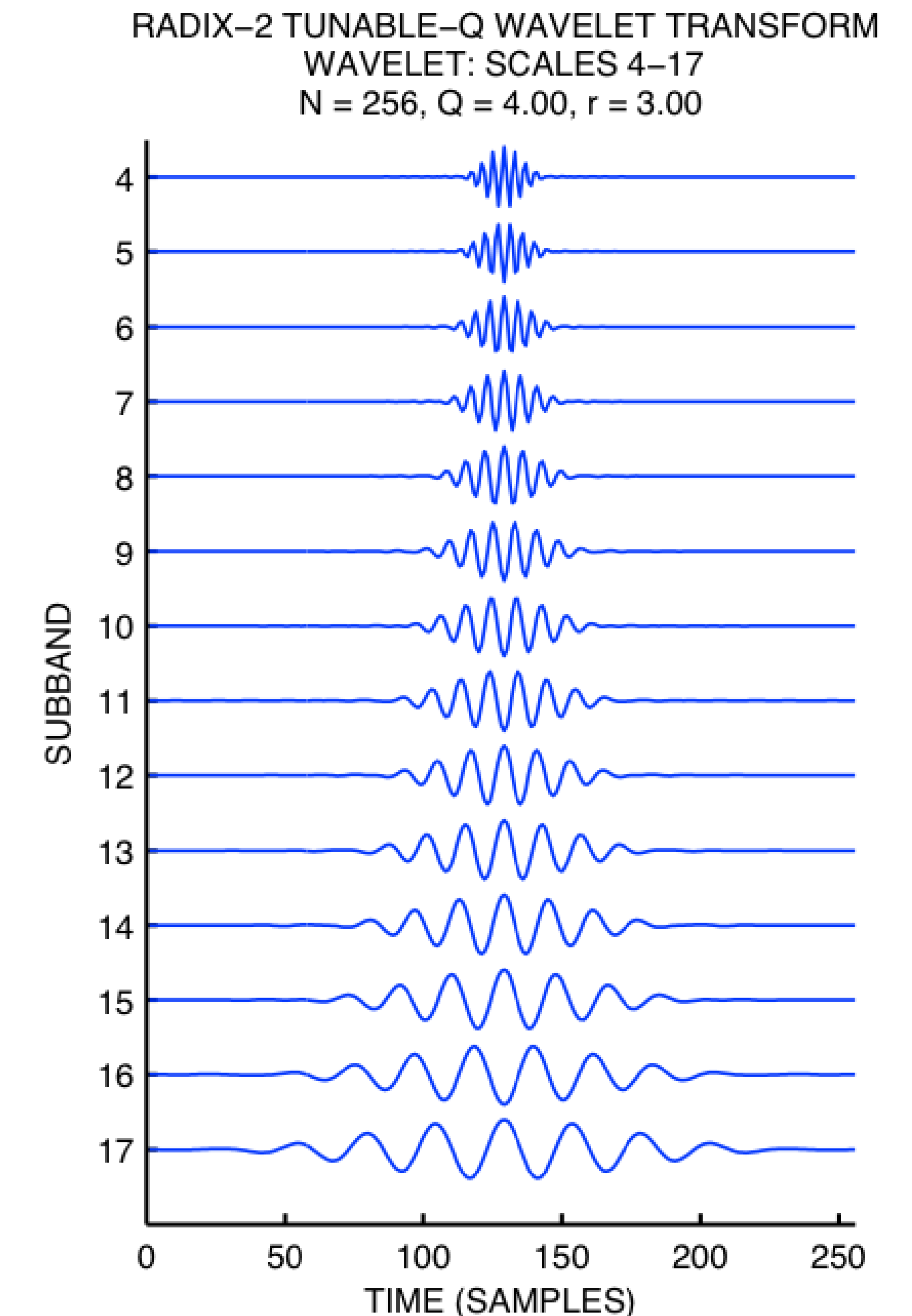
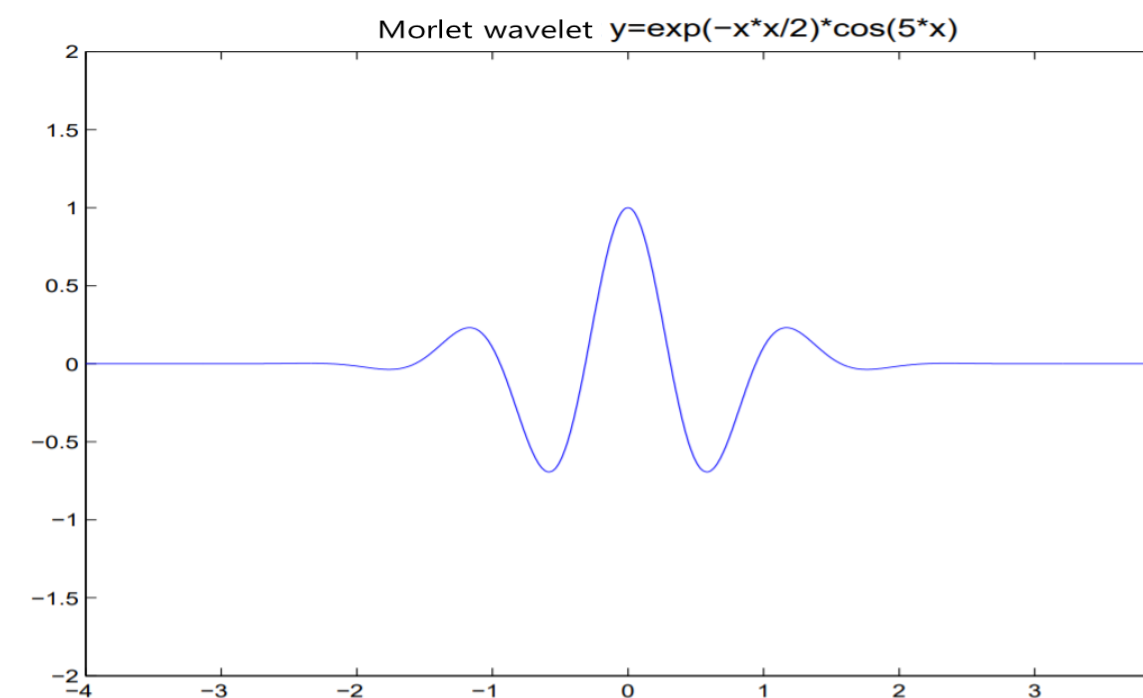
- Discrete wavelet basis definition

- discretization in wavelet (time-frequency) domain: $\tau \rightarrow n$ and $s \rightarrow k$
- discretization in time domain: $t \rightarrow m$

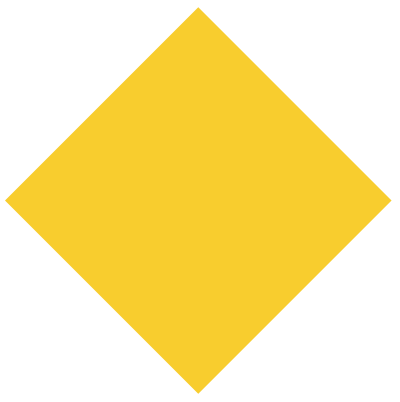
$$\psi_{k,n}[m] = s_0^{-\frac{k}{2}} \psi(s_0^{-k}m - n\tau_0) \quad \text{with } s_0 \text{ and } \tau_0 \text{ constants}$$

- Discrete wavelet transform definition

$$W_\psi[k, n] = \frac{1}{\sqrt{N}} \sum_{m=0}^{N-1} x[m] \psi_{k,n}[m] \quad \text{with } n \text{ time index and } k \text{ frequency bin}$$



REPRESENTATION & PURPOSE



Constant-Q transform (CQT)

(behavior of WT but no orthogonality CQ functions)

- Frequency resolution changing with frequency considered

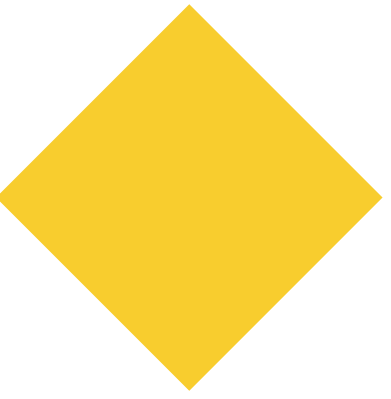
$$N_{[k]} = Q \frac{f_s}{f_k}$$

- A pitch translation corresponds to a frequency translation on CQT

$$X_{[k,n]}^{CQ} = \sum_{m=n-\frac{N_k}{2}}^{n+\frac{N_k}{2}} \frac{h\left(\frac{n-j}{N_k} - 1/2\right)}{N_k} e^{-i2\pi Q\left(\frac{n-j}{N_k} - 1/2\right)}$$

