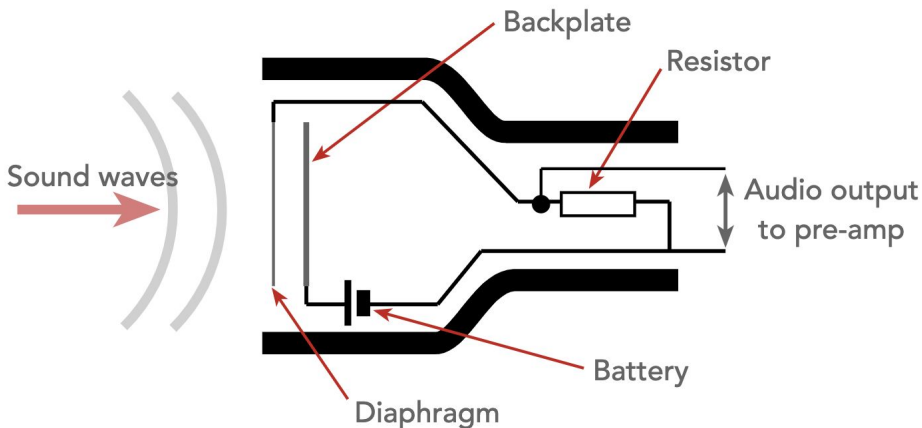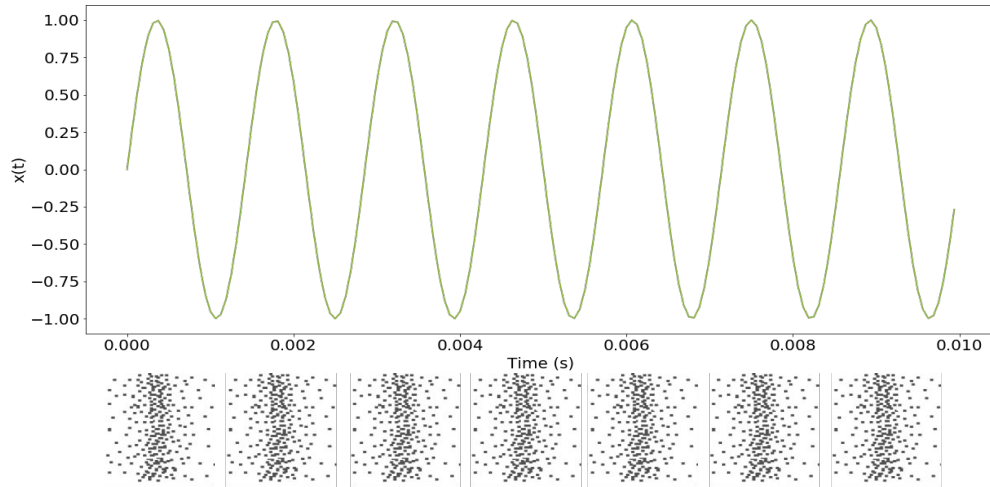# Introduction to DSP

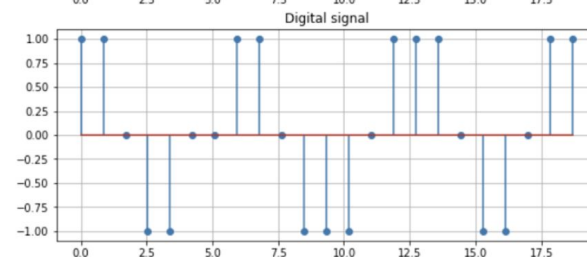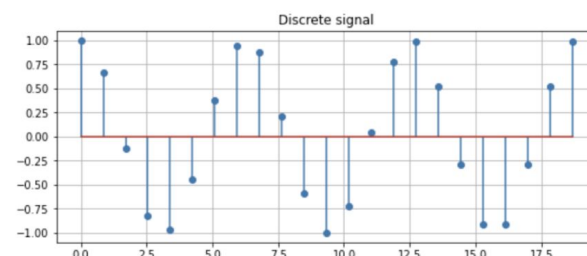Елена Кантонистова

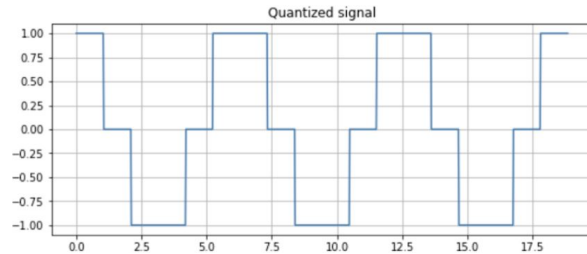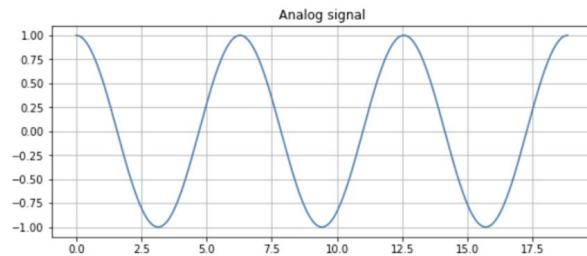По мотивам материалов Александра Марковича

# What is sound?

- **Sound wave** is the pattern of **oscillations** caused by the movement of energy traveling through the air

- **Microphone** picks up these air **oscillations** and converts them into electrical vibrations

