

# Deep Learning for Audio

## Lecture 1

Pavel Severilov

AI Masters

2023

# Outline

1. Voice Technologies: tasks
2. Speech Recognition: history, basics
3. Speech Synthesis: history, basics
4. Sound: characteristics
5. Fourier Transform
6. Spectrograms

# Organisation

1. 11 lectures + seminars
2. 4 homeworks (2 points for every)
3. Final test (2 points)
4. Grades:
  - ▶ 5: 8-10 points
  - ▶ 4: 6-7.9 points
  - ▶ 3: 4-5.9 points
5. Discussion: telegram chat (contact @severilov and @Oorgien)

# Outline

1. Voice Technologies: tasks
2. Speech Recognition: history, basics
3. Speech Synthesis: history, basics
4. Sound: characteristics
5. Fourier Transform
6. Spectrograms

# Voice Technologies: Applications

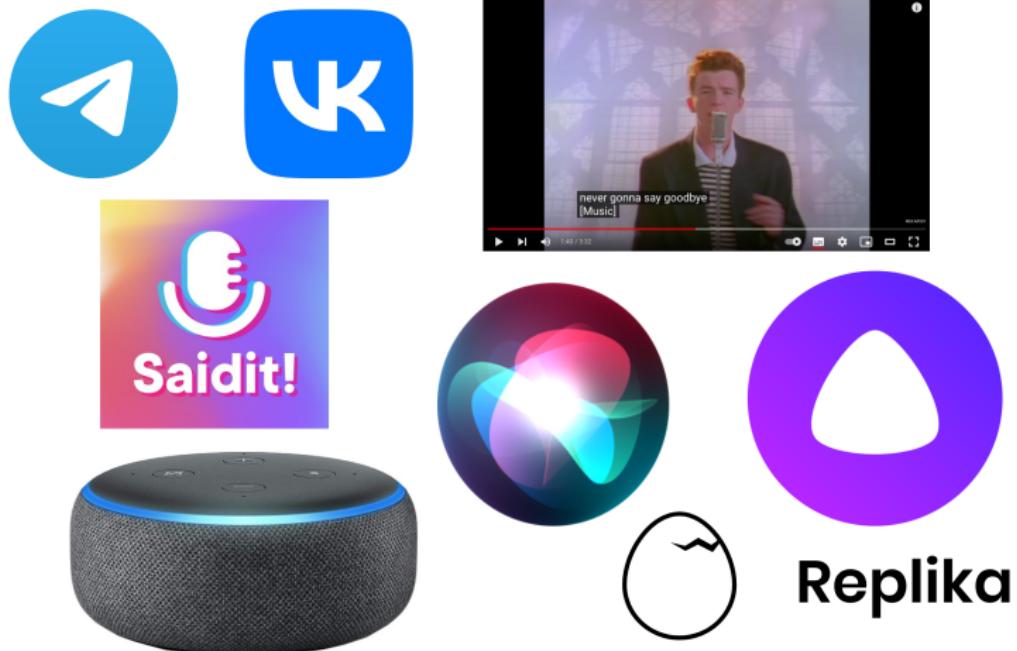
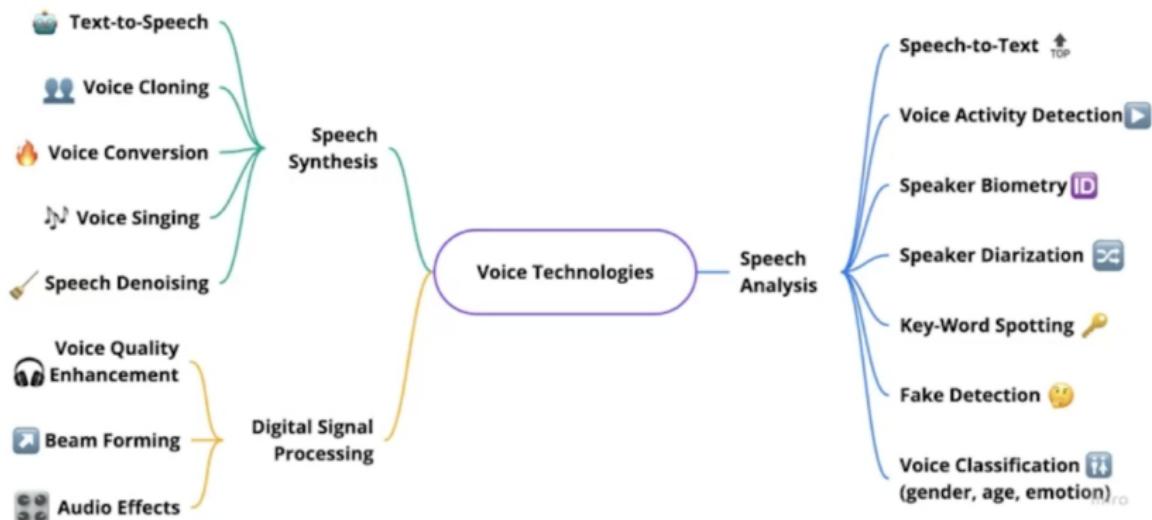
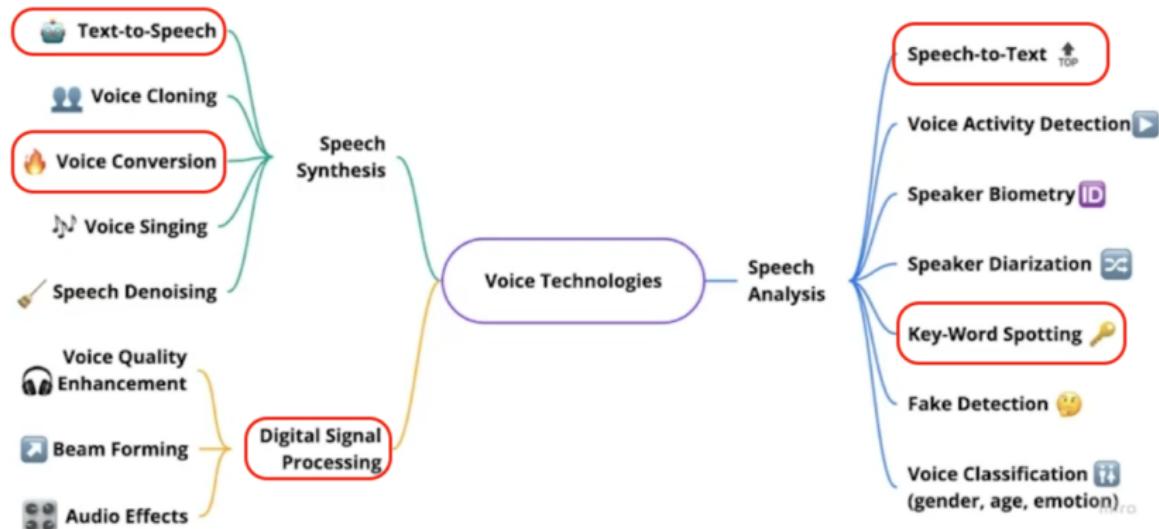