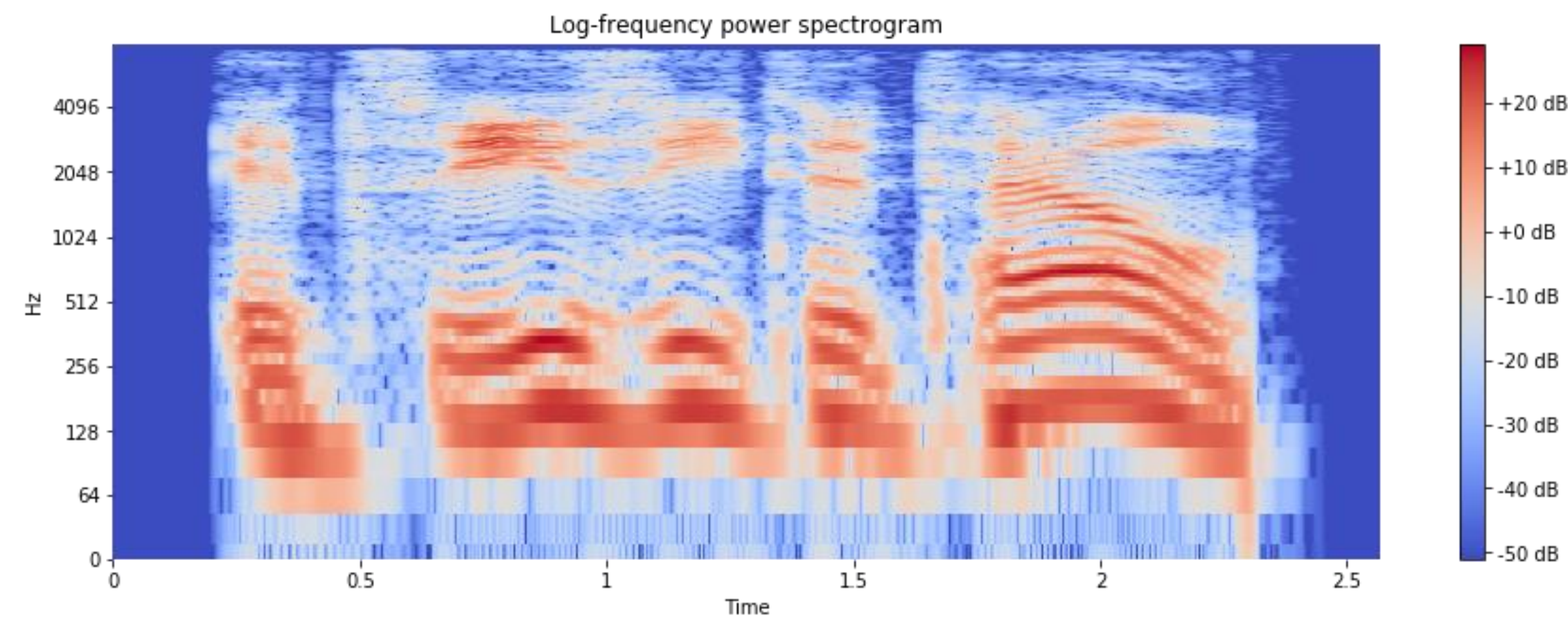
h is a window (Hanning, Blackmann Harris, ...)

REPRESENTATION & PURPOSE

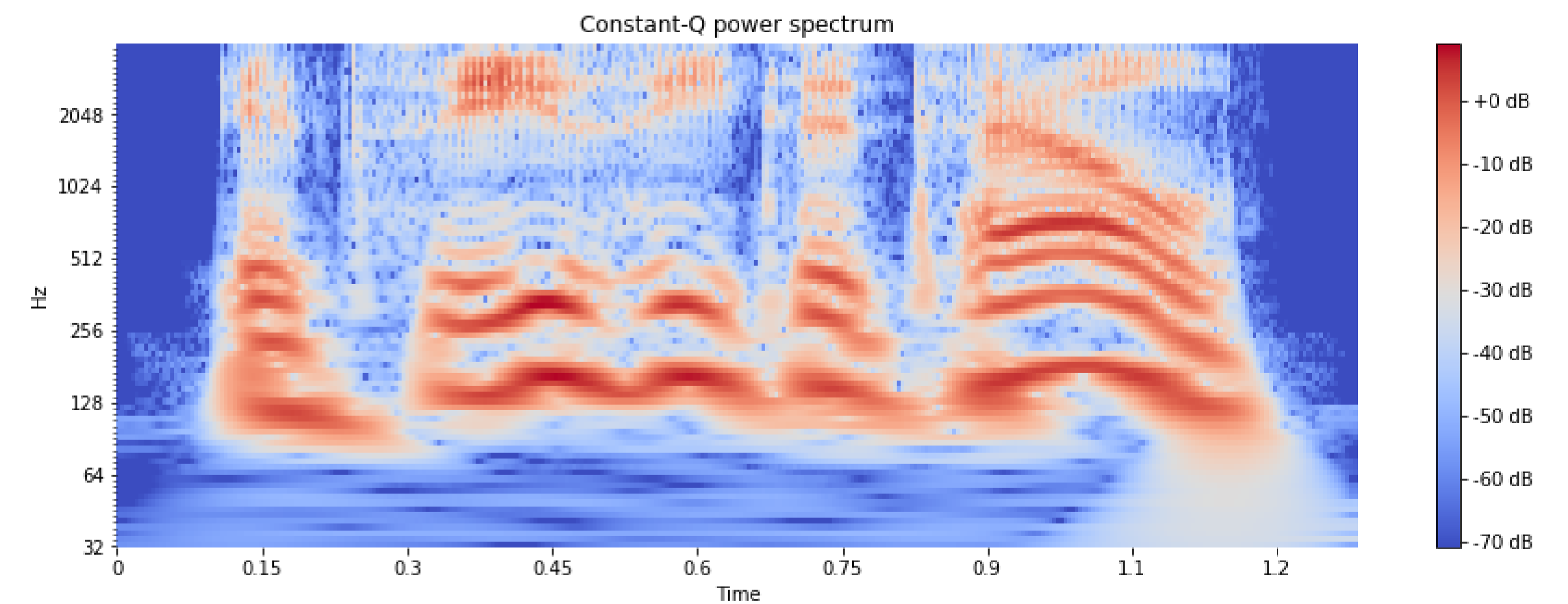


FT vs CQT

Fourier transform



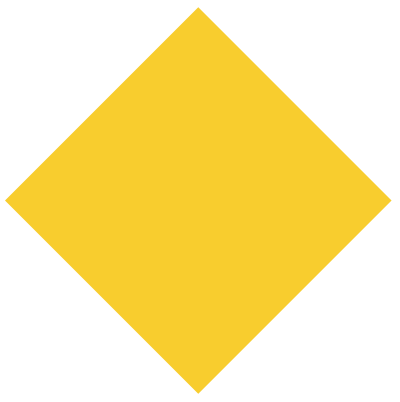
Constant-Q transform



SIGNAL, AUDIO, SPEECH ENCODING.

Signal characteristics

CHARACTERISTICS



Digital waveform characteristics

- Sample rate
- Bits depth
- level dBfs
- mean RMS (dBfs)
- Compression(dBA, dBC, SPL)

Artifacts

- noise (SNR)
- amplitude distortion (signal balance) target curve
- phase distortion (unprecise onset, spatialization)
- non-linearities (clipping for instance) THD
- delay (real-time, digitalization or processing)

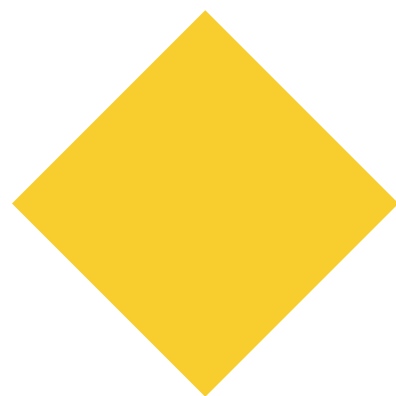
Sources of artifacts

- Sound or signal capture,
- sound or signal production
- sound propagation in air
- signal propagation in wires,
- processing

SIGNAL, AUDIO, SPEECH ENCODING.