- These **oscillations** are converted into an **analog** signal and then a **digital** signal
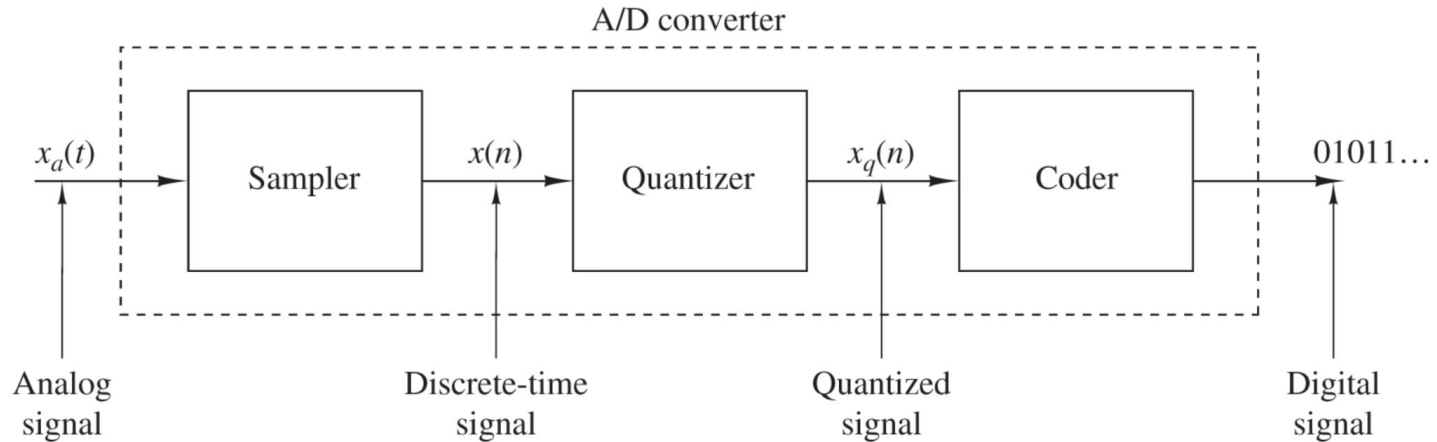
# How is sound stored in the computer?

- The **analog** signal is discretized, quantized and encoded

- An analog signal is **discretized** in that the signal is represented as a sequence of values taken at discrete points in time **t** with step **d**

- **Quantisation** of a signal consists in splitting the range of signal values into **N** levels in increments of **d** and selecting for each reference the level that corresponds to it

- Signal **encoding** is just a way of presenting the signal in a more compact form

https://github.com/markovka17/dla/blob/master/week01/dsp.ipynb
https://web.sonoma.edu/esee/courses/ee442/archives/sp2019/lectures/lecture09_pcm.pdf

# Analog-to-Digital Conversion

- Converting analog signals to a sequence of numbers having finite precision

- Corresponding devices are called A/D converters (ADCs)

# What other characteristics are there?

- **Sample rate (SR)** - number of audio samples per one second (e.g. 8 kHz, 22.05 kHz, 44.1 kHz)

- **Sample size** - number of bits per one sample (e.g. 8, 16, 25, 32 bits)

- **Number of channels** -- how many signals we record in parallel (e.g. mono(1), stereo(2))

**8000 Hz**
The international G.711 ☐ standard for audio used in telephony uses a sample rate of 8000 Hz (8 kHz). This is enough for human speech to be comprehensible.

**44100 Hz**
The 44.1 kHz sample rate is used for compact disc (CD) audio. CDs provide uncompressed 16-bit stereo sound at 44.1 kHz. Computer audio also frequently uses this frequency by default.

**48000 Hz**
The audio on DVD is recorded at 48 kHz. This is also often used for computer audio.
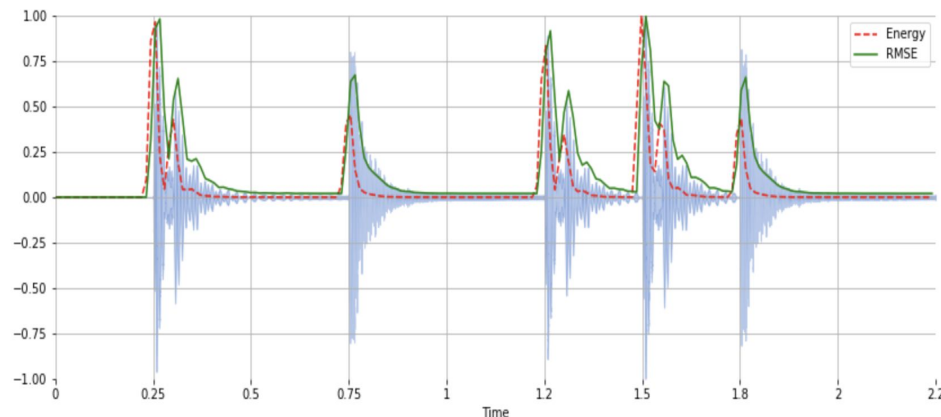
**96000 Hz**
High-resolution audio.

**192000 Hz**
Ultra-high resolution audio. Not commonly used yet, but this will change over time.
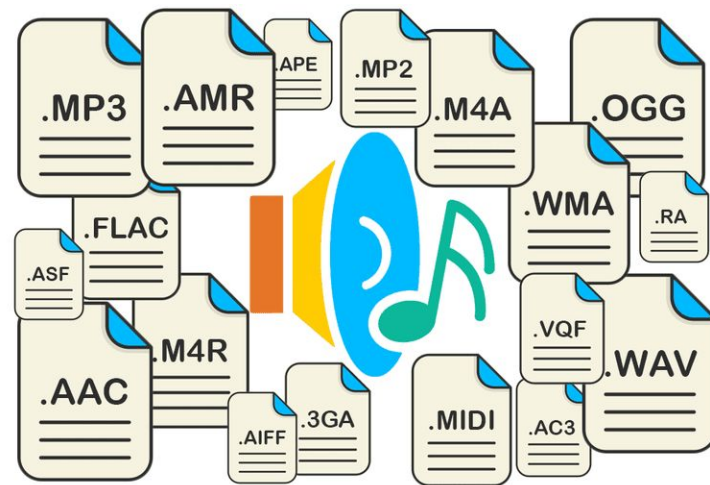
# What other characteristics are there?