Figure: Siri, Amazon Alexa, Alisa, Replika, Telegram, VK, YouTube

# Voice Technologies, Tasks: Mind Map



# Voice Technologies, Tasks: Course



# Outline

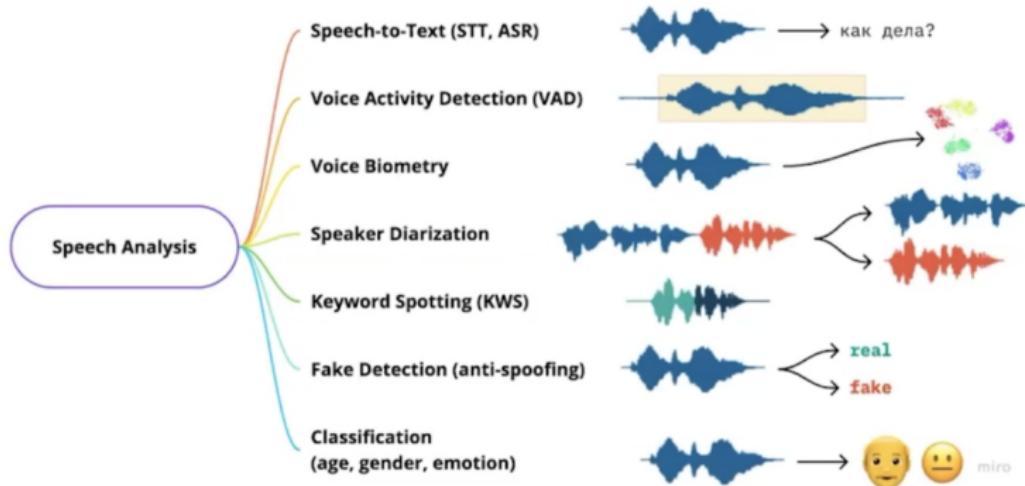
1. Voice Technologies: tasks
2. Speech Recognition: history, basics
3. Speech Synthesis: history, basics
4. Sound: characteristics
5. Fourier Transform
6. Spectrograms

