

# **Автоматическая обработка звуковых образов.**

Евгений Борисов

# Обработка звуковых образов

## Задачи

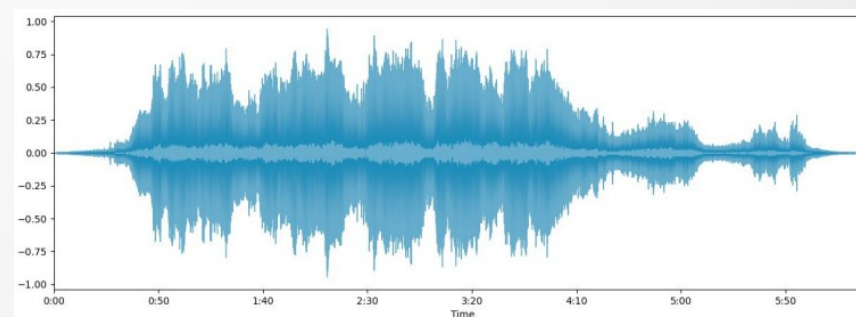
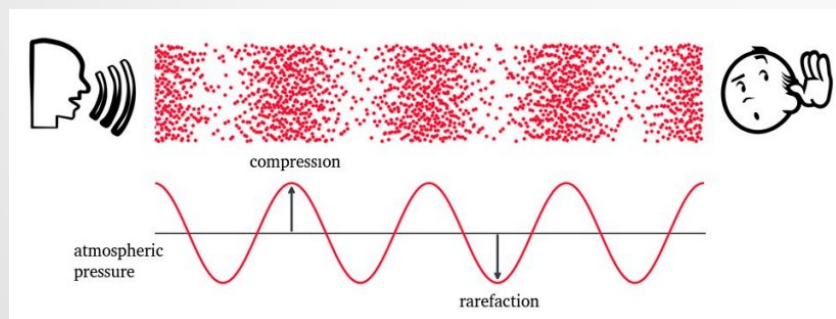
- Audio classification
- Speech recognition / speaker verification / speaker diarization
- Audio denoising / audio upsampling
- Music Information Retrieval ( Mood Classification )
- Audio styling
- Audio synthesis

# Обработка звуковых образов

## звуковые волны

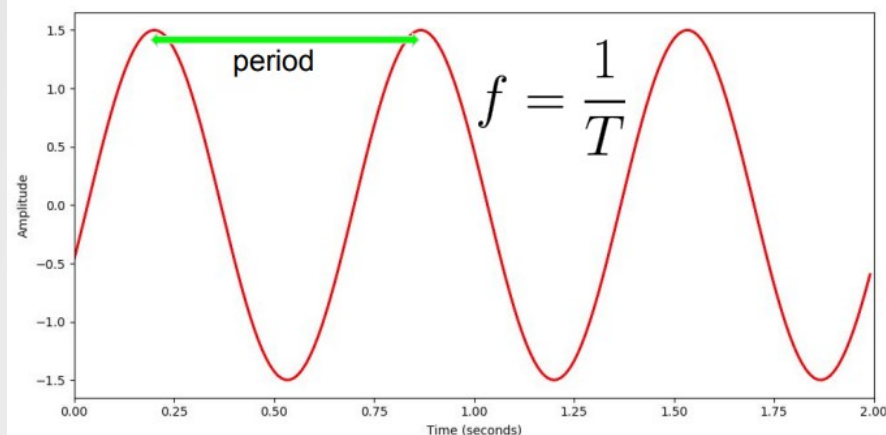
регулярные изменения давления воздуха

порождаются механической вибрацией объектов

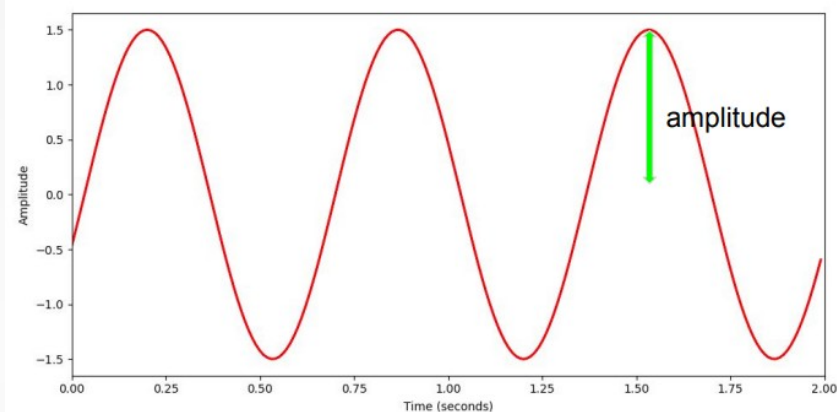


# Обработка звуковых образов

## частота



## амплитуда



$$y(t) = A \cdot \sin(2\pi f t + \phi)$$

$\phi$  – амплитуда ;

$f$  – частота ;

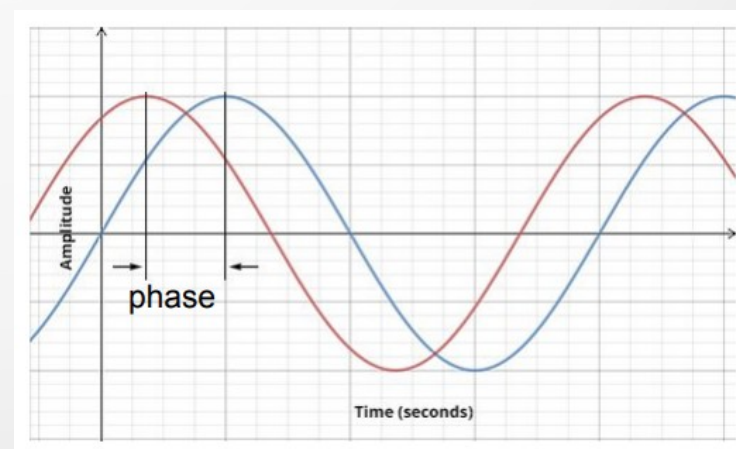
$t$  – время ;

$\phi$  – фаза ;