- Assume **f(n)** is our signal where **n** is time

- Power of signal is   $f^2(n)$

- Energy of signal is   $\sum f^2(n)$

- In practice estimated by some **window**

- Energy in **decibels:**   $10 \log_{10} E$

- $\text{SNR}_{dB} = 10 \log_{10} \dfrac{E_{\text{signal}}}{E_{\text{noise}}}$

# What about audio formats?

- Non-compressed formats: **WAV, AIFF, etc.**

- Lossless compression(2:1) : **FLAC, ALAC, etc.**

- Lossy compression(10:1) : **MP3, Opus, etc**

- **Bit rate** measure a degree of compression. Number of bit that are conveyed or processed per **unit of time**.
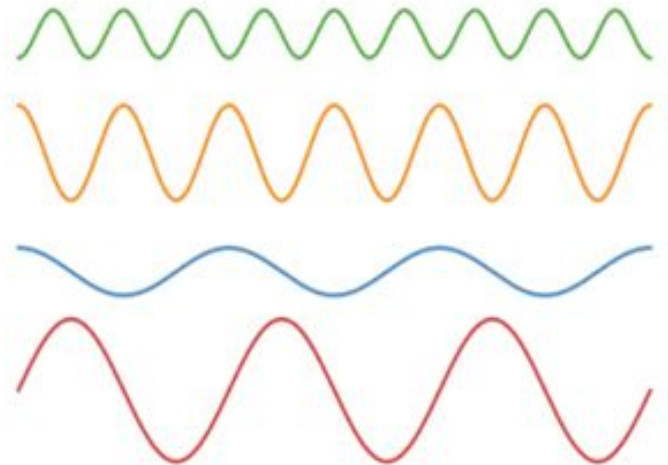
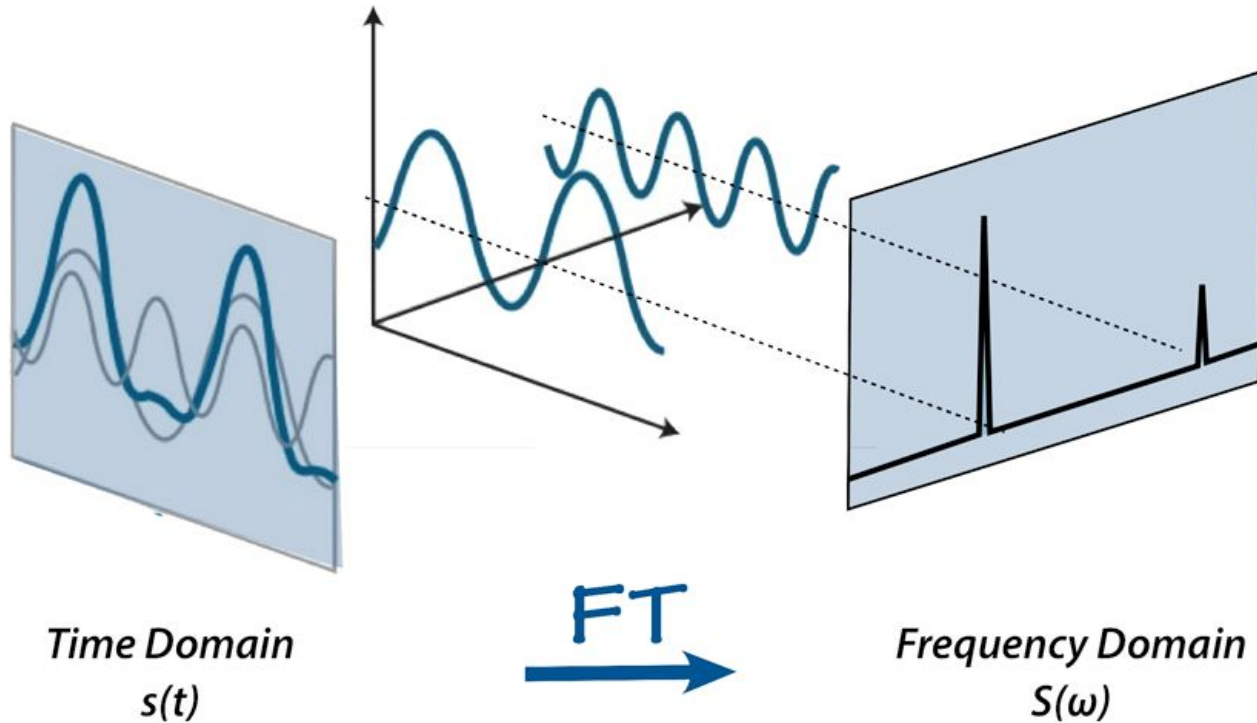# Fourier Transform



Complex sound wave

Find Fourier Series

Combination of simple sound waves

# Fourier Transform



**Time Domain**
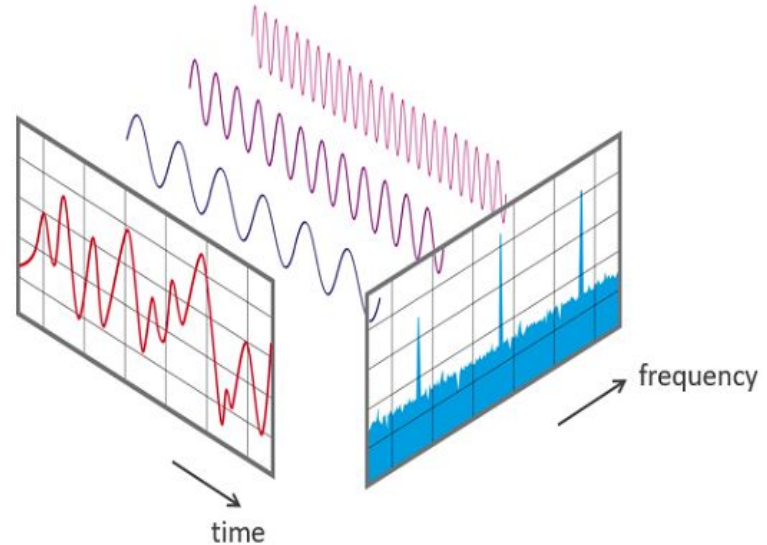*s(t)*

FT

**Frequency Domain**
*S(ω)*

# Fourier Transform

- The **Fourier transform(FT)** is a mathematical formula that allows us to decompose a signal into its individual **frequencies** and the frequency's **amplitude**