# History of Speech Recognition



- ▶ **50's:** 1952, Bell Laboratories, "Audrey" system, could recognize single voice speaking digits
- ▶ **60's:** 1961, IBM, "Shoebox", understood 16 words in English
- ▶ **70's:** DARPA, understood over 1000 words (Siri spin-out)
- ▶ **80's:** using HMM, understood several thousand words
- ▶ **90's:** became faster because of processors
- ▶ **00's-10's:** ML, DL, Big Data, GPUs

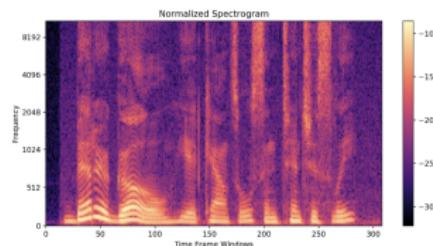
# Speech Analysis Tasks: Mind Map



# Speech Recognition & Deep Learning: Idea



Raw Audio



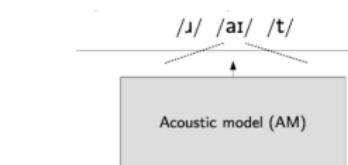
Feature: Spectrogram



"right"



Language Model



Acoustic model: phonemas

# Outline

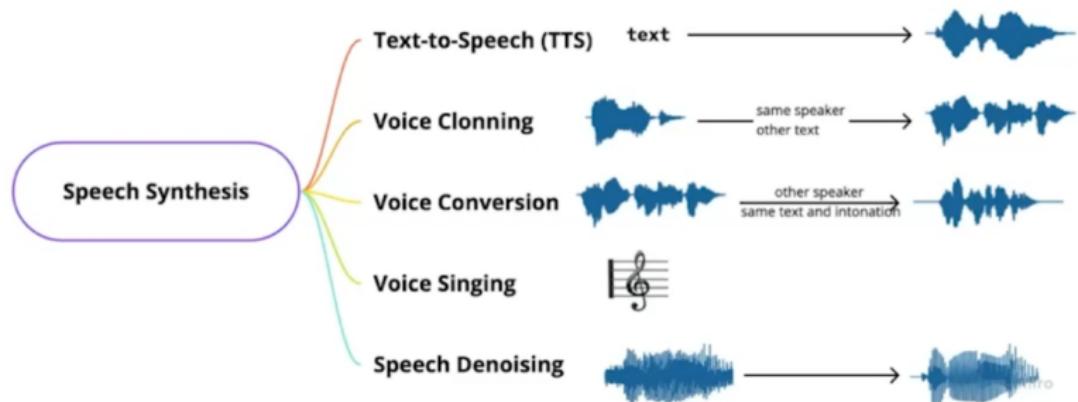
1. Voice Technologies: tasks
2. Speech Recognition: history, basics
3. Speech Synthesis: history, basics
4. Sound: characteristics
5. Fourier Transform
6. Spectrograms

# History of Speech Synthesis



- ▶ **30's:** 1939, Bell Laboratories, "Voder",
- ▶ **80's:** Format-based on rules, Atari/Sega
- ▶ **90's-00's:** Concatenative synthesis
- ▶ **10's:** ML, DL, Big Data, GPUs