Частота определяет высоту звука,

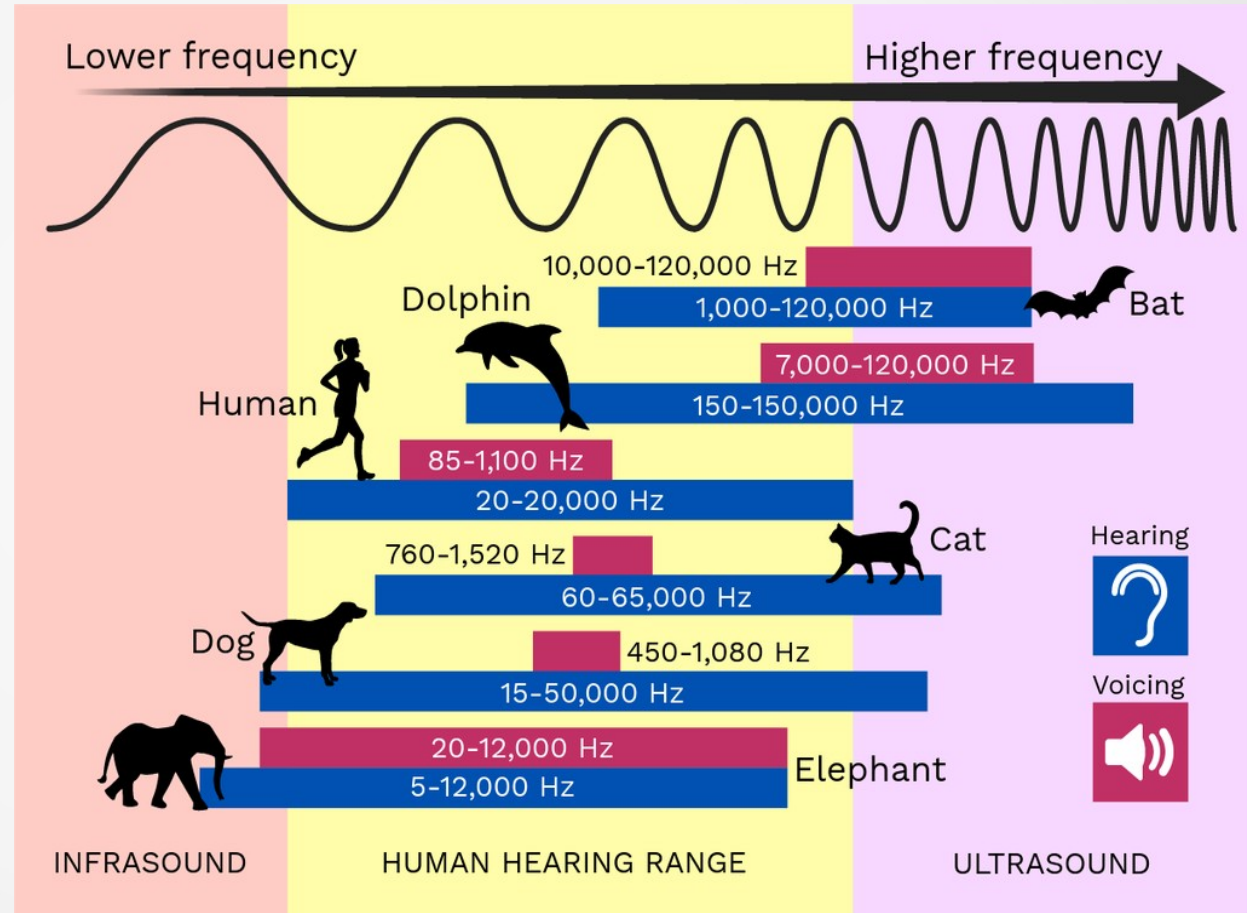
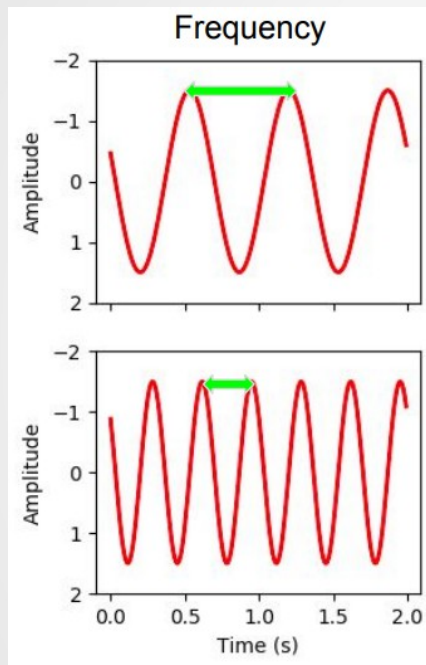
Амплитуда - громкость

## фаза



# Обработка звуковых образов

## ВЫСОКИЕ И НИЗКИЕ ЗВУКИ



# Обработка звуковых образов

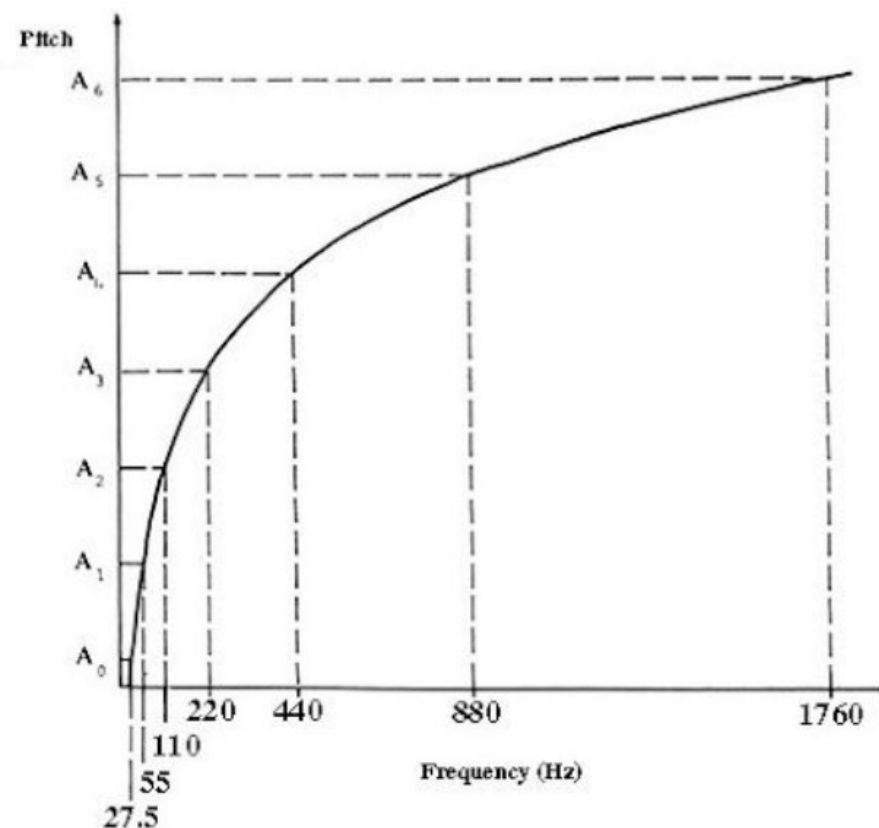
## ТОН И ЧАСТОТА

Тон - восприятие частоты человеком

логарифмическое восприятие,  
человек различает низкие звуки лучше чем высокие

$$F(p) = 440 \cdot 2^{\frac{p-69}{12}}$$

**Тембр** — многомерная характеристика составляющих звука



Note name	A0#	C1#	D1#	F1#	G1#	A1#	C2#	D2#	F2#	G2#	A2#	C3#	D3#	F3#	G3#	A3#	C4#	D4#	F4#	G4#	A4#	C5#	D5#	F5#	G5#	A5#	C6#	D6#	F6#	G6#	A6#	C7#	D7#	F7#	G7#	A7#																																																			
Midi number	21	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100	101	102	103	104	105	106	107	108
Note name	A0	B0	C1	D1	E1	F1	G1	A1	B1	C2	D2	E2	F2	G2	A2	B2	C3	D3	E3	F3	G3	A3	B3	C4	D4	E4	F4	G4	A4	B4	C5	D5	E5	F5	G5	A5	B5	C6	D6	E6	F6	G6	A6	B6	C7	D7	E7	F7	G7	A7	B7	C8																																			

440 Hz 880 Hz



# Обработка звуковых образов

**Громкость** — субъективное восприятие человеком, зависит от возраста, частоты и др.

**мощность звука** - энергия передаваемая от источника в единицу времени ( Вт, W )

**интенсивность звука** - мощность переданная на единицу площади ( Вт/м<sup>2</sup>, W/m<sup>2</sup> )

**уровень интенсивности** — логарифмическое восприятие интенсивности человеком (dB)

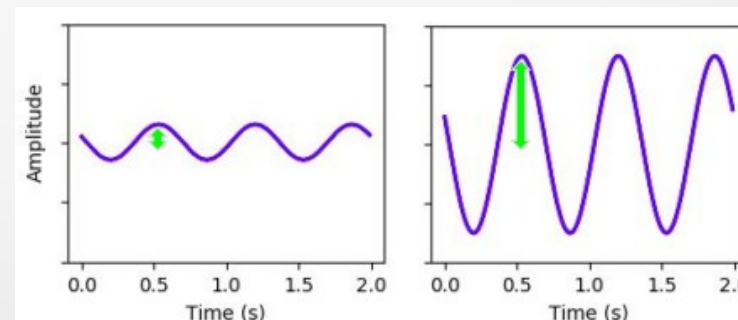
порог слышимости, threshold of hearing (TOH)  
 $TOH = 10^{-12} \text{ W/m}^2$