Signal models

SIGNAL MODELS



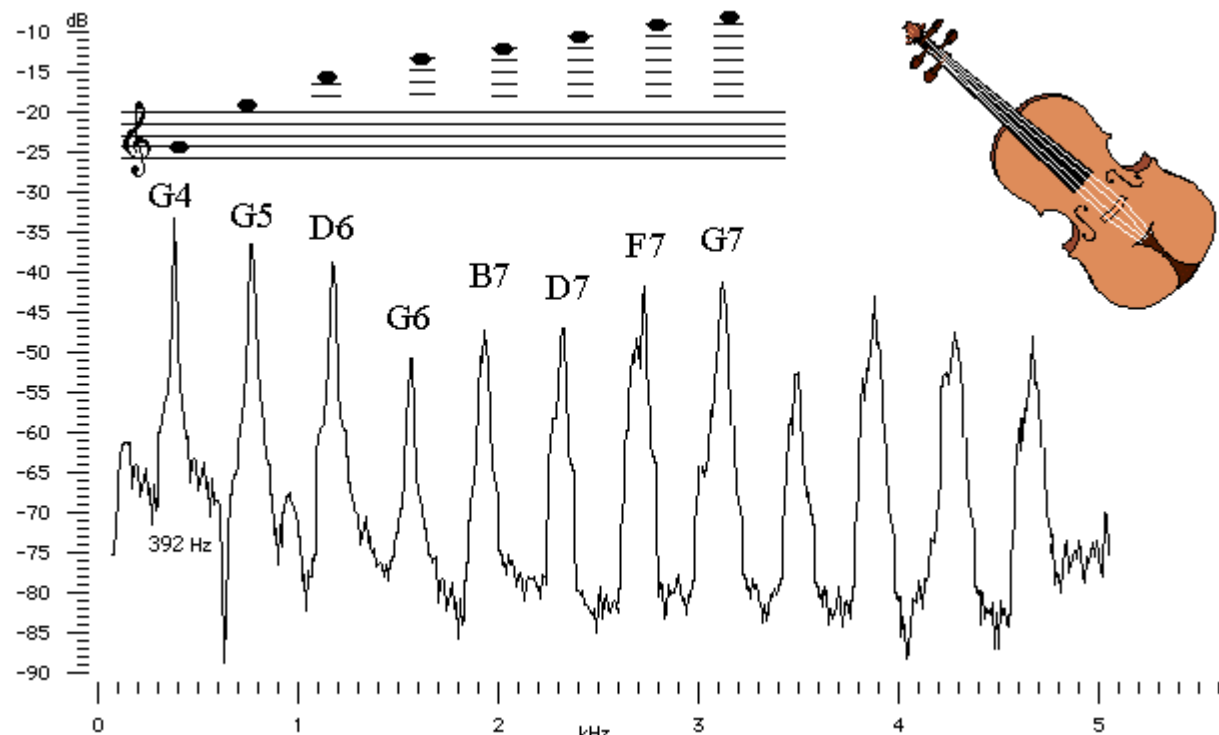
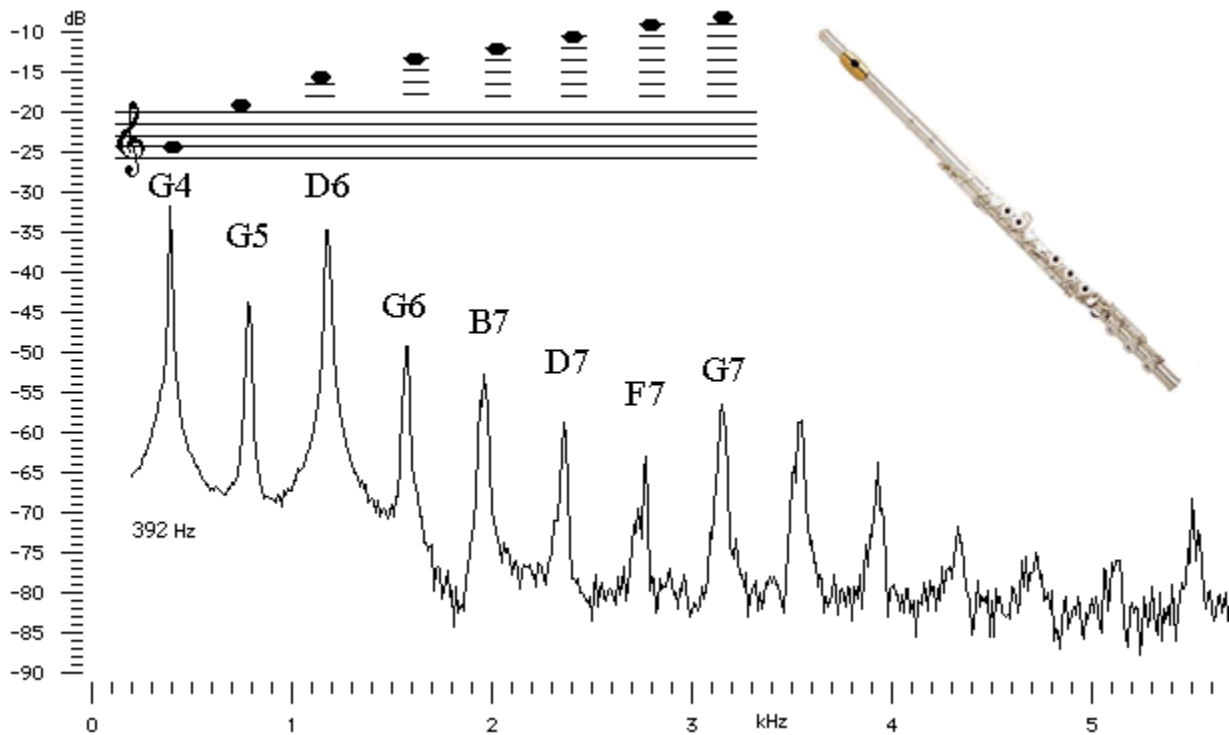
Sine wave harmonics and noise model:

$$x(t) = \sum_{h=1}^{H(t)} A_h(t) \cos(\phi_h(t)) + b(t)$$

Assumption that the audio signals follows this model.

The spectral envelope $\{A_h\}$ defines the timbre of the instrument.

Example of voice or instruments harmonic spectrum (spectral envelop)



Frequency	Order	Name	Standing wave representation
$1 \times f = 440 \text{ Hz}$	$n = 1$	1st harmonic (fundamental tone)	
$2 \times f = 880 \text{ Hz}$	$n = 2$	2nd harmonic (1 st overtone)	
$3 \times f = 1320 \text{ Hz}$	$n = 3$	3rd harmonic (2 nd overtone)	
$4 \times f = 1760 \text{ Hz}$	$n = 4$	4th harmonic (3 rd overtone)	