# Speech Synthesis Tasks: Mind Map



# Speech Synthesis & Deep Learning: Idea

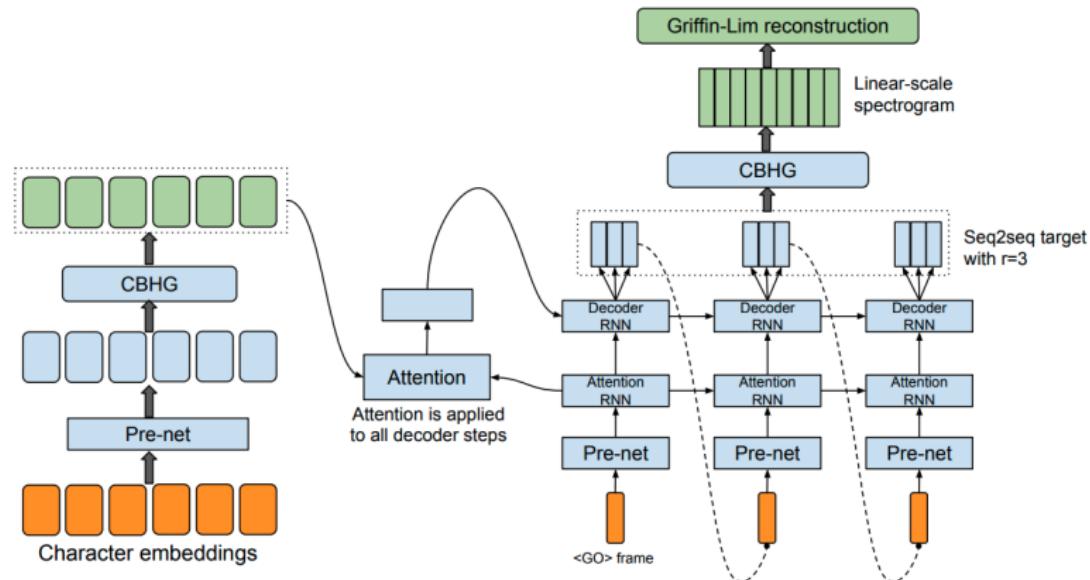


Figure: Example of Deep Learning approach to speech synthesis:  
encoder-decoder structure with recurrent parts

Wang, Yuxuan et al. "Tacotron: Towards End-to-End Speech Synthesis." INTERSPEECH (2017), Google Inc.

# Outline

1. Voice Technologies: tasks
2. Speech Recognition: history, basics
3. Speech Synthesis: history, basics
4. Sound: characteristics
5. Fourier Transform
6. Spectrograms

## Mind experiment: what is sound?



Figure: If a tree falls in the forest, and there's nobody around to hear, does it make a sound?

# Ear structure

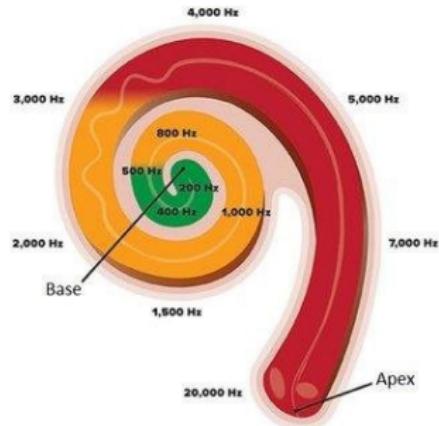
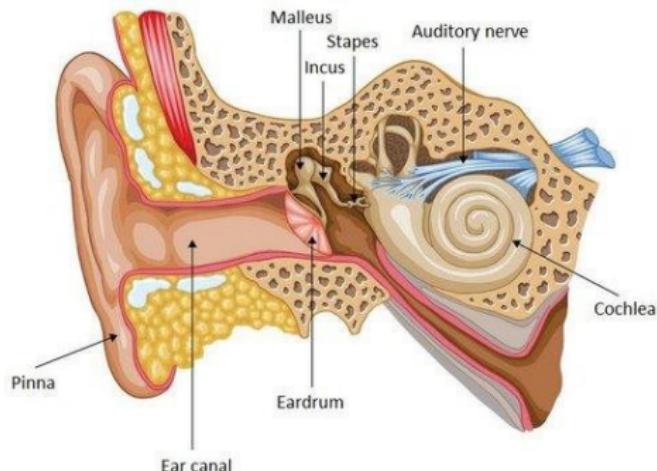