болевой порог громкости, threshold of pain (TOP)  
 $TOP = 10 \text{ W/m}^2$



$$dB(I) = 10 \cdot \log_{10} \frac{I}{TOH}$$

$$dB(TOH) = 10 \cdot \log_{10} \frac{TOH}{TOH} = 10 \cdot \log_{10} 1 = 0$$



# Обработка звуковых образов

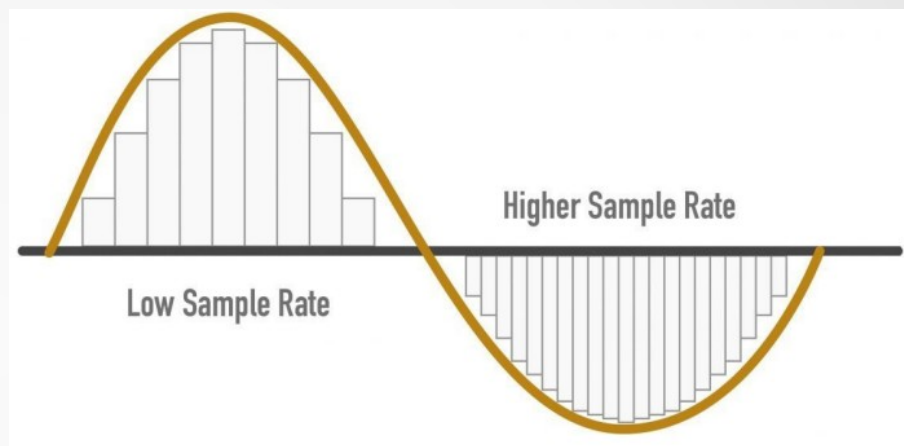
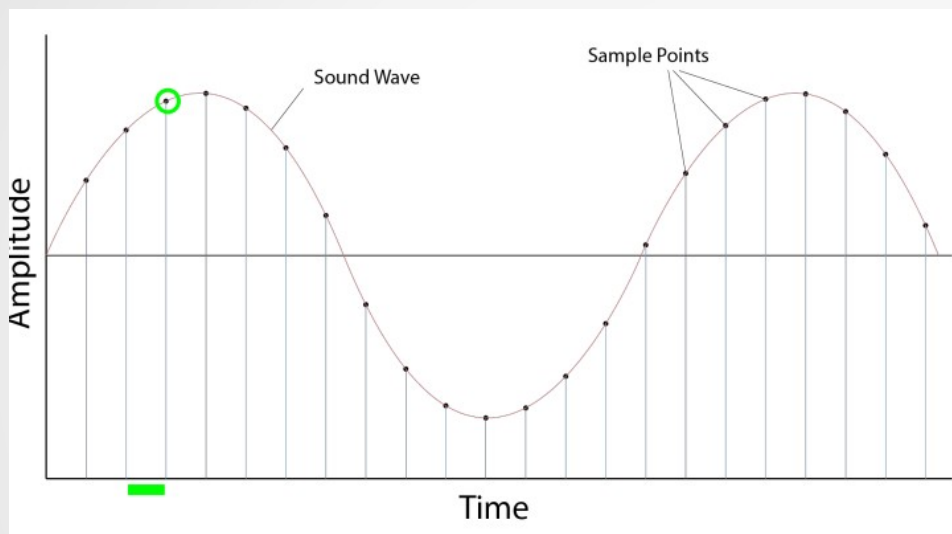
представление звукового сигнала





# Обработка звуковых образов

представление звукового сигнала — частота дискретизации



Sampling rate

$$Sr = \frac{1}{T}$$

Теорема Найквиста – Шенона – Котельникова

$$f_N = \frac{Sr}{2}$$

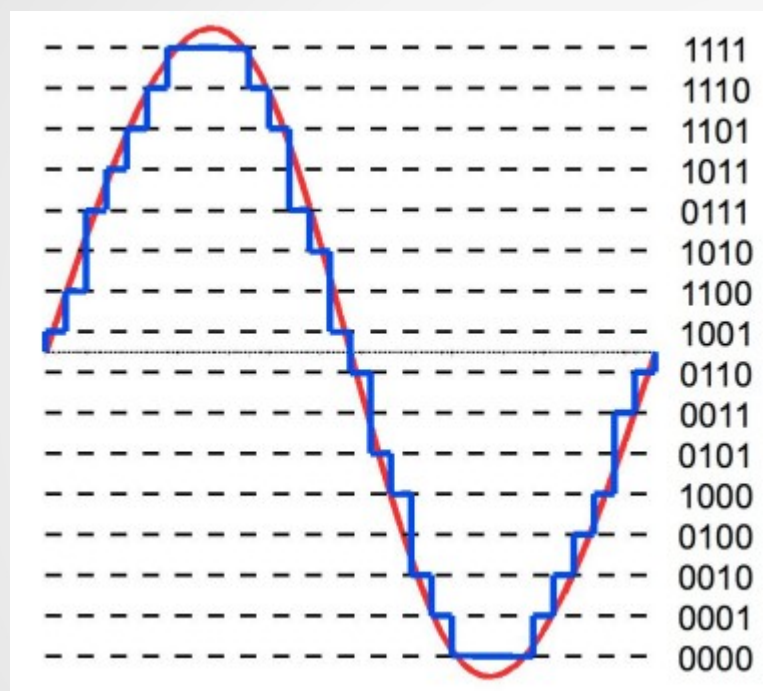
Если максимальная частота воспринимаемая человеком 22050Гц,

То для воспроизведения без (существенных) потерь минимальная частота дискретизации звукового сигнала (sample rate) должна быть 44100Гц ( CD качество музыки )