SIGNAL MODELS

Source/filter model

- Periodical pulse convolved with a resonant filter

$$x(t) = e(t) * h(t)$$

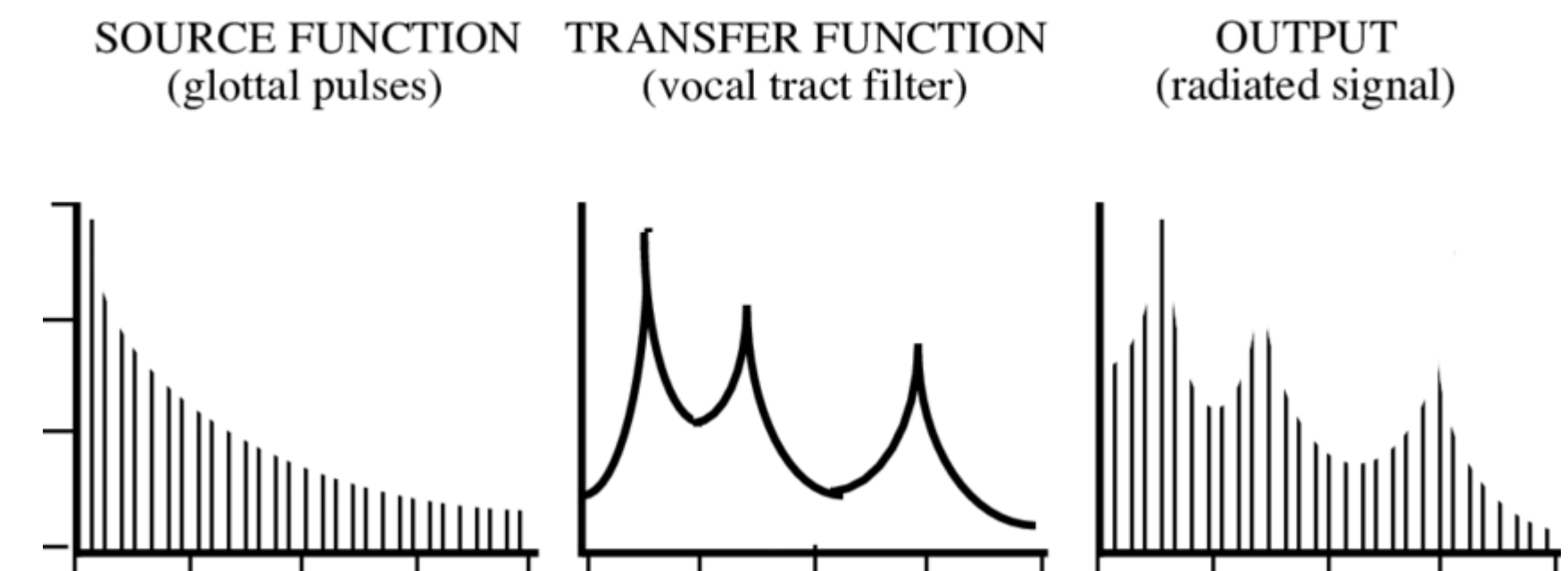
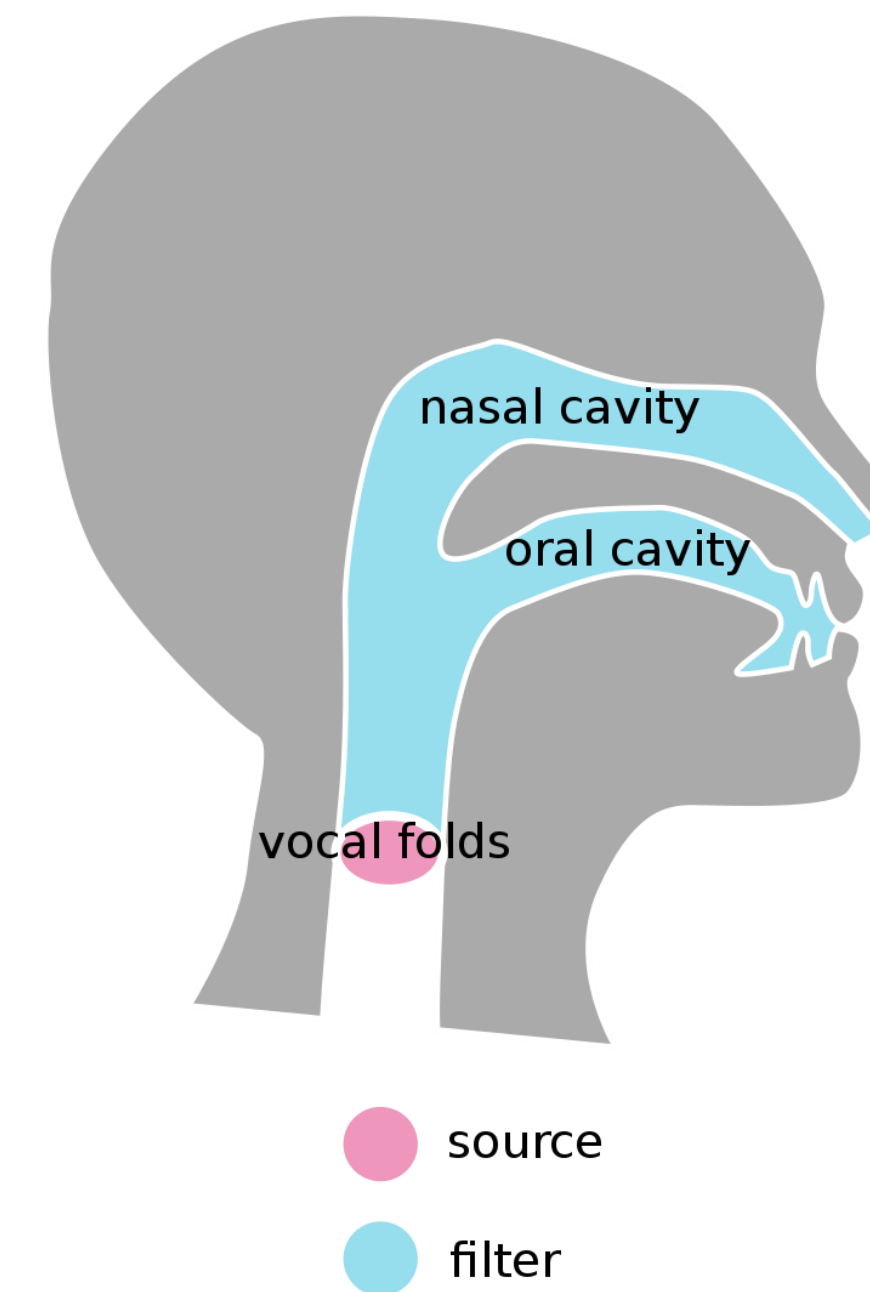
$$X(f) = E(f) \cdot H(f)$$

- Resonant filter is an AR filter:

$$H(z) = \frac{G}{1 + \sum_{k=1}^K a_k z^{-k}}$$

$x(n)$ can be predicted with a linear combination of previous values

$$x(n) \approx \sum_{k=1}^K a_k x(n - k)$$

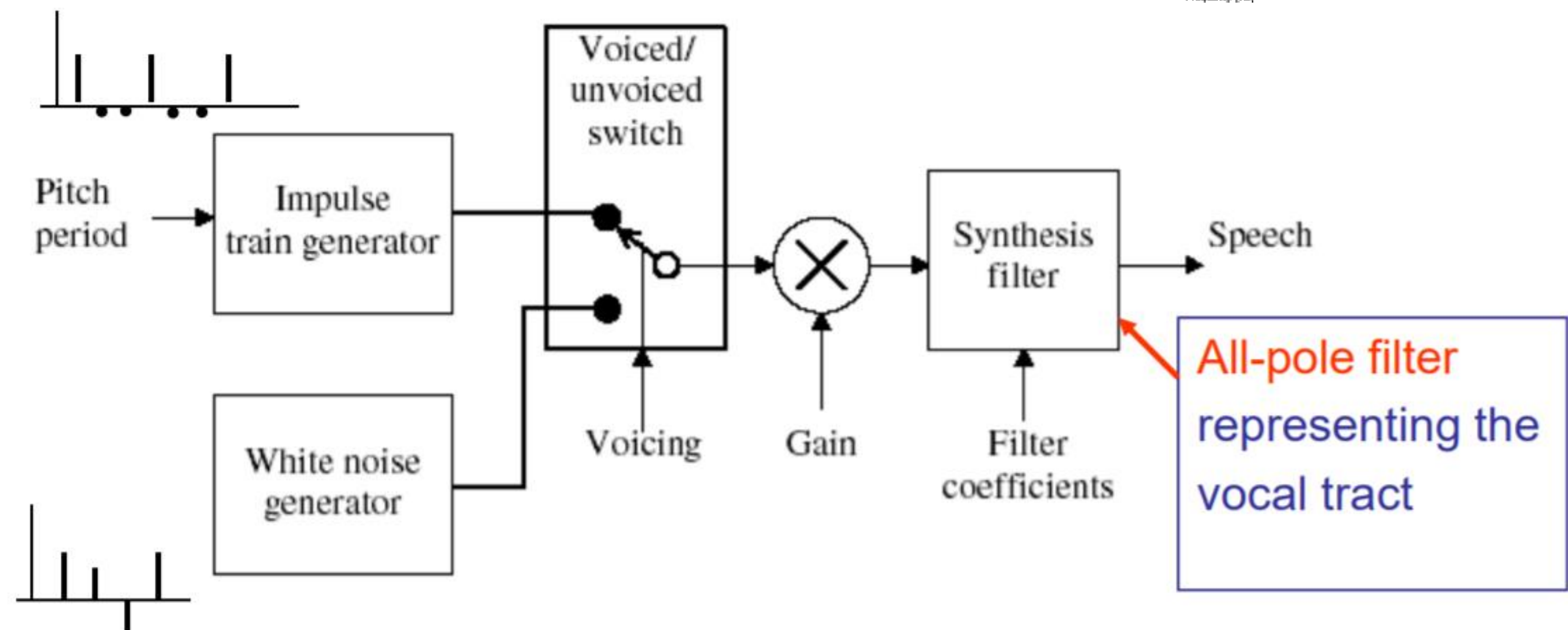
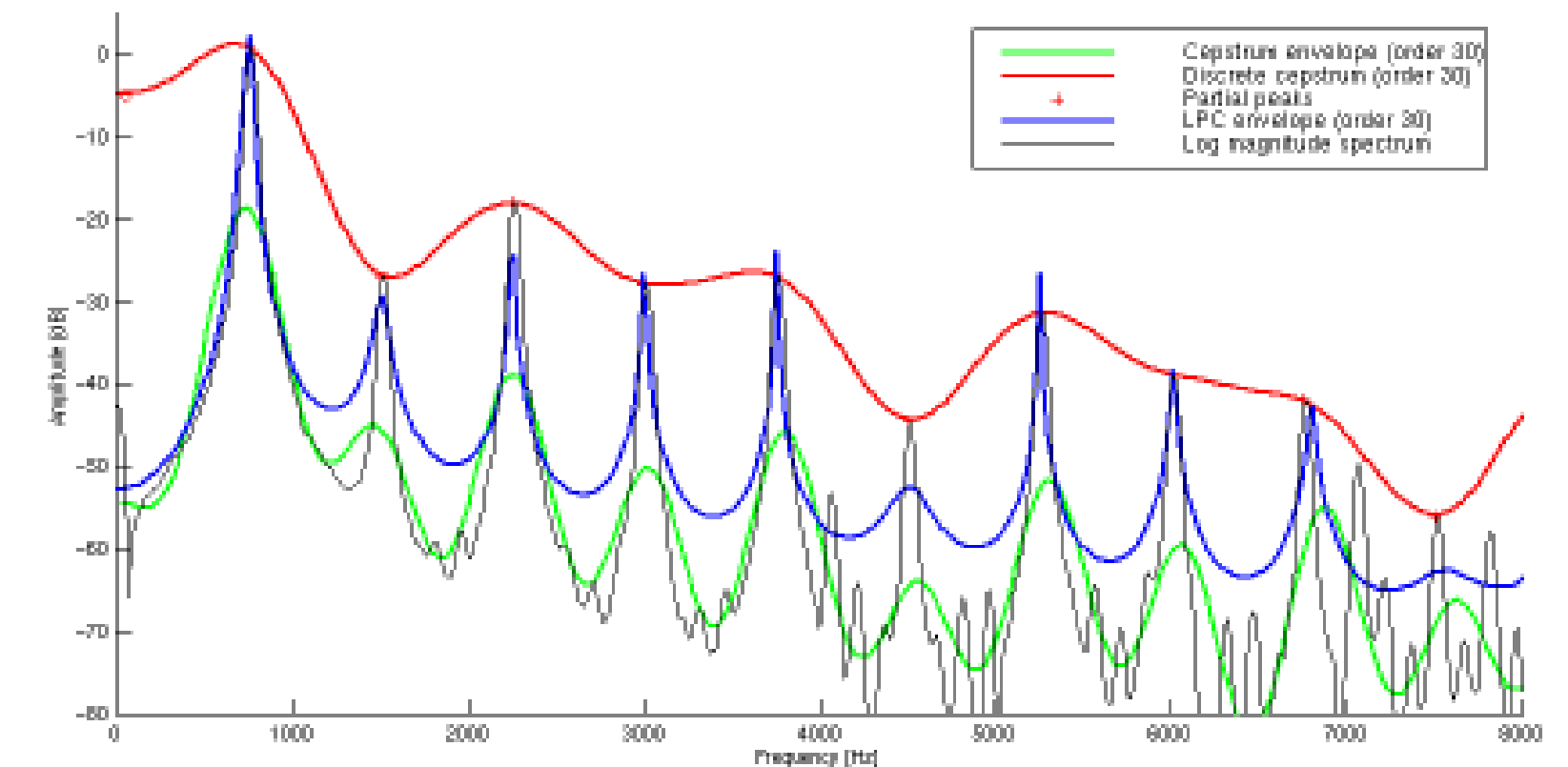