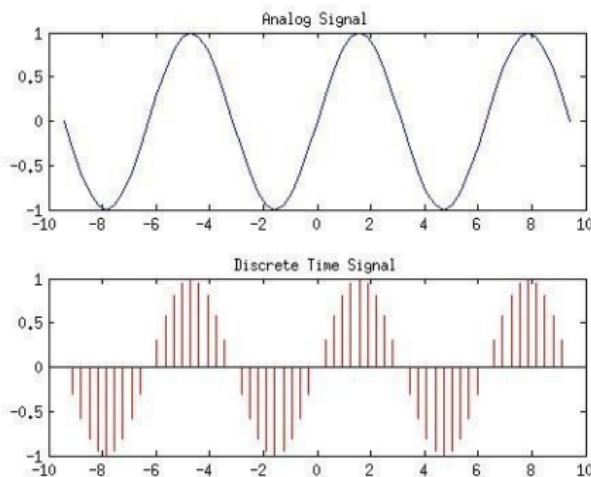
Figure: Fourier decomposition in our body

Fourier transform of a signal  $x(t)$

$$\mathcal{F}\{x(t)\} = X(j\omega) = \int_{-\infty}^{\infty} x(t)e^{-j\omega t} dt$$

## Sound types: analog vs digital

- ▶ For any  $T > 0$ , we may sample a **Continuous-Time** (CT) signal  $x(t)$  to generate the **Discrete-Time** (DT) signal  $x[n] = x(nT)$
- ▶ DT: signal  $x(t)$  is evaluated at uniformly spaced points on the t-axis. The number  $T$  is the **sampling period**
- ▶ **Sampling frequency:**  $f_s = \frac{1}{T}$ , [Hertz or samples/sec.]



## Sound characteristics

- ▶ Signal  $x(t)$
- ▶ Signal energy  $\int_{-\infty}^{\infty} |x(t)|^2 dt$  (used for normalizing signals, augmentations)
- ▶ Sample rate (SR) – number of audio samples per one second
  - ▶ 8 kHz or 16 kHz - standard for audio in telephony
  - ▶ 44.1 kHz - CD audio/Computer Audio
  - ▶ 48 kHz - DVD audio/Computer Audio
  - ▶ 96 kHz - High resolution Audio
- ▶ Number of channels – how many signals we record in parallel (mono: 1, stereo: 2)

# The Nyquist Theorem

- $\Sigma_{CT}$  – set of all CT signals  $x(t)$ ,  $\Sigma_{DT}$  – set of all DT signals  $x[n]$ . Procedure of sampling for a given sampling period  $T$ :

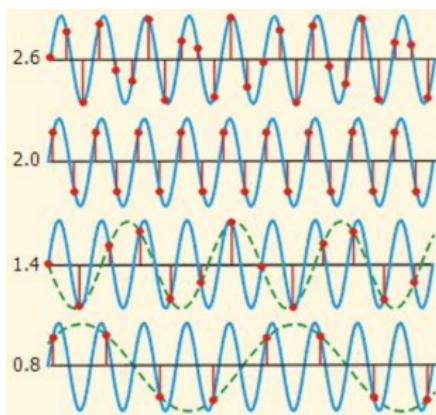
$$\Sigma_{CT} \longleftrightarrow \Sigma_{DT} \Leftrightarrow x(t) \mapsto x[n] = x(nT)$$

- CT signal  $x(t)$  is bandlimited if there exists  $\omega_B < \infty$  such that  $X_{CT}(j\omega) = 0$  for  $|\omega| > \omega_B$

## Nyquist Theorem

The sampling map is bijection on  $\Sigma_{\omega_B}$  iff  $\omega_s > 2\omega_B$ .

$\omega_s = 2\pi f_s = \frac{2\pi}{T}$  – radian sampling frequency



# Outline

1. Voice Technologies: tasks
2. Speech Recognition: history, basics
3. Speech Synthesis: history, basics
4. Sound: characteristics
- 5. Fourier Transform**
6. Spectrograms