# Обработка звуковых образов

представление звукового сигнала - quantization, bit depth



**bit depth** - насколько точно будем кодировать уровень сигнала

**CD качество**

Sample rate = 44100 Hz

bit depth = 16 bits

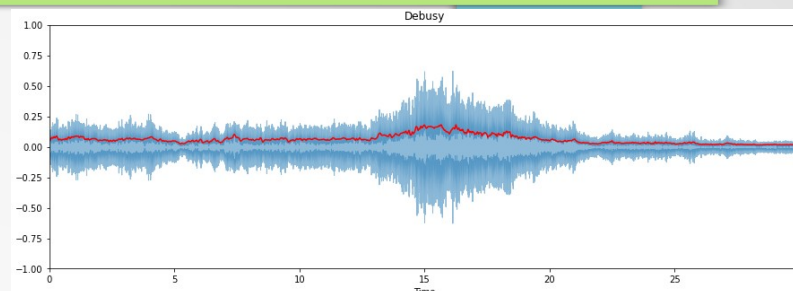
$$2^{16} = 65536$$



# Обработка звуковых образов

## извлечение признаков из звуковых сигналов

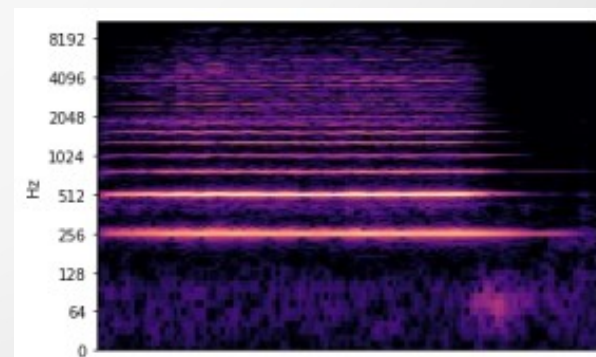
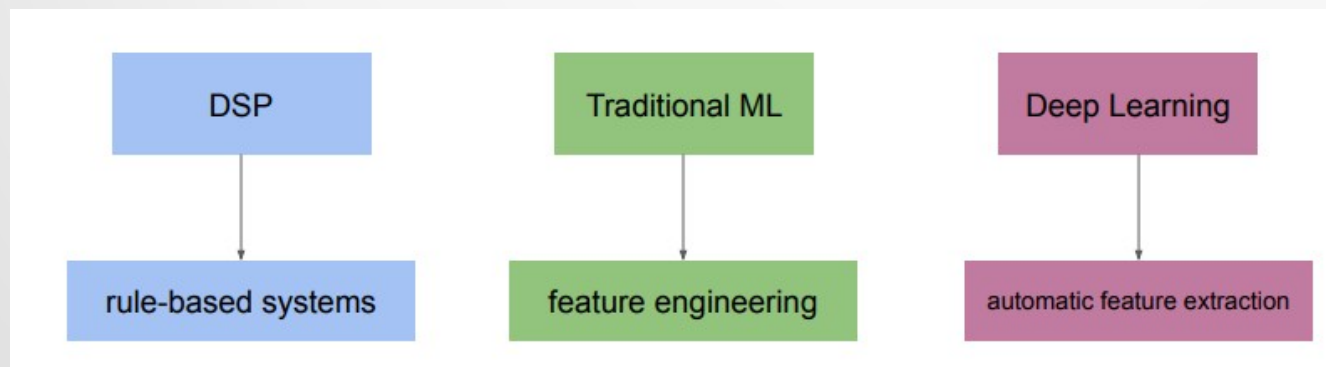
- моментальная характеристика (точечная)
- характеристика сегмента (на отрезке)
- общая характеристика (всё целиком)



## различные типы звуковых сигналов - разные наборы признаков

- Обработка музыкальных образов (Music Information Retrieval, MIR)
- Задача распознавания речи
- Прочее

- Amplitude envelope
- Root-mean square energy
- Zero crossing rate
- Spectrogram
- Spectral centroid
- Mel-spectrogram



# Обработка звуковых образов

## **Типы аудиопризнаков:**

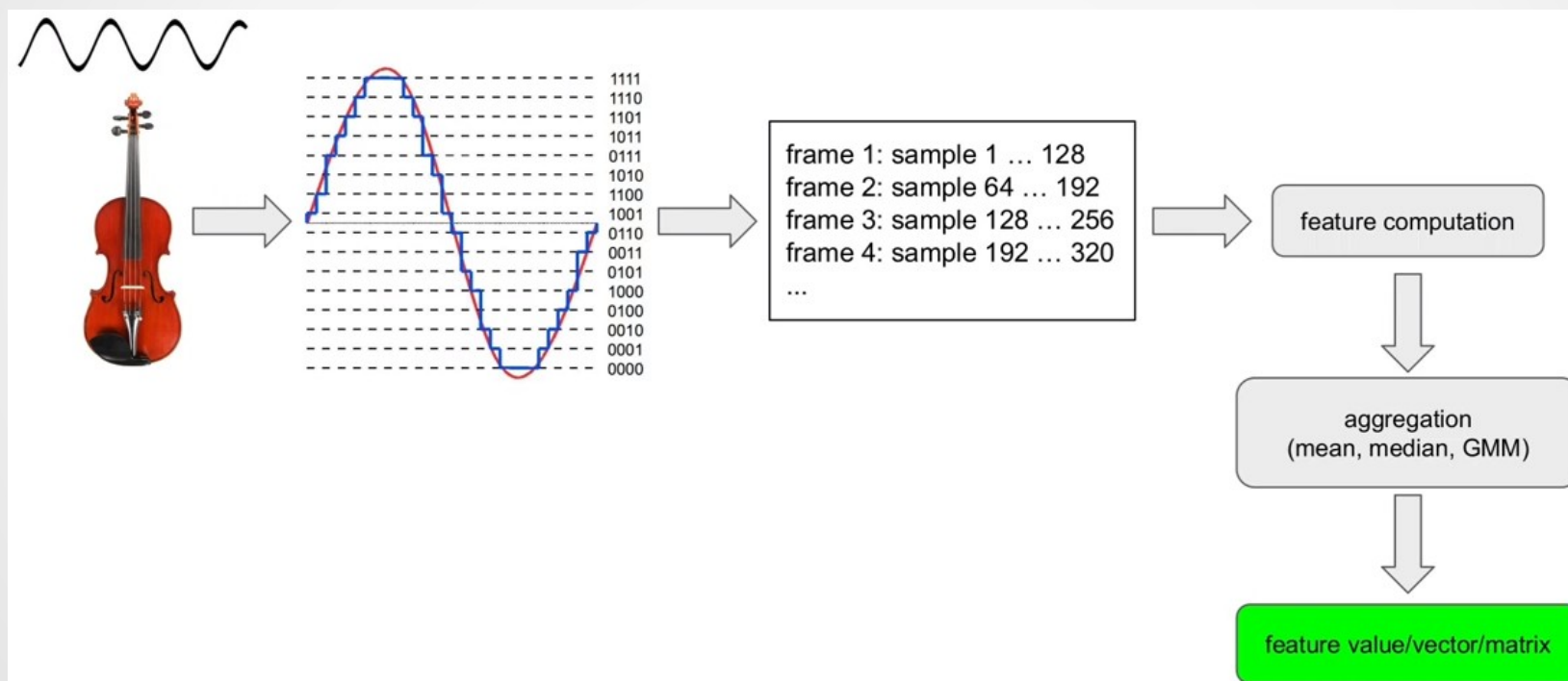
Временные характеристики (Time-domain features)

Частотные характеристики (frequency domain features)

# Обработка звуковых образов

## Временные характеристики (Time-domain features)

считаем характеристику для каждого фрейма и агрегируем результаты



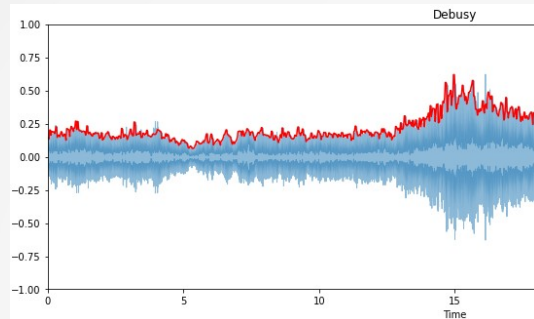
# Обработка звуковых образов

## Временные характеристики (Time-domain features)

считаем характеристику для каждого фрейма и агрегируем результаты

*Amplitude envelope (AE)*

$$AE(t) = \max_{k=t \cdot K}^{(t+1) \cdot K - 1} s(k)$$

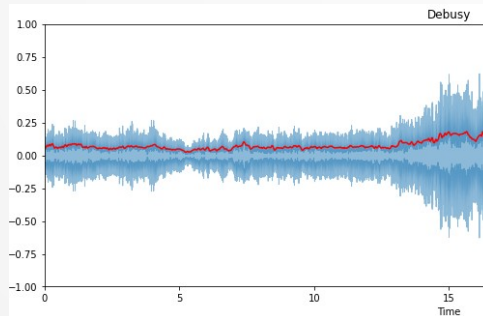


характеристика чувствительна к выбросам

можно использовать для определения жанра музыки

*Root-mean-square energy (RMS)*

$$RMS(t) = \sqrt{\frac{1}{K} \sum_{k=t \cdot K}^{(t+1) \cdot K - 1} s^2(k)}$$



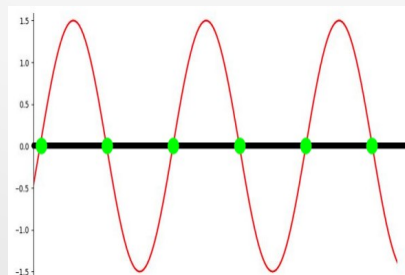
определяет уровень громкости

менее чувствительна к выбросам чем AE

можно использовать для сегментации аудио

*Zero-crossing rate (ZCR)*

$$ZCR(t) = \frac{1}{2} \sum_{k=t \cdot K}^{(t+1) \cdot K - 1} |\text{sign}(s(k)) - \text{sign}(s(k+1))|$$



