# NNs for audio

air pressure



$t$



$2\pi$ sec. $\rightarrow$ 1 osc.

1 sec. $\rightarrow \frac{1}{2\pi}$ osc.

Human perception:

16 Hz $\rightarrow$ 20 kHz

Wav-file:

bit rate    44.1 kHz

data width:  8-bit or 16-bit

number of channels (mono or stereo)

$$B = \frac{1}{2\pi} Hz$$

$$\frac{1}{2B} = \frac{2\pi}{2} = \pi$$

44 kHz $\rightarrow$ 16 kHz

# Th (Nyquist – Shannon)

$f(t)$ has largest freq. $B$ Hz $\Rightarrow$

$f(t)$ can be restored from sampled signal, if it is sampled $< \dfrac{1}{2B}$ samples per sec.

$$f(t) = A \sin(\omega t + \varphi)$$
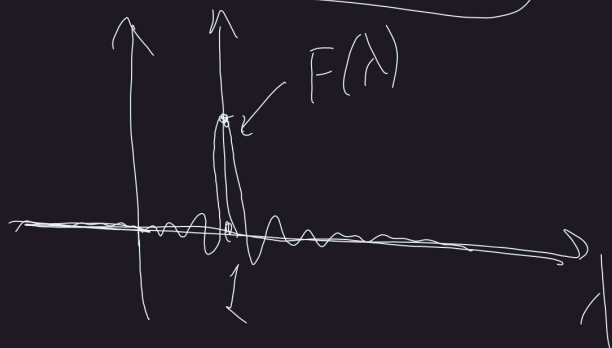
$A$ — amplitude

$\omega$ — freq.

$\varphi$ — phase

$$\boxed{f(t) = \sin(t)}$$

$$f(t) \longrightarrow A, \omega, \varphi$$

$$F(\lambda) := \int_{-\infty}^{+\infty} f(t) \sin(\lambda t)\, dt =$$

$$F(\lambda) = \int_0^T f(t)\sin(\lambda t)\, dt \approx \frac{1}{T}\sum_{n=0}^{T} f\left(\frac{n}{T}\right)\sin\left(\lambda \frac{n}{T}\right)$$

$= \delta(\lambda - 1)$

$$f(t) = \sin(t + \varphi)$$

$$F(\lambda) = \int f(t) \sin(\lambda t) dt$$



$$\leftarrow |F(\lambda)|^2 \quad \lambda$$

$$\in \mathbb{C}$$

$$F(\lambda) = \int f(t) \exp(-i\lambda t) dt$$

$$A = |F(\lambda)|^2$$

$$\varphi = \text{Arg } F(\lambda)$$

$$\begin{cases} F_1(\lambda) = \int f(t) \sin(\lambda t) dt \\ F_2(\lambda) = \int f(t) \cos(\lambda t) dt \end{cases}$$

$$A = F_1^2(\lambda) + F_2^2(\lambda)$$

doesn't depend on $\varphi$

$$\varphi = \arctan\left(\frac{F_1(\lambda)}{F_2(\lambda)}\right)$$
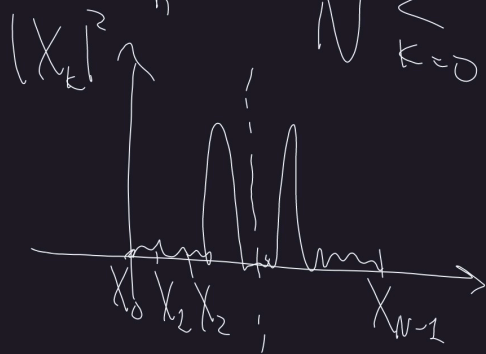
# Discrete Fourier Transform:

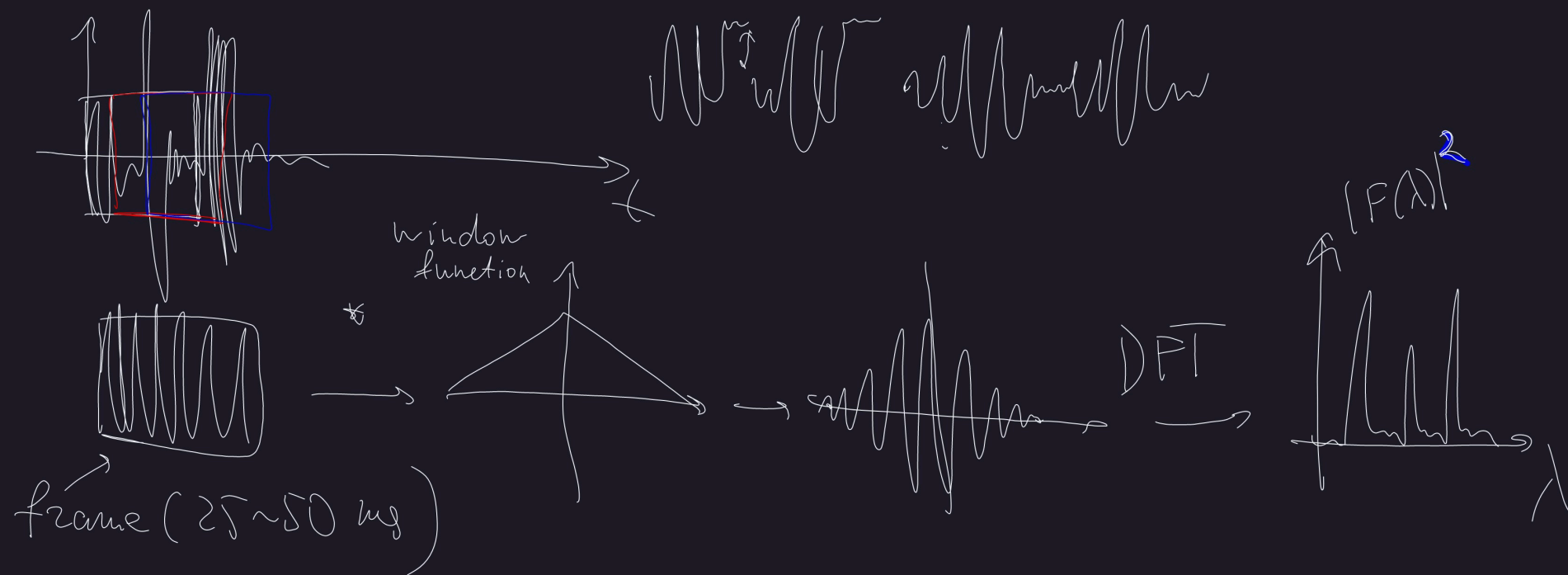$$x_0, x_1, \ldots, x_{N-1} \xrightarrow{\text{DFT}} X_0, X_1, \ldots, X_{N-1}$$
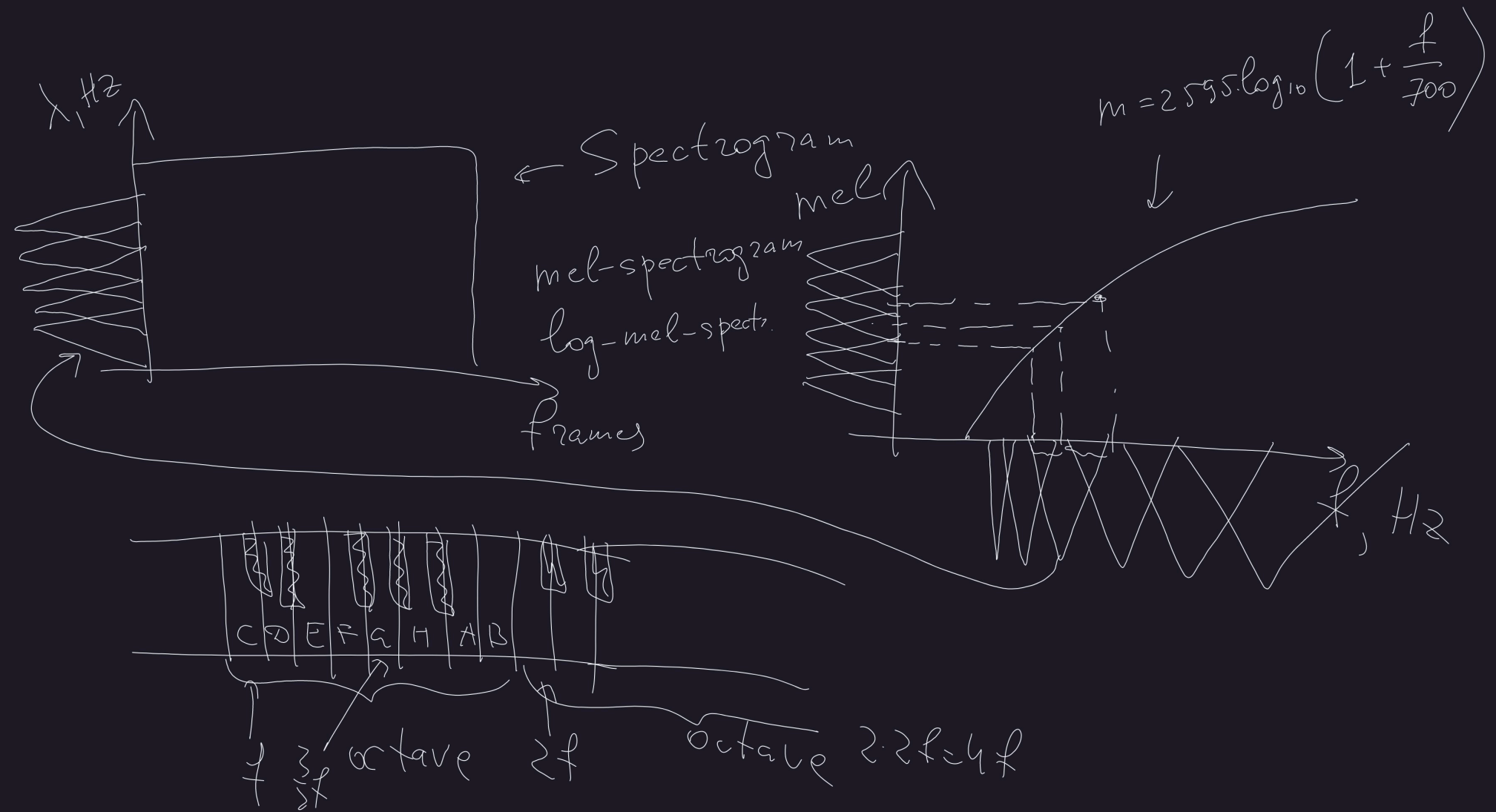
$$X_k = \sum_{n=0}^{N-1} x_n \exp\left(-i 2\pi \frac{k}{N} n\right) \quad \forall k = 0, 1, \ldots, N-1$$