# Fourier Transform: motivation

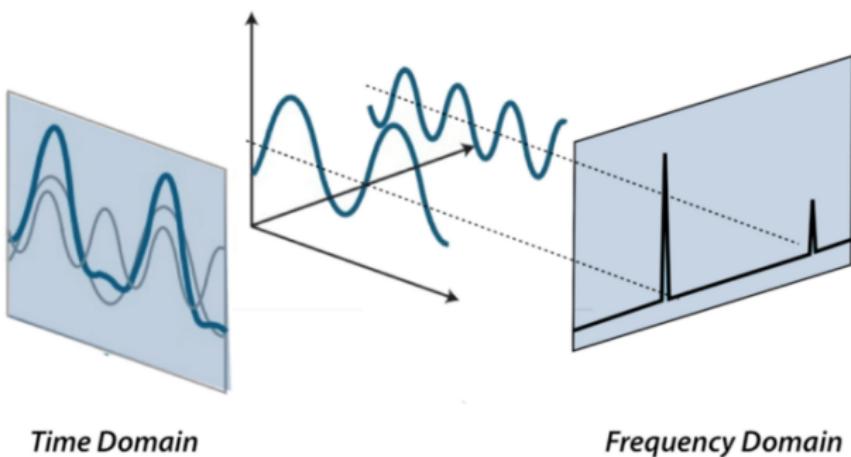


Figure: Fourier Transform transfer a signal from the **time domain** to the **frequency domain**

## Fourier Transform

- ▶ CT Fourier transform (CTFT) of a CT signal  $x(t)$  is

$$\mathcal{F}\{x(t)\} = X(j\omega) = \int_{-\infty}^{\infty} x(t)e^{-j\omega t} dt$$

- ▶ CT unit impulse  $\delta(t)$ ,  $\int_{-\infty}^{\infty} \delta(t)dt = 1$ ,
- ▶ Define signal by CT impulse  $x(t) = \sum_{n=-\infty}^{\infty} x[n]\delta(t - n)$
- ▶ Taking Fourier transform for DT signal:

$$\begin{aligned} X(j\omega) &= \int_{-\infty}^{\infty} \sum_{n=-\infty}^{\infty} x[n]\delta(t - n)e^{-j\omega t} dt \\ &= \sum_{n=-\infty}^{\infty} x[n]e^{-j\omega n} \int_{-\infty}^{\infty} \delta(t - n)e^{-j\omega(t-n)} dt \\ &= \sum_{n=-\infty}^{\infty} x[n]e^{-j\omega n} = \mathbf{M}x(t). \end{aligned}$$

## DTFT: example

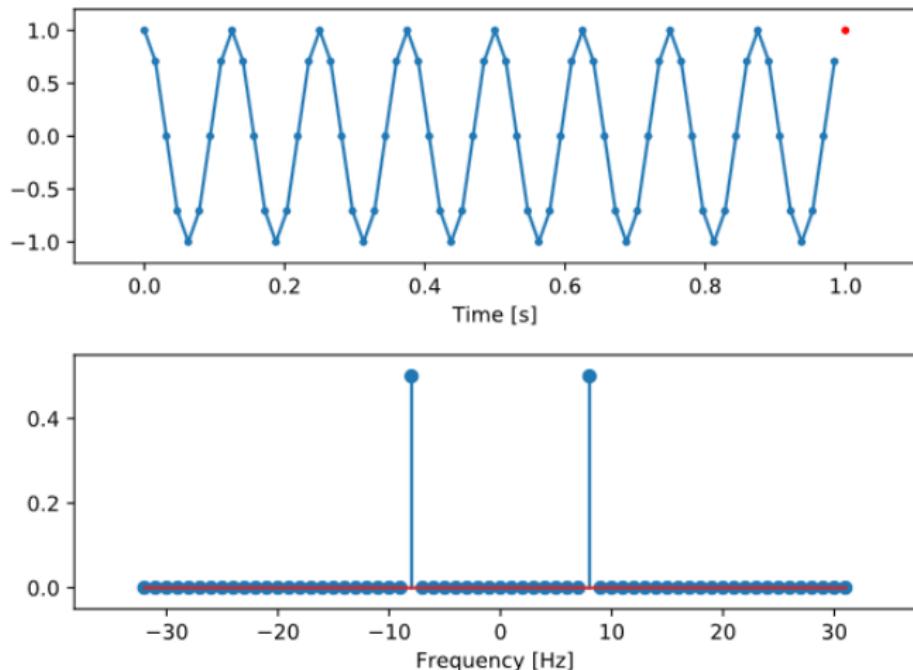
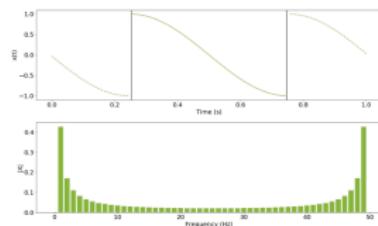