- FT transfer a signal from the **time domain** to the **frequency domain**

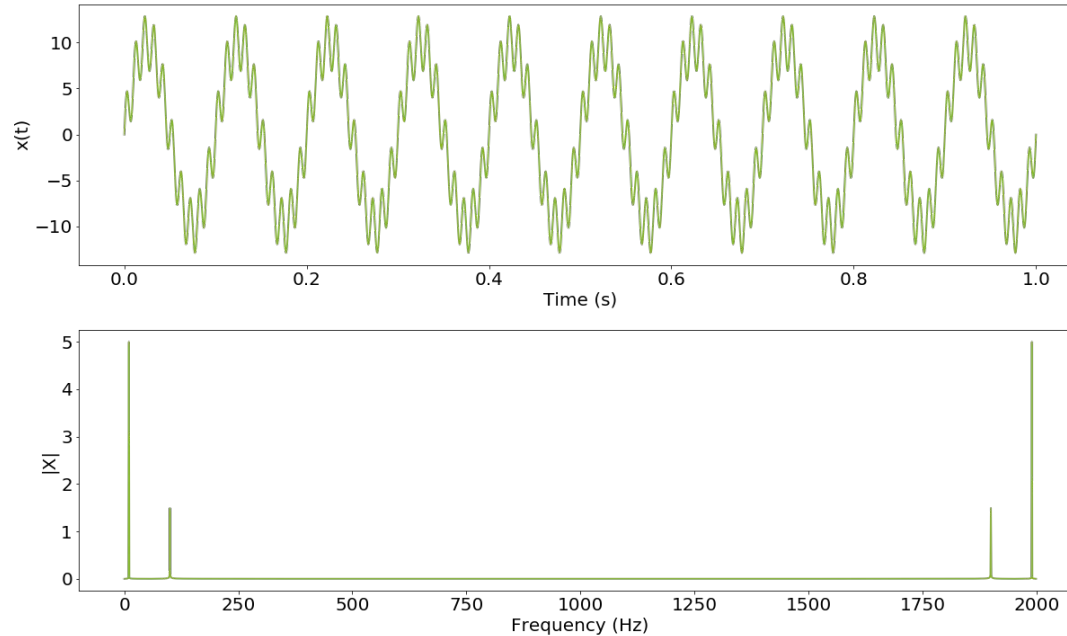- $F(y) = \int_{-\infty}^{\infty} f(x)e^{-2\pi ixy}dx$

  $time \rightarrow frequency$
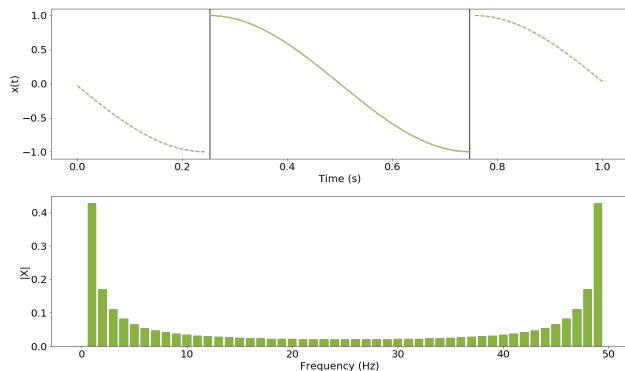
# Example of DFT

$$F = 2kHz$$

$$f(t) = 10\sin(2\pi 10 t) + 3\sin(2\pi 100 t)$$
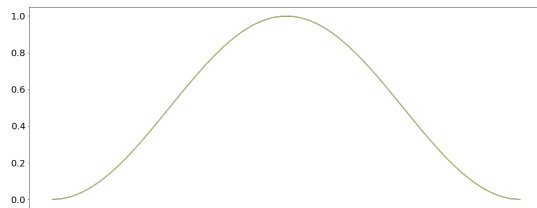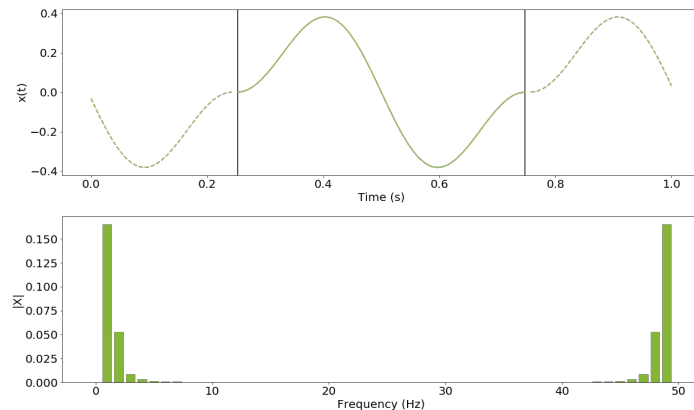
# Short-time Fourier transform

FFT + Windowing



Sliced signal

Window

Windowed signal

# Short-time Fourier transform

FFT + Windowing

# Spectrograms

# Mel Scale

$$m = 2595 \log_{10}\left(1 + \frac{f}{700}\right) = 1127 \ln\left(1 + \frac{f}{700}\right)$$

$$f = 700\left(10^{\frac{m}{2595}} - 1\right) = 700\left(e^{\frac{m}{1127}} - 1\right)$$

# CNN-approach



Mel-Spectrogram · Convolution Layer (32 x 3 x 3) · Max Pooling (2 x 2) · Convolution Layer (64 x 3 x 3) · Max Pooling (2 x 2) · Flatten · Dense (64) · Dense (1) · Output

# Automatic Speech Recognition

Task:

- Transform speech (audio) to text

Also known as:

- STT - Speech To Text

# WER

## Word Error Rate

**How to compute?**

Edit path from reference to prediction
- **S** – substitution count
- **D** – deletions count
- **I** – insertions count
- **C** – correct count
- **N** – **S** + **D** +**C** - Total word count in refrence

```
True: quick brown    fox jumped over  a lazy dog
Pred: quick brow  an fox jumped over    lazy dog
```

$$WER = \frac{S + D + I}{N} = \frac{S + D + I}{S + D + C}$$

# CER

Character Error Rate

**The same, but  on character level**

Questions:

1)    what is more important?
2)    what is more difficult to minimize?