SIGNAL MODELS

Classical speech production model: LPC

- Impulse train or white noise (voiced/unvoiced speech frame)
- Gain: energy level of frame
- AR synthesis filter
- Pitch period (pitch height)

$$y(n) = \sum_{i=1}^p a_i y(n-i) \pm Gx(n)$$



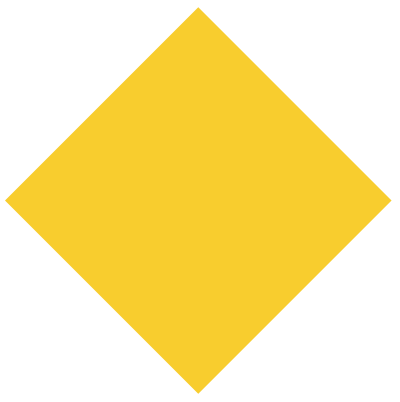
Variants: LPC-10, CELP, MELP, RELP, VSELP, ASELP, LD-CELP

Ref: Nimrod Peleg/IRCAM

SIGNAL, AUDIO, SPEECH ENCODING.

Classical ML approaches

CLASSICAL ML APPROACHES



Hand design features

Knowledge a priori for the design of audio features

MFCC

- Source-filter model $X(f) = E(f) \cdot H(f)$
- Cepstrum

$$cepstr(x)(\tau) = DCT^{-1}(\log(|X(f)|^2))$$

$$f(\tau) = DCT^{-1}(\log(|H(f)|^2)) + DCT^{-1}(\log(|E(f)|^2))$$

- Filter/source contribution separation
- Perception frequency scale: Mel

Chromas

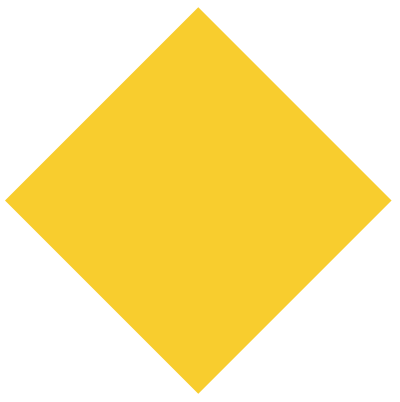
Contribution of frequency height class

$$f(C) = X^2(f_{C_1}) + X^2(f_{C_2}) + \dots$$

What are MFCC ?

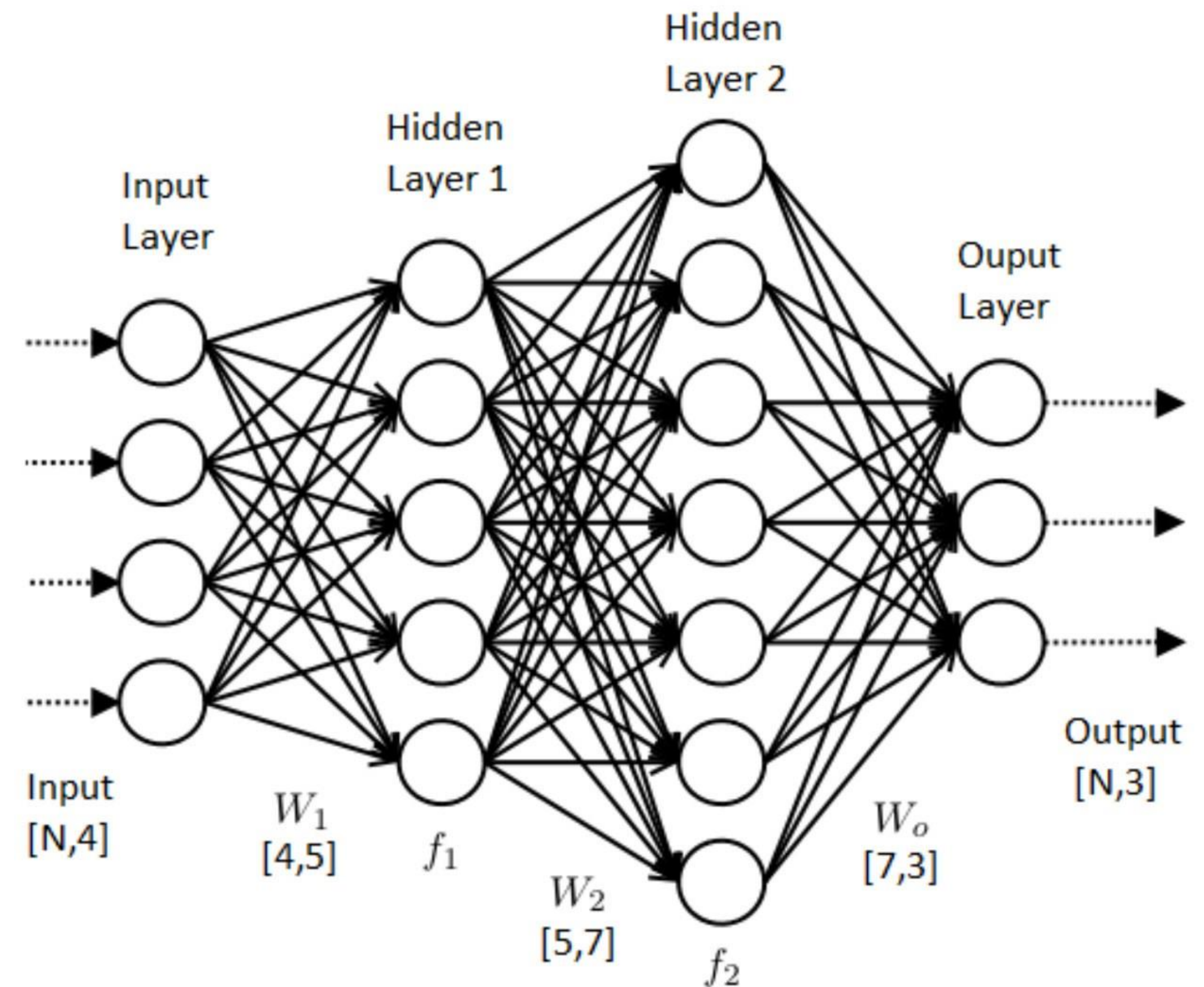
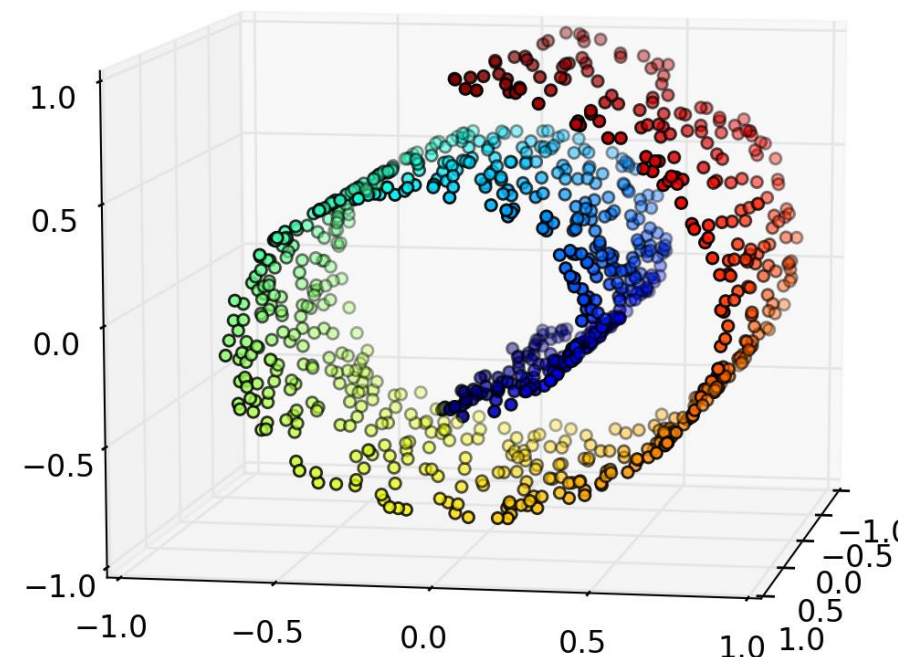
What is the difference between MFCC and filter bank ?