детектор ритмичной музыки

детектор голоса



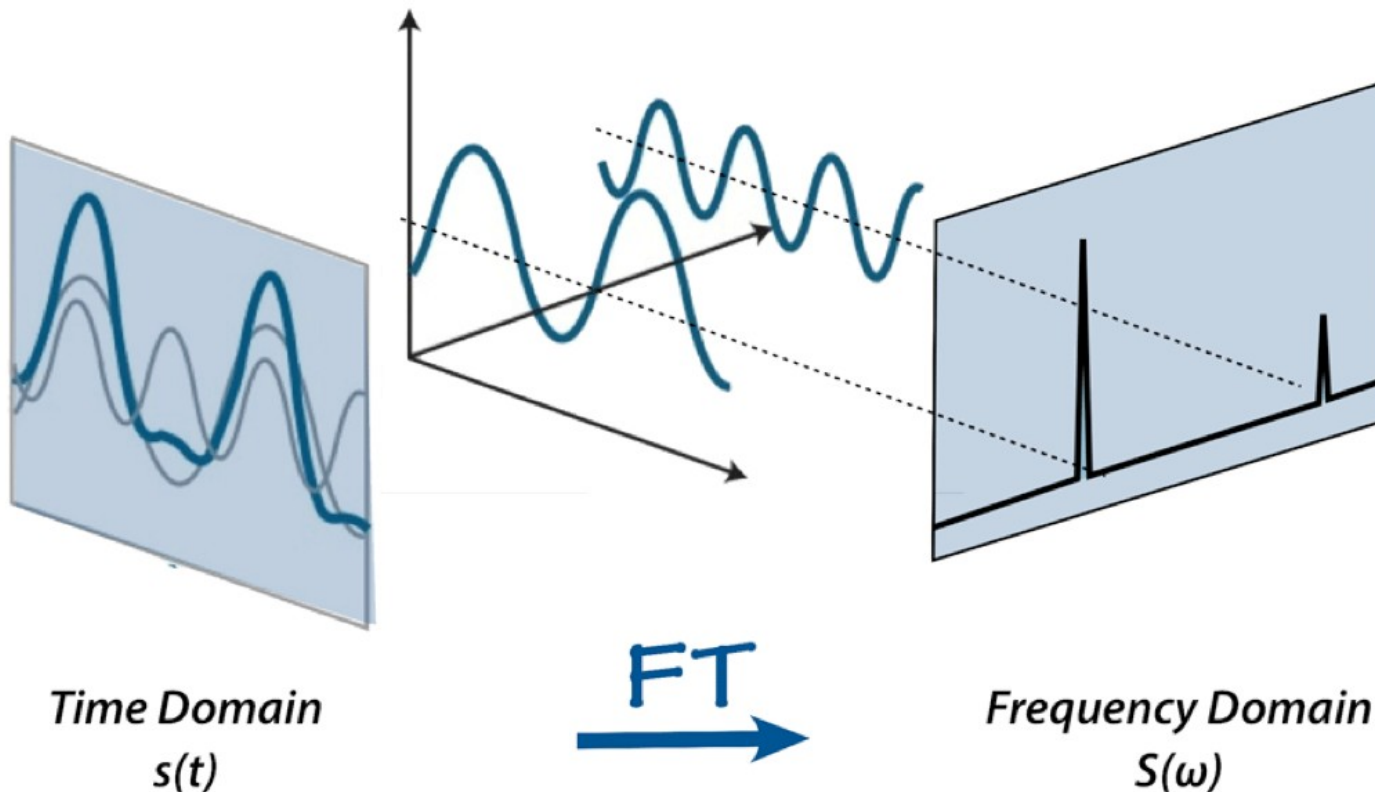
# Обработка звуковых образов



## Частотные характеристики (frequency domain features)

преобразование Фурье - разбираем сигнал на частотные составляющие

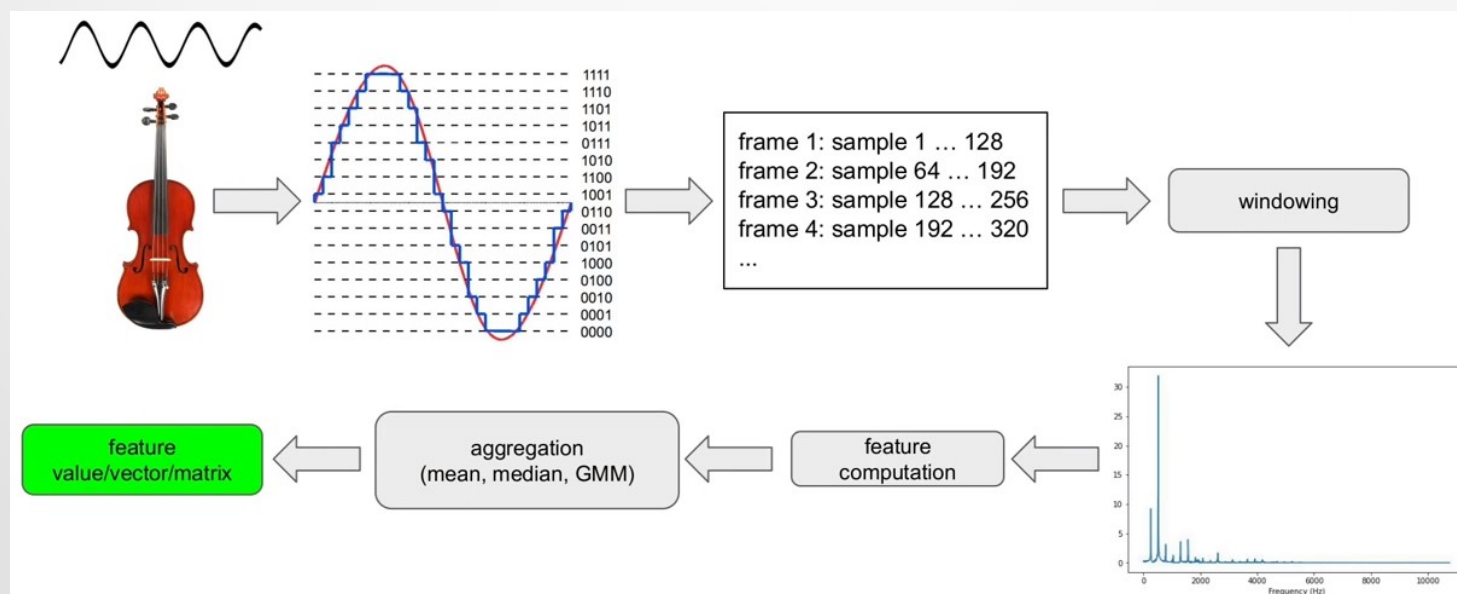
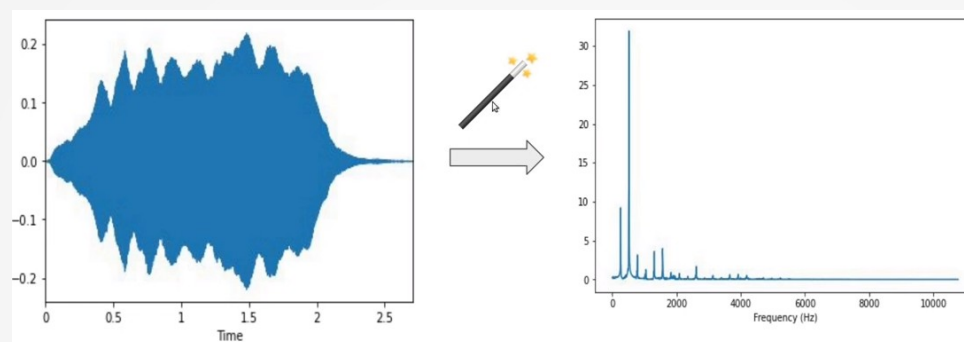
$$G(f) = \int_{-\infty}^{\infty} g(t) e^{-2\pi i f t} dt.$$