# LibriSpeech

Domain:
- audiobooks

Parts:
- clean - low-WER speakers
- other - high-WER speakers

Features:
- 10-20s audio
- long sentences
- complex language

| subset | hours | per-spk minutes | female spkrs | male spkrs | total spkrs |
|---|---|---|---|---|---|
| dev-clean | 5.4 | 8 | 20 | 20 | 40 |
| test-clean | 5.4 | 8 | 20 | 20 | 40 |
| dev-other | 5.3 | 10 | 16 | 17 | 33 |
| test-other | 5.1 | 10 | 17 | 16 | 33 |
| train-clean-100 | 100.6 | 25 | 125 | 126 | 251 |
| train-clean-360 | 363.6 | 25 | 439 | 482 | 921 |
| train-other-500 | 496.7 | 30 | 564 | 602 | 1166 |

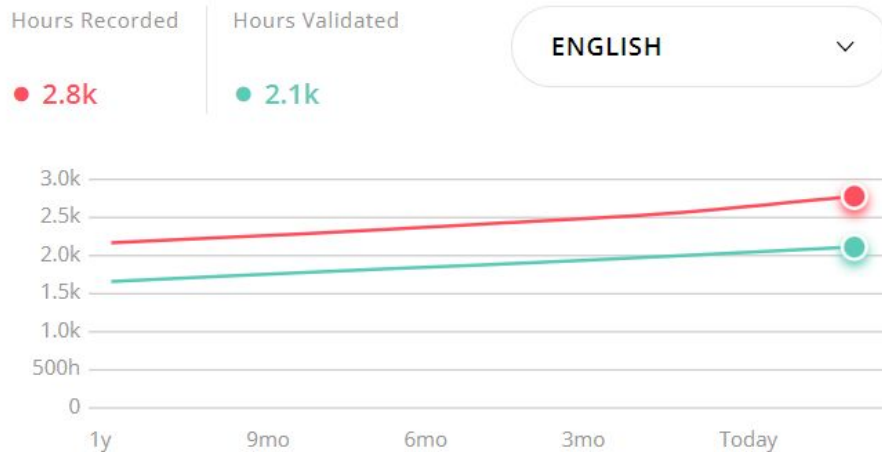[1]

Human WER:
- test-clean: **5.83**
- test-other: **12.69**

# Mozilla Common Voice

Domain:
- Short random phrases
- Multiple languages

Features:
- Crowdsourced
- Simple language
- Short phrases
- Frequently updated and validated

| Hours Recorded | Hours Validated | ENGLISH ⌄ |
|---|---|---|
| ● 2.8k | ● 2.1k | |

```
3.0k
2.5k                                              ●
2.0k                                              ●
1.5k
1.0k
500h
   0
      1y      9mo      6mo      3mo      Today
```
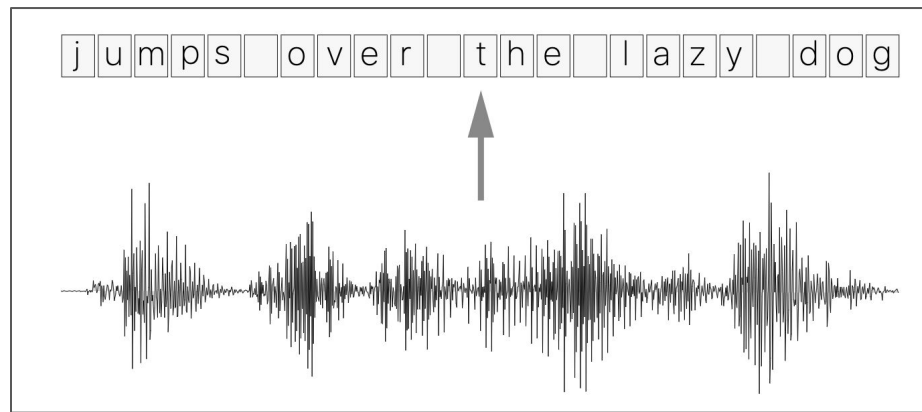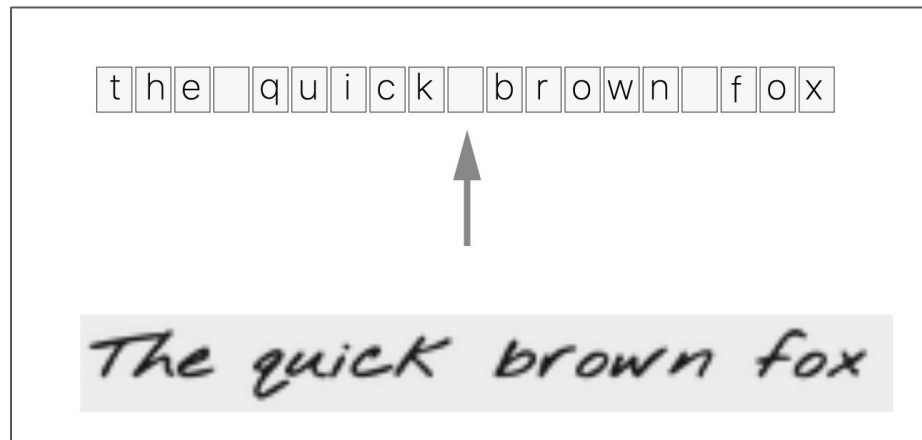
# Fundamental problem

Alignment problem
- Variable length input
- Variable length output
- No alignment

How to train?



[3]

# Connectionist Temporal Classification

2006
Idea:

- Split input into frames
- Classify each frame into *num_letters* classes
- Merge consecutive letters

Issues?

| $x_1$ $x_2$ $x_3$ $x_4$ $x_5$ $x_6$ | input ($X$) |
| c c a a a t | alignment |
| c a t | output ($Y$) |

[3]

# Connectionist Temporal Classification

Idea:

- Split input into frames
- Classify each frame into *num_letters* classes
- Merge consecutive letters

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | $x_6$ | input ($X$) |
|---|---|---|---|---|---|---|
| c | c | a | a | a | t | alignment |
| c | | a | | | t | output ($Y$) |

[3]

Issues?