CLASSICAL ML APPROACHES

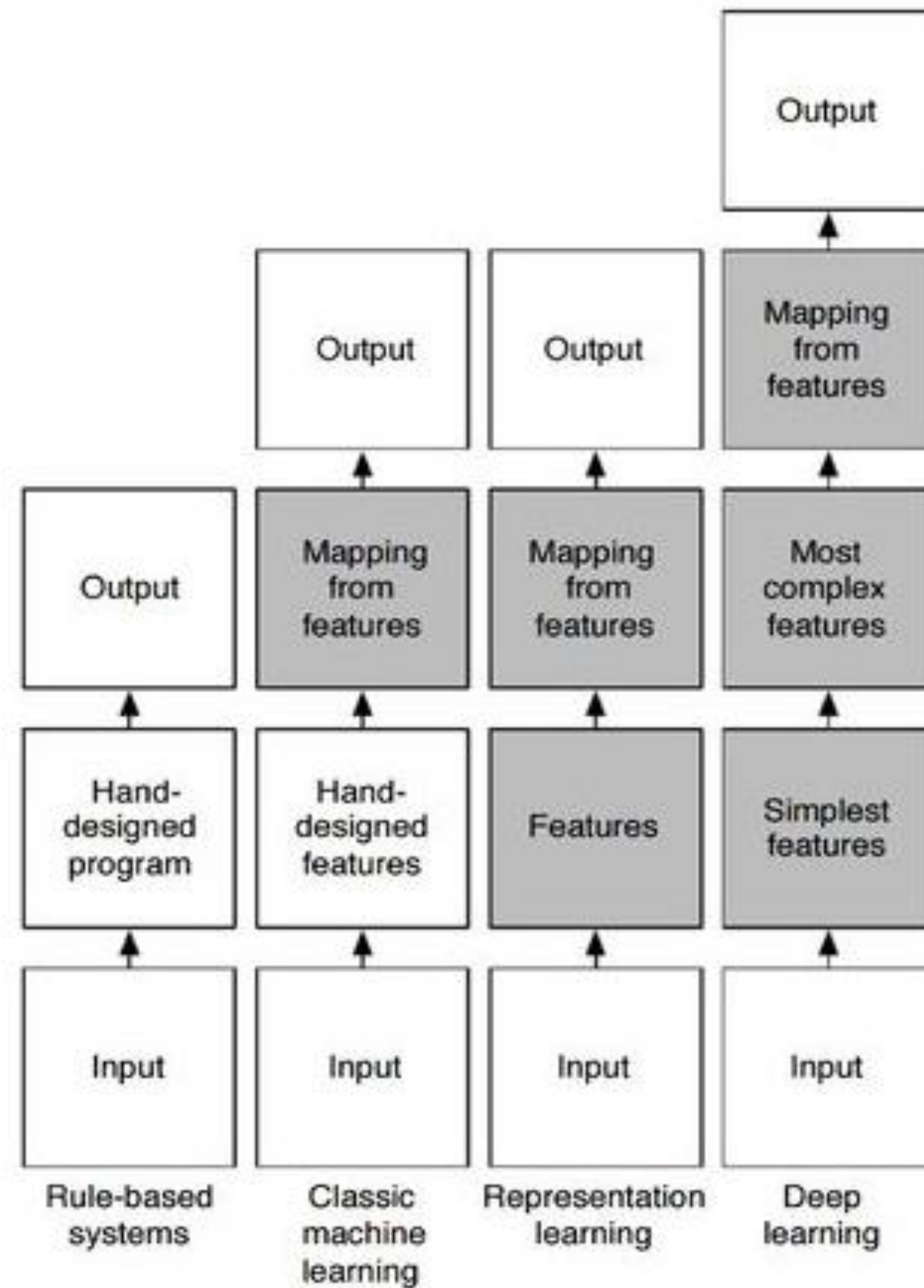
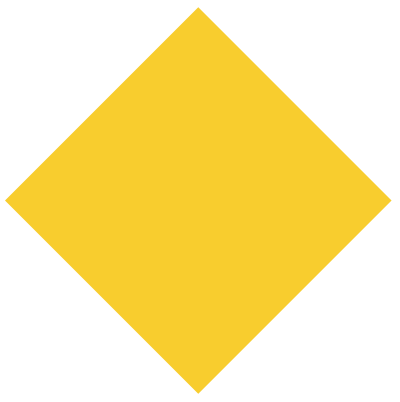


Mapping from features

- Preactivation (linear)
- Activation (non-linear)
- Non-linear operator enabling a linear separation
- The Manifold Hypothesis states that real-world high-dimensional data lie on low-dimensional manifolds embedded within the high-dimensional space.



CLASSICAL ML APPROACHES



J. Humphrey, J. P. Bello and Y. Lecun.
Moving beyond feature design: Deep
architectures and automatic feature
learning in music informatics. ISMIR 2012.

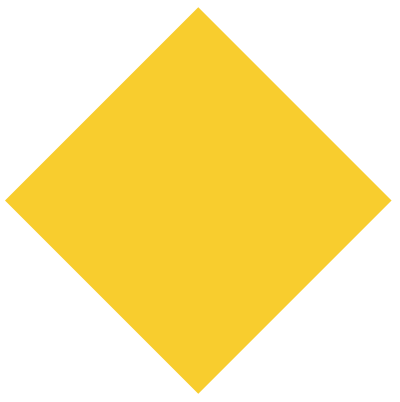
SIGNAL, AUDIO, SPEECH ENCODING.

Thank you for your attention.