# Обработка звуковых образов

## Частотные характеристики (frequency domain features)

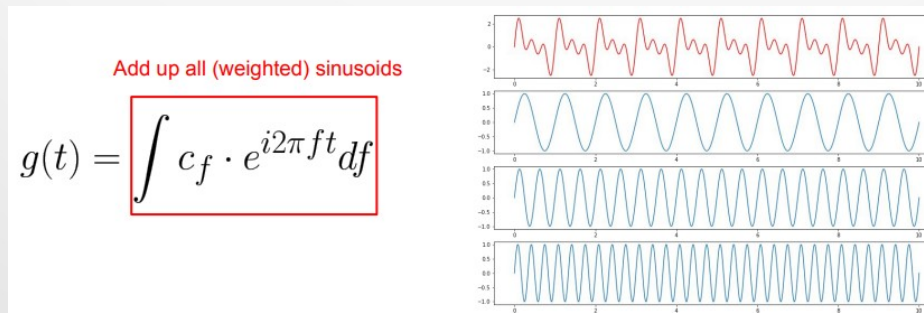
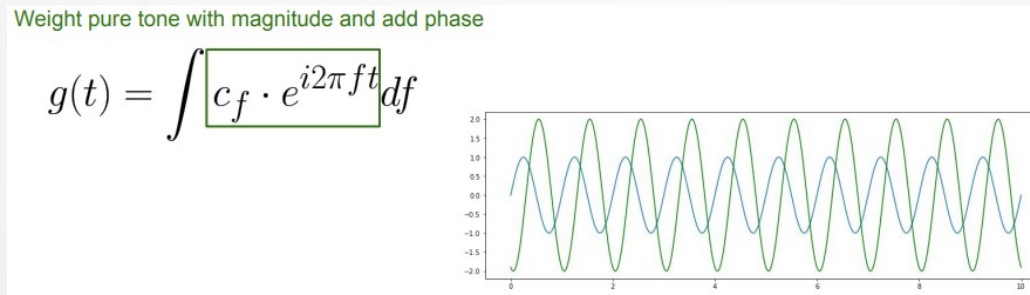
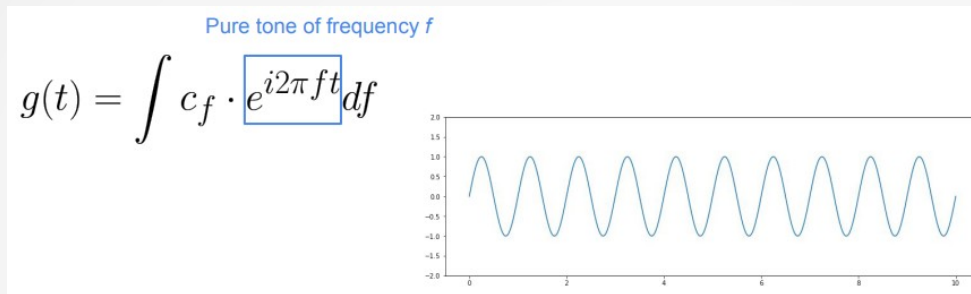
разбираем сигнал на частотные составляющие, считаем характеристику и агрегируем результаты



# Обработка звуковых образов

## Частотные характеристики (frequency domain features)

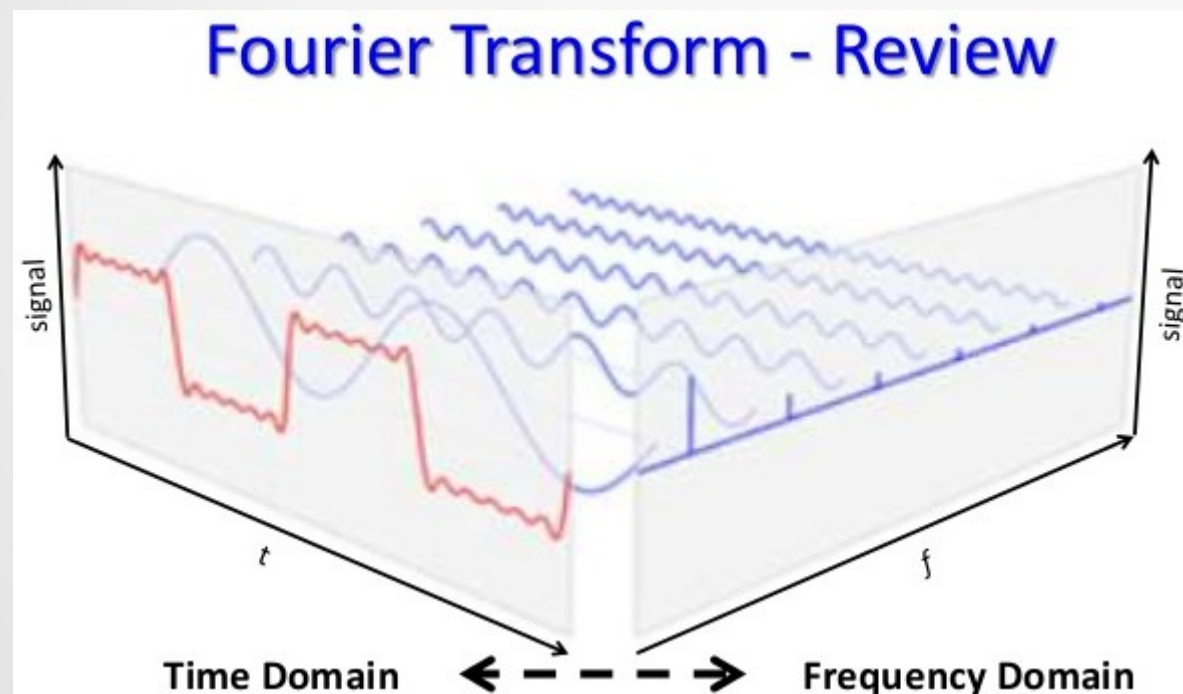
преобразование Фурье - разбираем сигнал на частотные составляющие



# Обработка звуковых образов

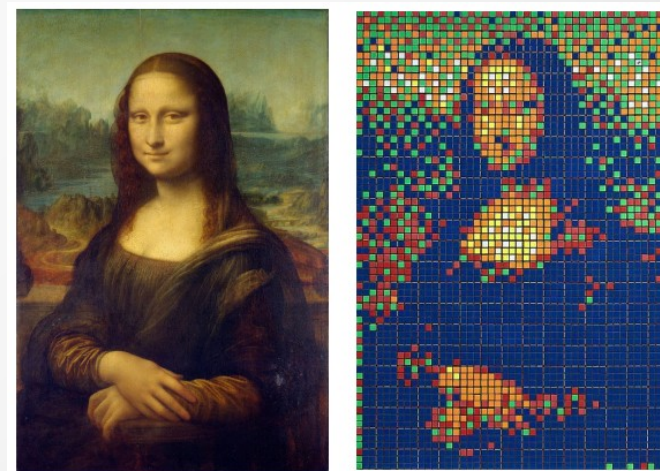
## Частотные характеристики (frequency domain features)