$$x_n = \frac{1}{N} \sum_{k=0}^{N-1} X_k \exp\left(i 2\pi \frac{k}{N} n\right) \quad \forall n = 0, 1, \ldots, N-1$$

## Complexity of DFT

Direct comp. $O(N^2)$

FFT $O(N \log N)$

ASR: Audio → Spectrogram $\xrightarrow{NN}$ Text



window function

frame (25~50 ms)

DFT

$|F(\lambda)|^2$

$x, Hz$

← Spectrogram

mel-spectrogram

log-mel-spectr.

Frames

mel

$m = 2595 \cdot \log_{10}\left(1 + \frac{f}{700}\right)$

$f, Hz$

C D E F G H A B

$f$ $\frac{3}{2}f$ octave $2f$ octave $2 \cdot 2f = 4f$

MFCC    spectrogram

Mel_frequency_cepstrum_coefs

spectrum

$$MFCC = \left| F \left( \log \left| F(f(t)) \right|^2 \right) \right|^2$$

Audio →

spectr.

$$x_t = z_t \cdot \sigma(z_{<t}) + \mu(z_{<t}) \qquad \left( \triangle \square \right) \qquad = \sum_t \log p(z_t) - \sum_t \log \sigma(z_{<t})$$

$$z \sim p(z) = N(z|0,I)$$

$$\log p(x) = \log p(z) - \log \left| \det \frac{\partial x}{\partial z} \right| = \sum_t \log p(z_t) - \log \prod_{t=1}^{T} \sigma(z_{<t}) =$$

$$KL\left(p_S(x|\theta) \;\|\; p_T(x)\right) \longrightarrow \min_\theta$$

$$\| \\[4pt]
\mathbb{E}_{p_S(x|\theta)} \log \frac{p_S(x|\theta)}{p_T(x)} = \mathbb{E}_{p_S(x|\theta)}\left( \sum_t \log p(z_t) - \sum_t \log \sigma(z_{<t},\theta) \right) -$$

$$- \mathbb{E}_{p_S(x|\theta)} \sum_t \log p_T(x_t|x_{<t}) = \left\{ \begin{array}{l} z \sim p(z),\; x = f(z) \\[6pt] \mathbb{E}_{p(x)}\, g(x) = \mathbb{E}_{p(z)}\, g(f(z)) \end{array} \right\} =$$

$$= \mathbb{E}_{p(z)}\left( \underbrace{\sum_t \log p(z_t)}_{\text{const}} - \sum_t \log \sigma(z_{<t},\theta) \right) - \mathbb{E}_{p(z)} \sum_t \log p_T\left( \underbrace{z_t \cdot \sigma(z_{<t}) + \mu(z_{<t})}_{x_t} \Big| \ldots \right)$$