Referencences:

- Geoffroy Peeters, Telecom Paris Tech
- Ian Goodfellow, Deep learning, 2018

ANNEXE 1: FOURIER TRANSFORM



The Fourier transform (FT) is an integrable function $\hat{f}: \mathbb{R} \rightarrow \mathbb{C}$, defined as following (one among several convention):

$$F[\omega] = \left[\int_{-\infty}^{\infty} f(t) e^{-j\omega t} dt \right]$$

with pulsation: $\omega = 2\pi \cdot f_{req}$

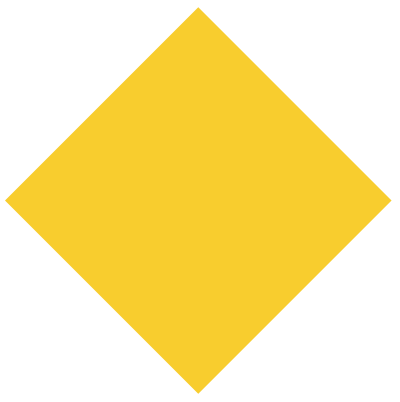
The Fourier inverse function is the following:

$$f(t) = \int_{-\infty}^{\infty} F[\omega] e^{j\omega t} d\omega$$

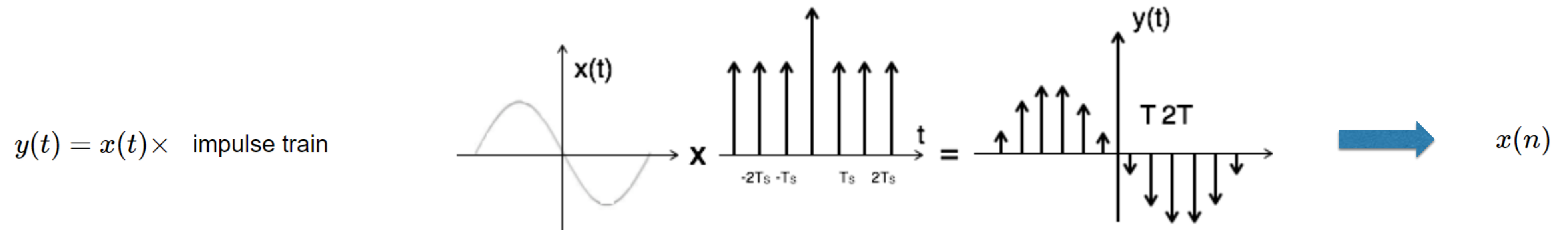
Usefull properties (let be x, y real functions):

- Linearity $ax(t) + by(t) \xleftrightarrow{\text{F.T}} aX(\omega) + bY(\omega)$
- FT of an impulse function $\delta(\omega) = 1$
- Time scaling $x(at) \xleftrightarrow{\text{F.T}} \frac{1}{|a|} X\left(\frac{\omega}{a}\right)$
- Time shifting / frequency shifting $x(t - t_0) \xleftrightarrow{\text{F.T}} e^{-j\omega t_0} X(\omega)$ and $e^{j\omega_0 t} \cdot x(t) \xleftrightarrow{\text{F.T}} X(\omega - \omega_0)$
- Multiplication and convolution $x(t) \cdot y(t) \xleftrightarrow{\text{F.T}} X(\omega) * Y(\omega)$ and $x(t) * y(t) \xleftrightarrow{\text{F.T}} \frac{1}{2\pi} X(\omega) \cdot Y(\omega)$
- Vectorized representation of FT values $X(\omega) = |X(\omega)| e^{j\theta(\omega)}$ where $\theta(\omega) = \arg X(\omega)$
and $|X(\omega)|, \theta(\omega)$ are called magnitude and phase spectrum of $X(\omega)$.

ANNEXE 2: SIGNAL SAMPLING



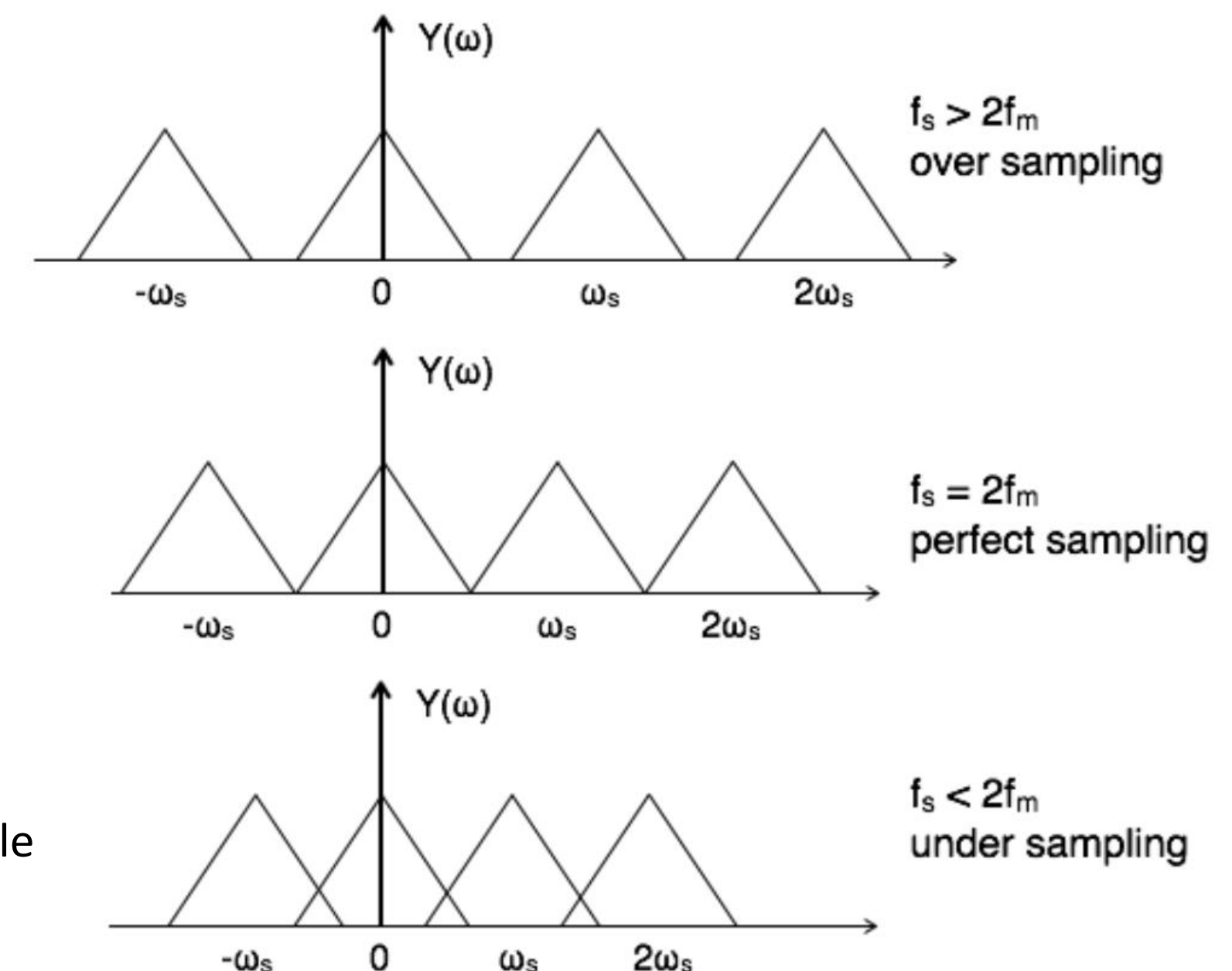
Sampling of a continuous signal into a discrete signal:



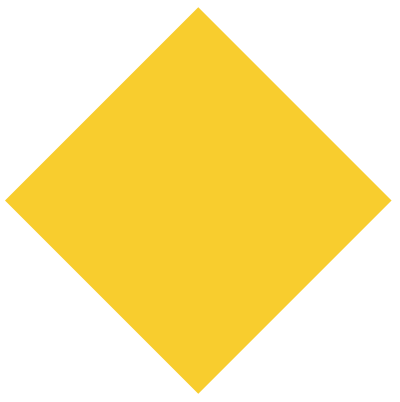
Nyquist theorem: A continuous time signal can be represented in its samples and can be recovered back when sampling frequency f_s is greater than or equal to the twice the highest frequency component of message signal. i. e.

$$f_s \geq 2f_m$$

frequency overlap here is called aliasing



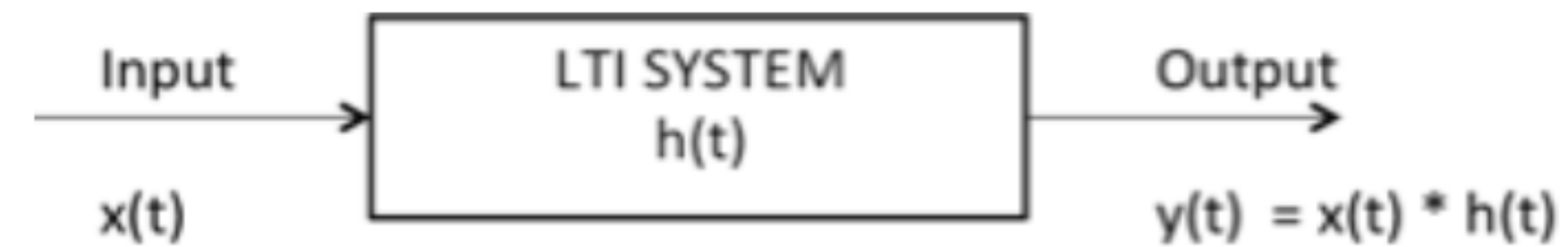
ANNEXE 3: CONVOLUTION & CORRELATION



Convolution of x by H (or H by x):

$$y(t) = x(t) * h(t)$$

$$= \int_{-\infty}^{\infty} x(\tau)h(t - \tau)d\tau$$



Property of Fourier Transform (see annexe 1)