- Multiple consecutive letters in target word (e.g.: he**ll**o)
- Silence between words and letters

23

# CTC: Empty token



h h e $\epsilon$ $\epsilon$ l l l $\epsilon$ l l o

h e $\epsilon$ l $\epsilon$ l o

h e l l o

h e l l o

[3]

First, merge repeat characters.

Then, remove any $\epsilon$ tokens.

The remaining characters are the output.

# CTC: Empty token

**Valid Alignments**

| ϵ | c | c | ϵ | a | t |

| c | c | a | a | t | t |

| c | a | ϵ | ϵ | ϵ | t |

[3]

**Invalid Alignments**

| c | ϵ | c | ϵ | a | t |

corresponds to
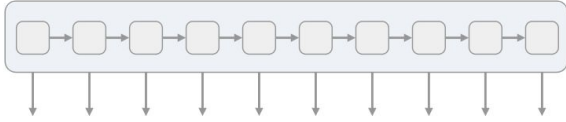$Y$ = [c, c, a, t]

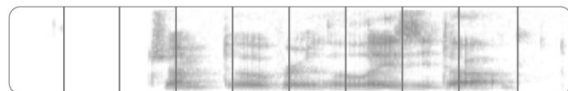| c | c | a | a | t | |

has length 5

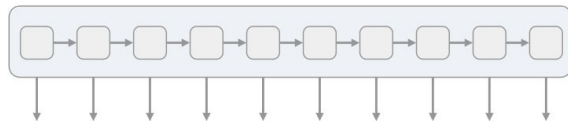| c | ϵ | ϵ | ϵ | t | t |

missing the 'a'

We start with an input sequence, like a spectrogram of audio.
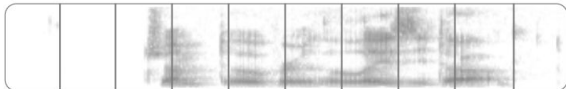


The input is fed into an RNN, for example.

[3]

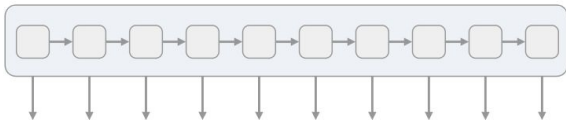We start with an input sequence, like a spectrogram of audio.

The input is fed into an RNN, for example.

The network gives $p_t\,(a \mid X)$, a distribution over the outputs $\{h, e, l, o, \epsilon\}$ for each input step.

[3]

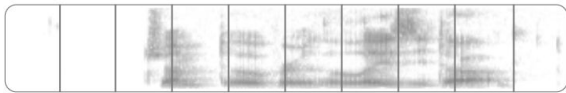We start with an input sequence, like a spectrogram of audio.


The input is fed into an RNN, for example.

| h | h | h | h | h | h | h | h | h | h |
| e | e | e | e | e | e | e | e | e | e |
| l | l | l | l | l | l | l | l | l | l |
| o | o | o | o | o | o | o | o | o | o |
| $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ | $\epsilon$ |

The network gives $p_t(a \mid X)$, a distribution over the outputs {h, e, l, o, $\epsilon$} for each input step.

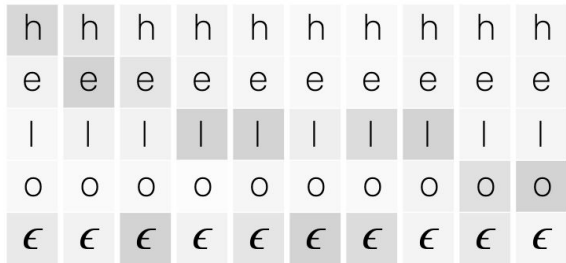| h | e | $\epsilon$ | l | l | $\epsilon$ | l | l | o | o |
| h | h | e | l | l | $\epsilon$ | $\epsilon$ | l | $\epsilon$ | o |
| $\epsilon$ | e | $\epsilon$ | l | l | $\epsilon$ | $\epsilon$ | l | o | o |

With the per time-step output distribution, we compute the probability of different sequences
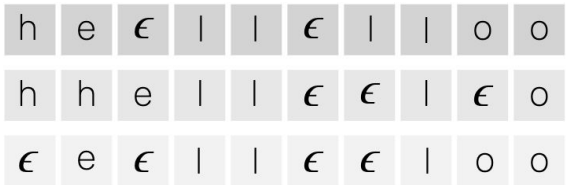
[3]

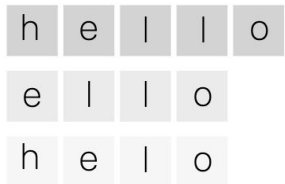We start with an input sequence, like a spectrogram of audio.

The input is fed into an RNN, for example.

The network gives $p_t\,(a \mid X)$, a distribution over the outputs $\{h, e, l, o, \epsilon\}$ for each input step.
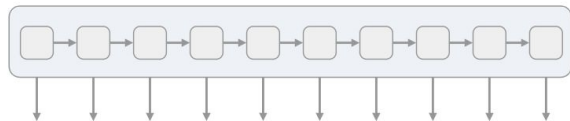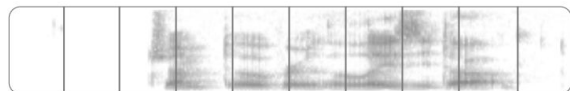
With the per time-step output distribution, we compute the probability of different sequences