Прямое и обратное преобразование Фурье



$$\hat{g}(f) = \int g(t) \cdot e^{-i2\pi ft} dt$$

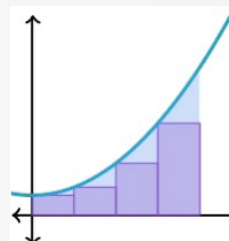
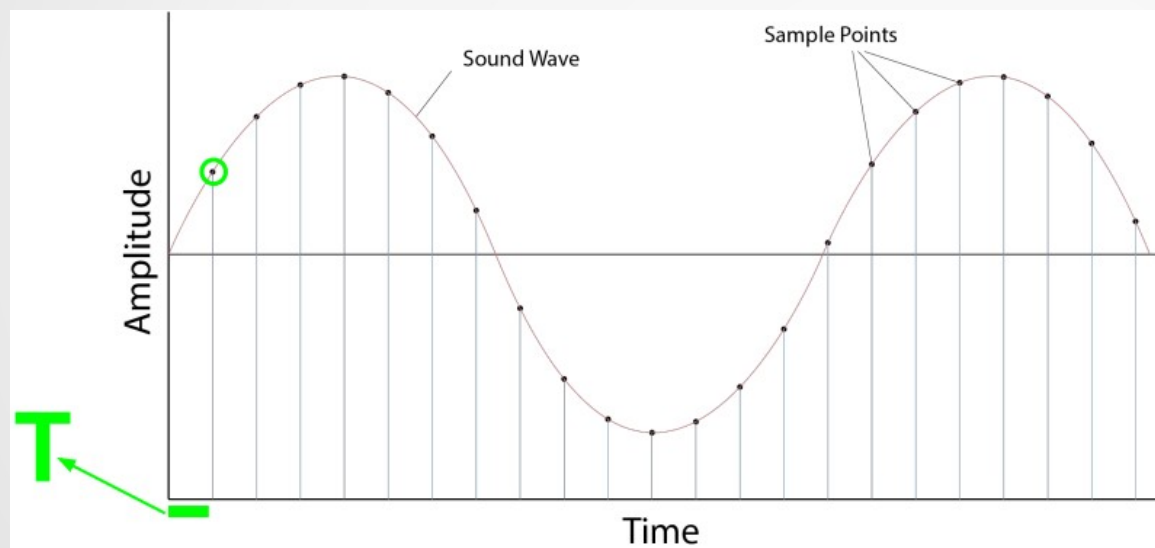
$$g(t) = \int c_f \cdot e^{i2\pi ft} df$$



# Обработка звуковых образов

## Частотные характеристики (frequency domain features)

дискретное преобразование Фурье DFT



$$\hat{g}(f) = \int g(t) \cdot e^{-i2\pi ft} dt$$

$$\hat{x}(f) = \sum_n x(n) \cdot e^{-i2\pi fn}$$



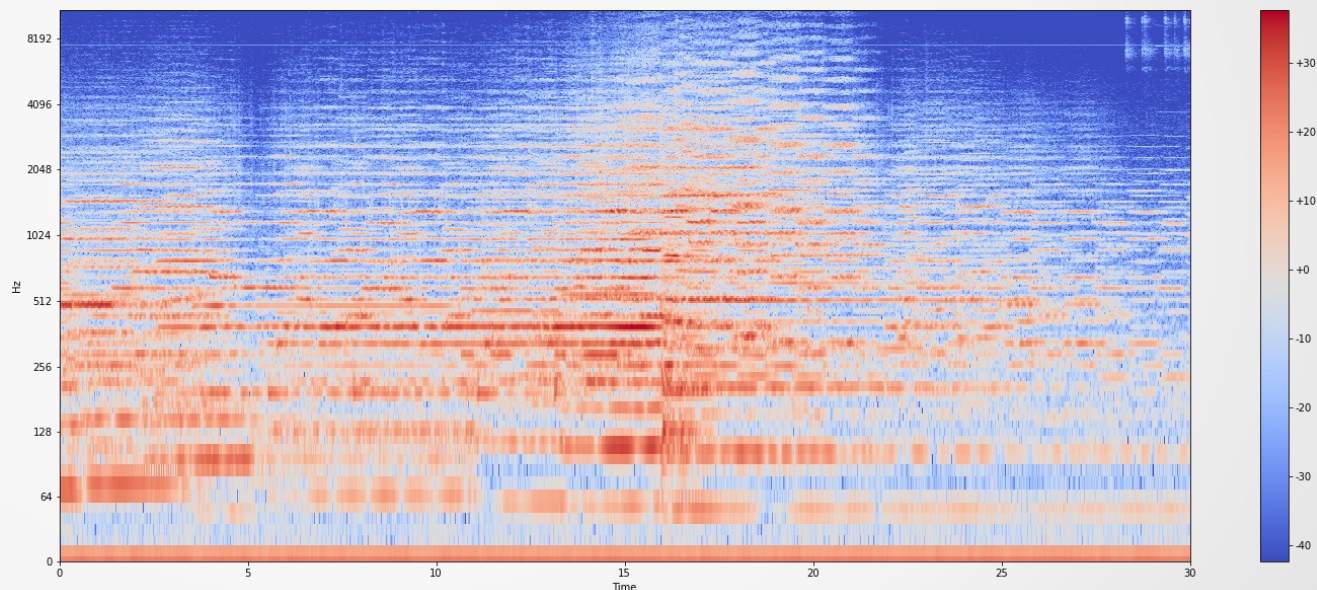
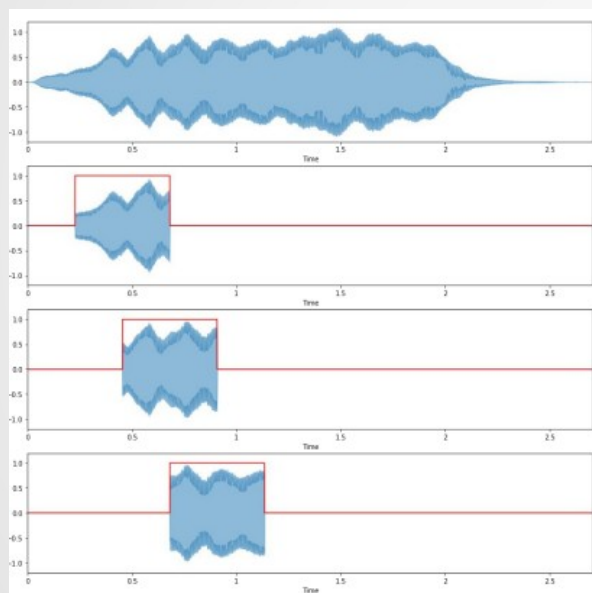
# Обработка звуковых образов



## Частотные характеристики (frequency domain features)

разбираем сигнал на частотные составляющие, считаем характеристику и агрегируем результаты