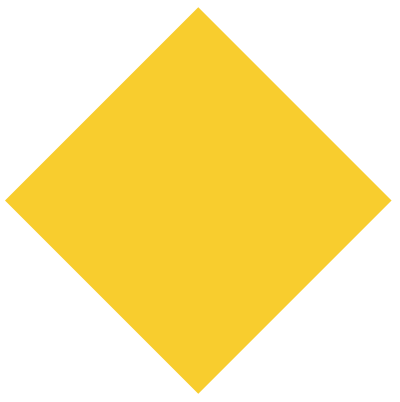
$$x(t) \cdot y(t) \xrightarrow{\text{F.T}} X(\omega) * Y(\omega)$$

$$x(t) * y(t) \xrightarrow{\text{F.T}} \frac{1}{2\pi} X(\omega) \cdot Y(\omega)$$

Correlation of 2 signals x_1 and x_2 :

$$\int_{-\infty}^{\infty} x_1(t)x_2(t - \tau)dt$$

ANNEXE 4: Z-TRANSFORM AND FILTERING



Definition of Z-Transform:

$$X(z) = \sum_{n=-\infty}^{\infty} x[n]z^{-n} \quad \text{with } z = re^{j\omega}$$

Possibility to create y signal that depends on x(n), x(n-1), x(n-2), ... and y(n-1), y(n-2), y(n-3),

This operation of y creation is called x filtering. And the Z-transform helps to characterize this filter (linear and invariant in time):

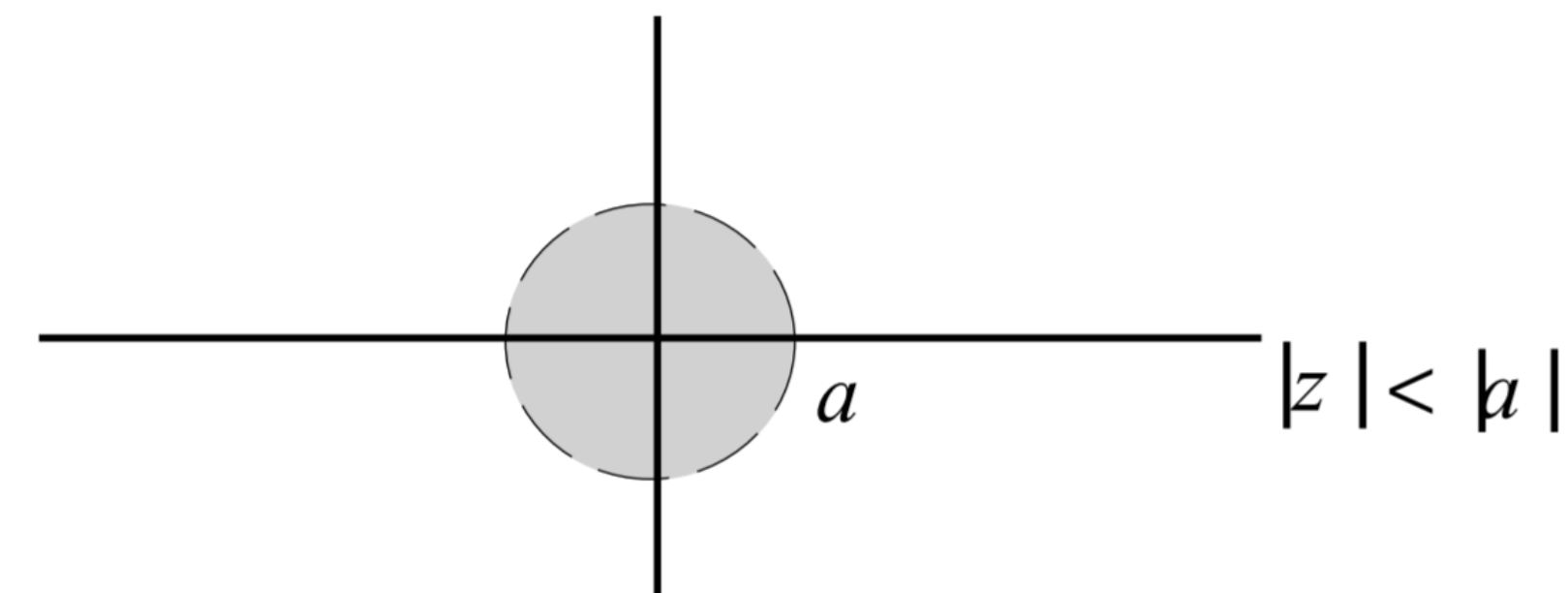
$$\sum_{l=-N}^N a_l y[n-l] = \sum_{k=-M}^M b_k x[n-k]$$

$$\sum_{l=-N}^N a_l z^{-l} Y(z) = \sum_{k=-M}^M b_k z^{-k} X(z)$$

$$H(z) = \frac{Y(z)}{X(z)} = \frac{\sum_{k=-M}^M b_k z^{-k}}{\sum_{l=-N}^N a_l z^{-l}}$$

Specific properties of H:

- if H denominators (complex) zeros have a modulus greater than 1, H is not stable.
- Resonances at pulsation close to argument of denominators zeros



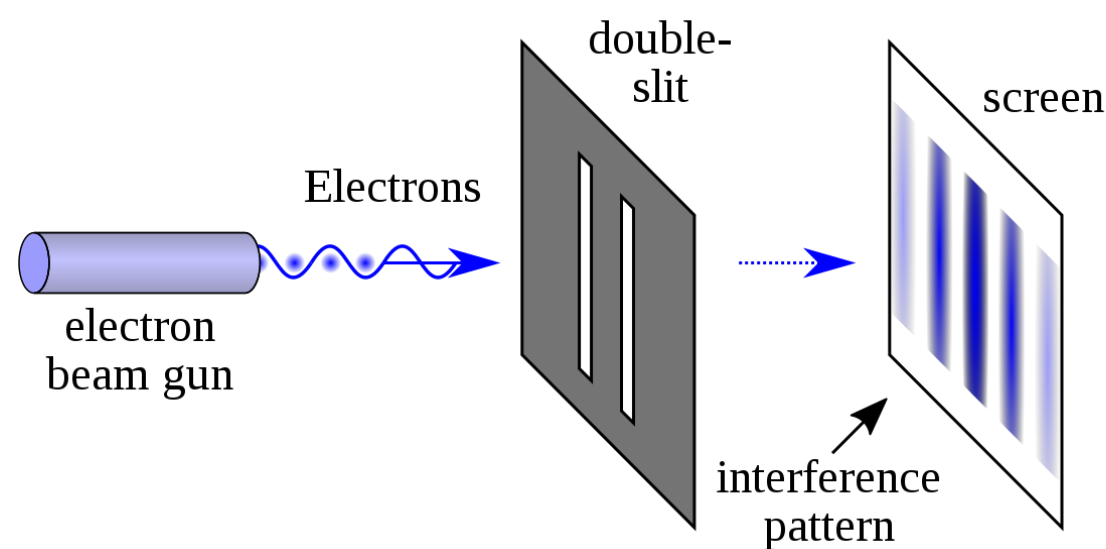
ANNEXE 5: WAVES INTERFERENCES & WAVEGUIDES

Sound wave

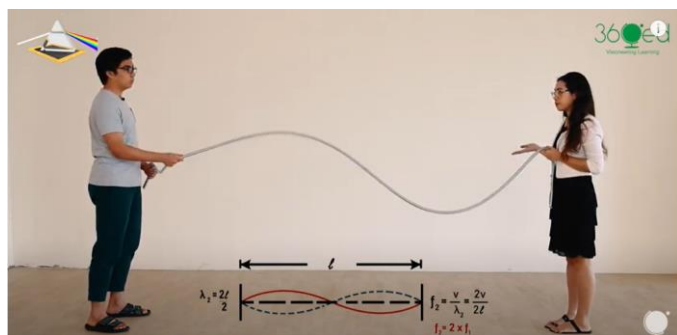
Propagation of a local oscillation of air pressure in a medium (air for instance).
It takes the form of longitudinal or transverse (for solids only) wave.

Wave interferences

Example of light wave interference: Young experiment.



Harmonics and wave guides



If λ is the wavelength (depending on frequency f) and v is sound speed (340 m/s): $\lambda = \frac{v}{f}$

Frequency	Order	Name	Standing wave representation
$1 \times f = 440 \text{ Hz}$	$n = 1$	1st harmonic (fundamental tone)	
$2 \times f = 880 \text{ Hz}$	$n = 2$	2nd harmonic (1st overtone)	
$3 \times f = 1320 \text{ Hz}$	$n = 3$	3rd harmonic (2nd overtone)	
$4 \times f = 1760 \text{ Hz}$	$n = 4$	4th harmonic (3rd overtone)	