By marginalizing over alignments, we get a distribution over outputs.
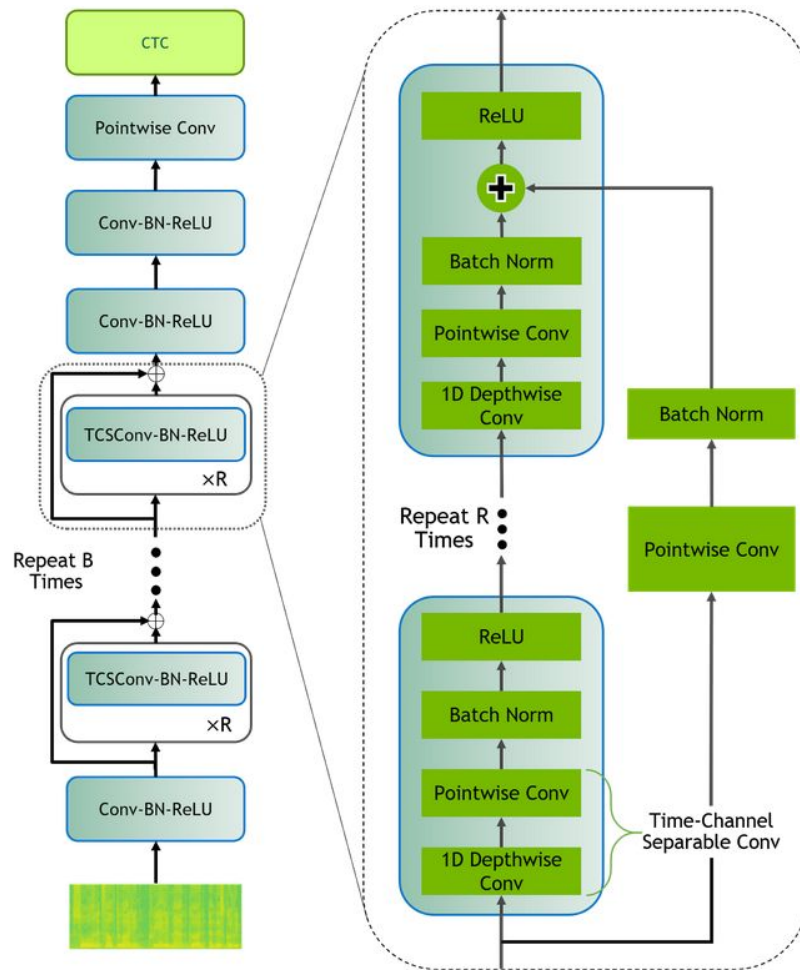
[3]

Loss Function:

$$p(Y \mid X) = \sum_{A \in \mathcal{A}_{X,Y}} \prod_{t=1}^{T} p_t(a_t \mid X)$$

The CTC conditional **probability**

**marginalizes** over the set of valid alignments

computing the **probability** for a single alignment step-by-step.

[3]

# QuartzNet

# Language Models (LM): motivation

Problem:

- Spelling of a word heavily depends on its context
- We always want **more** data
- Labeled audio data is difficult to obtain

hypo 1: let's go **two** a movie (score: 0.21)
hypo 2: let's go **to** a movie (score: 0.19)
hypo 3: let's go **too** a movie (score: 0.13)

Idea:

- ASR predicts texts
- Unlabeled text data is **very easy** to get
- Let's make use of it!

# LM: motivation

Language model - a model that
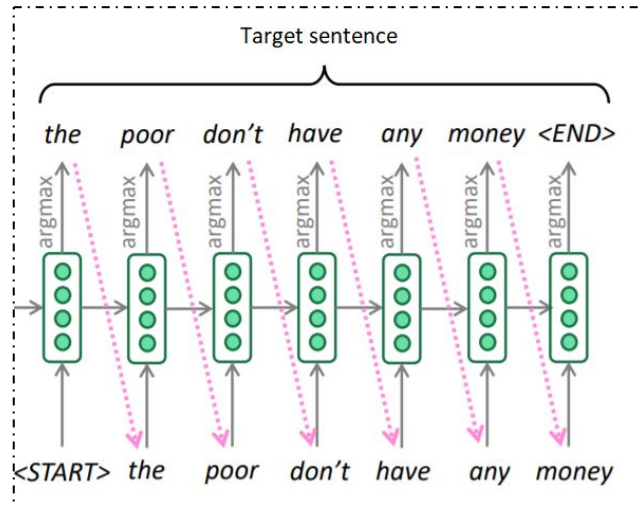estimates the probability of a text.

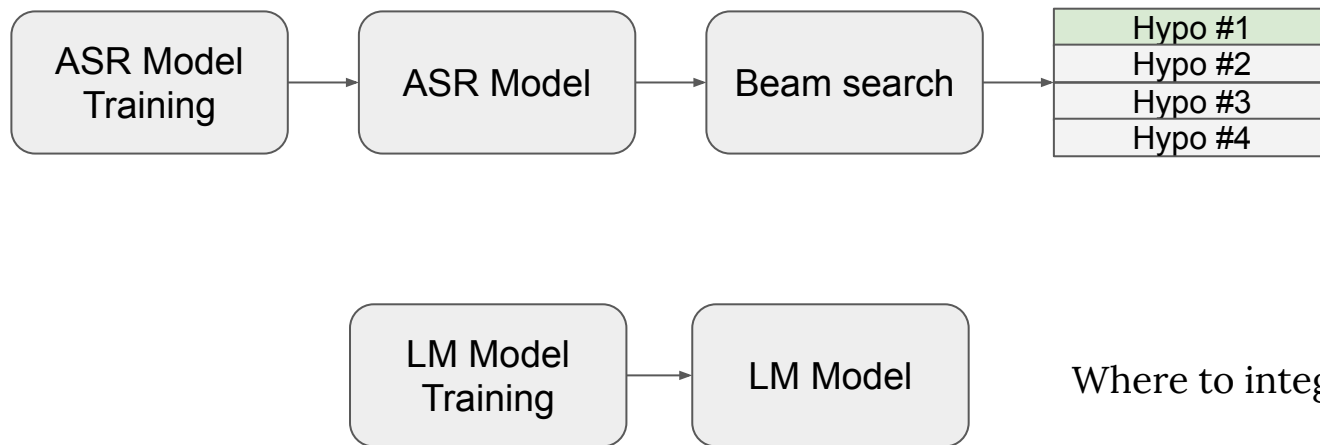Usually requires unlabeled text corpus to train.