Figure: Example of DTFT for cosine signal

# Problems with Fourier Transform in real life

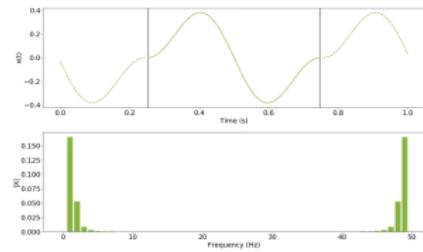
Sliced signal



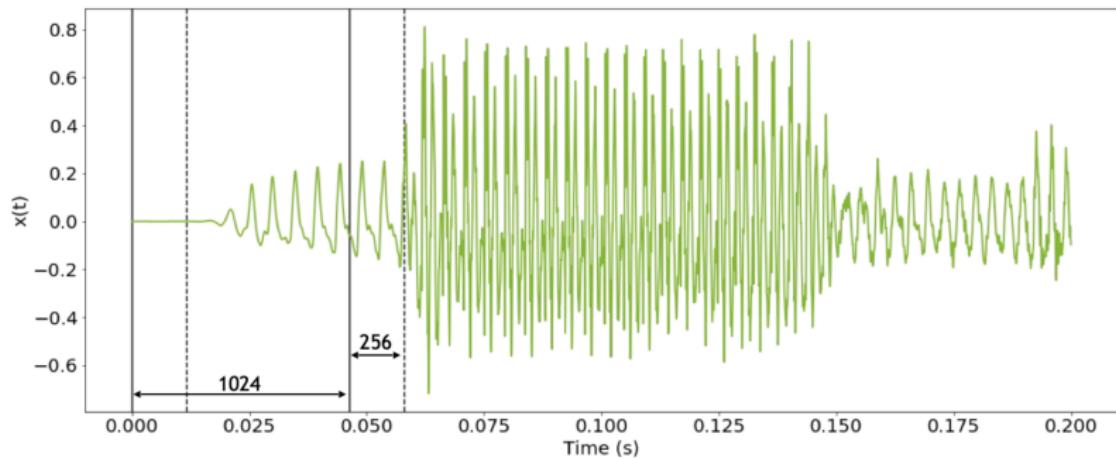
Window



Windowed signal



# Short-time Fourier Transform (STFT)



STFT (DTFT with Hamming window):

$$X[r, w] = \sum_{n=-\infty}^{\infty} w[r-n]x[n]e^{-j\omega n},$$

where  $w$  – window function,  $r$  – location of window along the time axis

# Outline

1. Voice Technologies: tasks
2. Speech Recognition: history, basics
3. Speech Synthesis: history, basics
4. Sound: characteristics
5. Fourier Transform
6. Spectrograms