**Спектрограммы** — разбиваем образ на несколько перекрывающихся окон, в каждом применяем FT





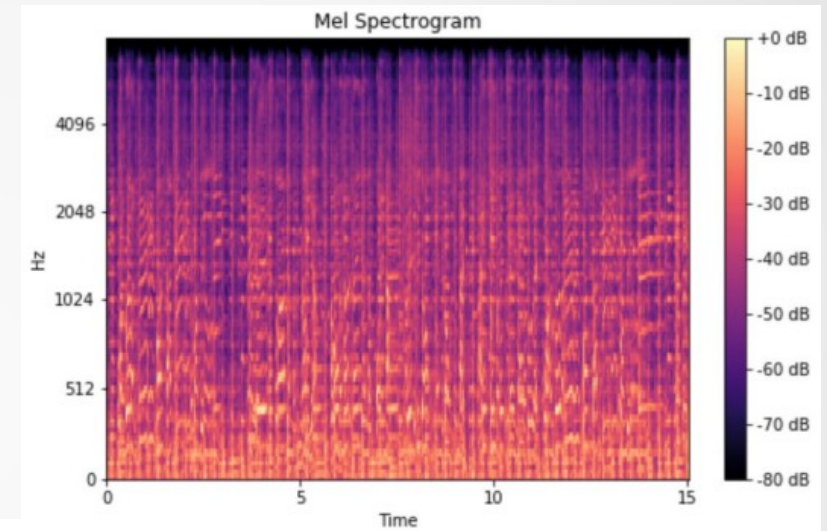
# Обработка звуковых образов

## Частотные характеристики (frequency domain features)

### Мэл-спектрограммы

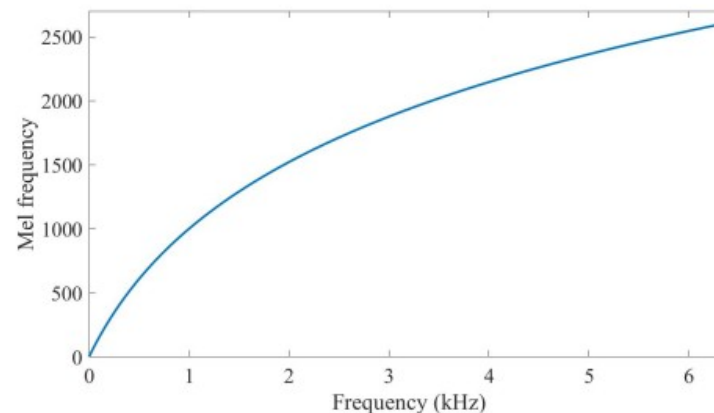
логарифмическое восприятие,  
человек различает низкие частоты лучше чем высокие

выполняем логарифмическое преобразование Hz  $\rightarrow$  Mel



$$m = 2595 \cdot \log\left(1 + \frac{f}{500}\right)$$

$$f = 700(10^{m/2595} - 1)$$

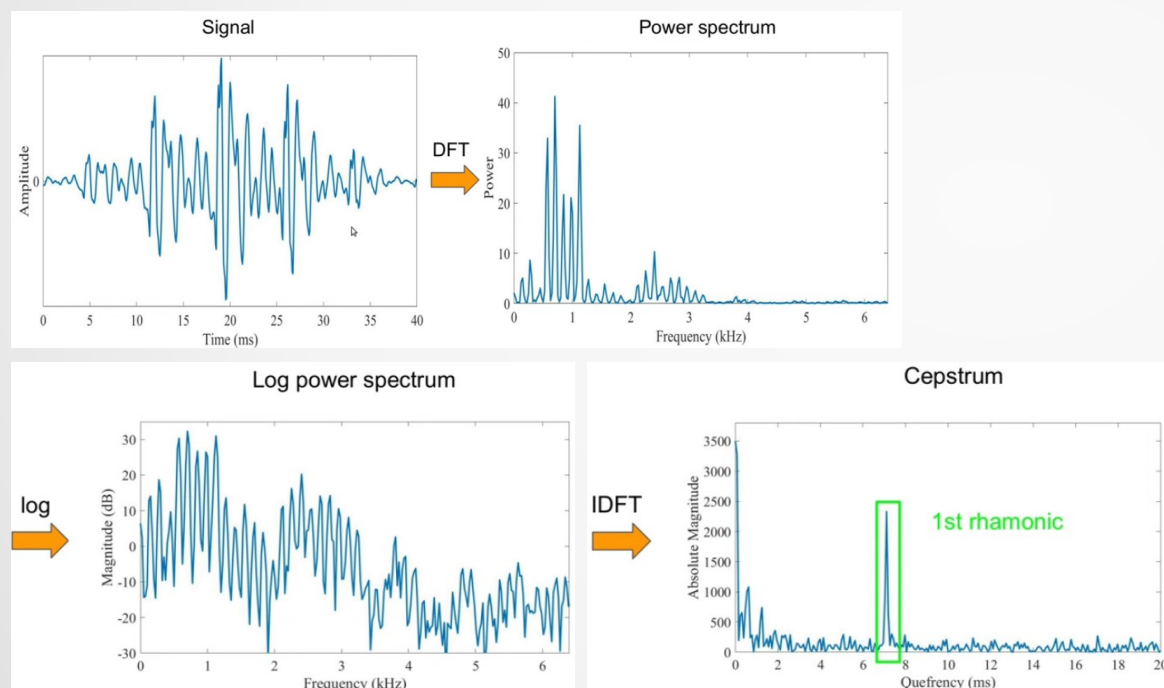


# Обработка звуковых образов

## Частотные характеристики (frequency domain features)

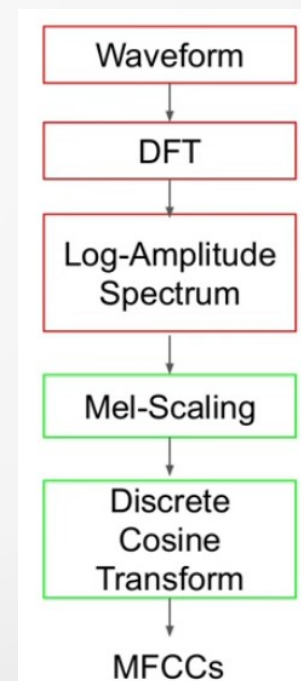
### Мэл-кепстральные коэффициенты ( Mel-Frequency Cepstral Coefficients, MFCC )

#### Spectrum - Cepstrum



$$C(x(t)) = F^{-1}[\log(F[x(t)])]$$

Time-domain signal  $x(t)$  is transformed to the Spectrum  $F[x(t)]$ , then to the Log spectrum  $\log(F[x(t)])$ , and finally to the Cepstrum  $C(x(t))$  via the inverse Fourier transform  $F^{-1}$ .



Mel-Frequency Cepstral Coefficients Explained Easily  
[https://www.youtube.com/watch?v=4\\_SH2nfbQZ8](https://www.youtube.com/watch?v=4_SH2nfbQZ8)

# Обработка звуковых образов

## Обработка музыкальных образов (Music Information Retrieval, MIR)

### задачи

- разделение звучания музыкальных инструментов
- классификация музыкальных образов
- оценка настроения музыкального образа

### характеристики

- ритм
- тон
- тембр
- гармоничность

# Обработка звуковых образов

## Литература

git clone [https://github.com/mechanoid5/ml\\_lectorium](https://github.com/mechanoid5/ml_lectorium)

Дауни Аллен Цифровая обработка сигналов на языке Python. 2017

Allen B.Downey Think DSP. Digital Signal Processing in Python. 2014

Meinard Müller, Brian McFee Instructional material for the Music Information Retrieval Workshop at CCRMA, Stanford University, 2014-18.

<https://github.com/bmcfee/stanford-mir>

Valerio Velardo The Sound of AI : Audio Signal Processing for Machine Learning

<https://www.youtube.com/playlist?list=PL-wATfeyAMNqIee7cH3q1bh4QJFAaeNv0>

<https://github.com/musikalkemist/AudioSignalProcessingForML>