Examples:

- N-gramms (simple)
- Neural networks (compex) e.g: Bert, GPT-3

$P$(let's go **two** a movie) = 0.01

$P$(let's go **to** a movie) = 0.6



[1]

# LM: ASR pipeline

```
┌──────────────┐     ┌──────────────┐     ┌──────────────┐     ┌──────────────┐
│  ASR Model   │ ──▶ │  ASR Model   │ ──▶ │ Beam search  │ ──▶ │   Hypo #1    │
│  Training    │     │              │     │              │     │   Hypo #2    │
└──────────────┘     └──────────────┘     └──────────────┘     │   Hypo #3    │
                                                               │   Hypo #4    │
                                                               └──────────────┘

┌──────────────┐     ┌──────────────┐
│  LM Model    │ ──▶ │  LM Model    │       Where to integrate LM?
│  Training    │     │              │
└──────────────┘     └──────────────┘
```

# LM: final hypothesis rescoring
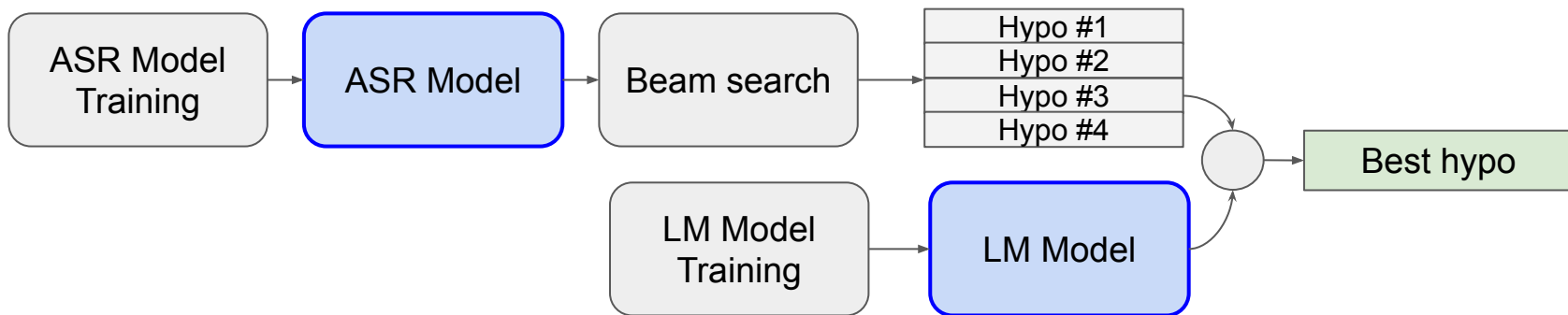
a.k.a.: *second-pass rescoring*

Prerequisites:
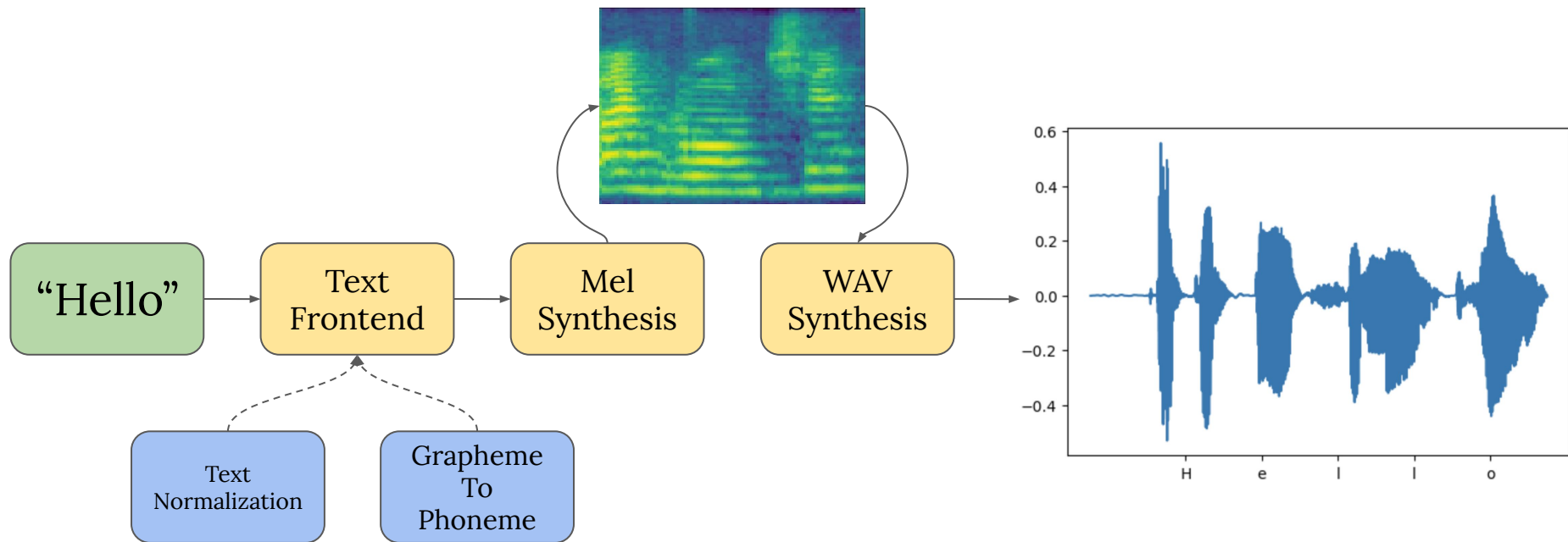- Trained ASR
- Trained LM

Concept: Rescore beam-search output with LM scores

$$P_{final} = P_{ASR}(Y|X) + \alpha \cdot P_{LM}(Y) + \beta \cdot len(Y)$$

***Oracle WER*** – lowest WER among all beam search output hypothesis

# Text-to-Speech (TTS)

# Tacotron 2