# Spectrograms

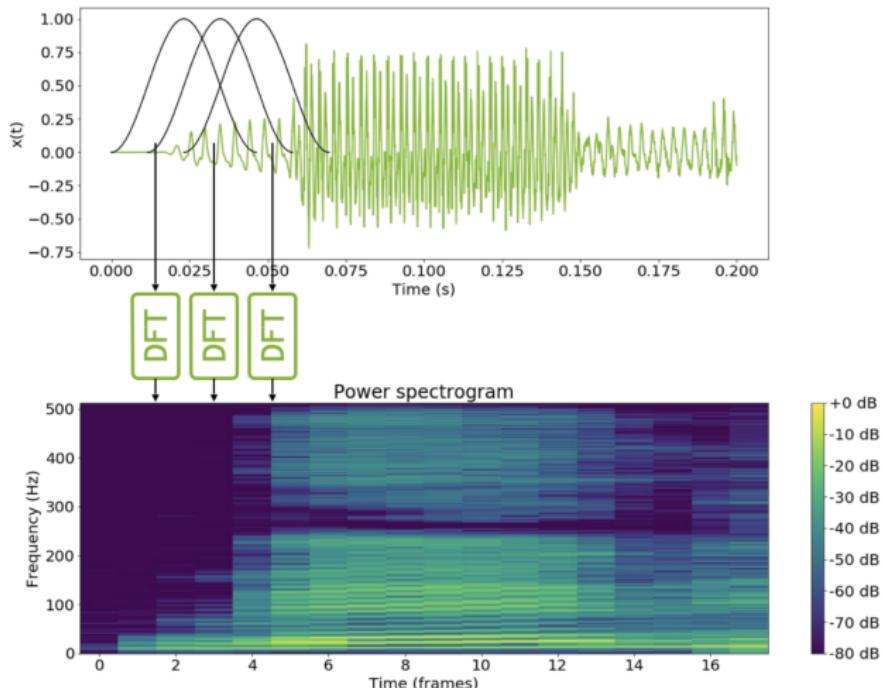


Figure: STFT+window Spectrogram

## Mel Scale

- ▶ Mel Scale: humans perceive sound on a log-scale, not linear
- ▶ A lot of formulas to convert  $f$  hertz into  $m$  mels. Popular example:

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

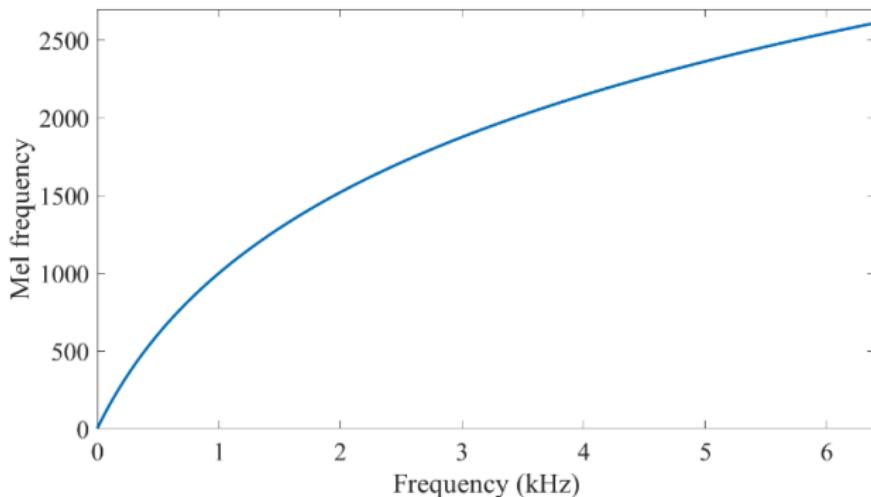


Figure: Mel scale

# Mel spectrogram

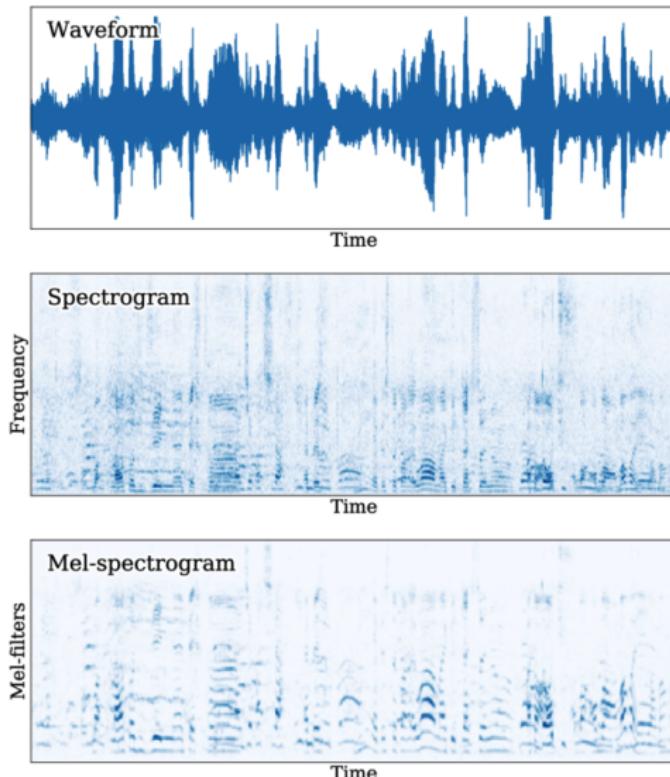


Figure: Waveform → STFT+window Spectrogram → STFT+window+Mel scale Spectrogram