# Audio features

level of abstract$^n$

high-level → normal people

mid-lev → expert

low-level → makes sense for machine

Temporal scope →
- Instantaneous ($\sim 50$ ms)
- Segment level (seconds)
- Global ($\sim 10$ sec)

Music aspect →
- Beat
- Timbre
- Pitch
- harmony

Time-frequency →
- Spectrogram
- Mel —"—
- Const Q-transform

(Signal Domain)

time domain $\xrightarrow{FT}$ Frequency domain

Fourier transform.

time domain:
- Amp envelope
- Rms energy
- Zero crossing rate

Frequency domain:
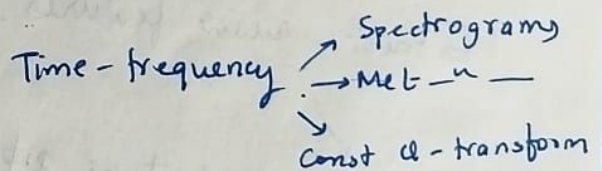- Band energy ratio
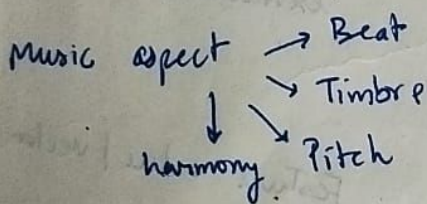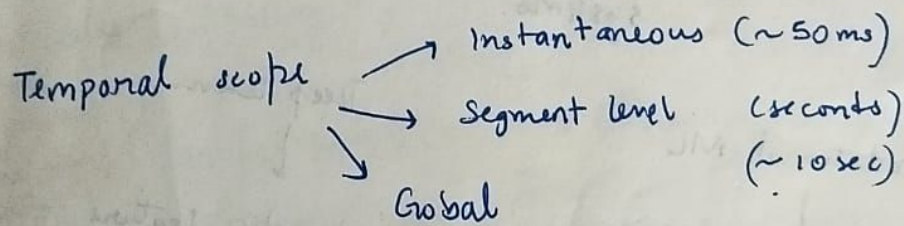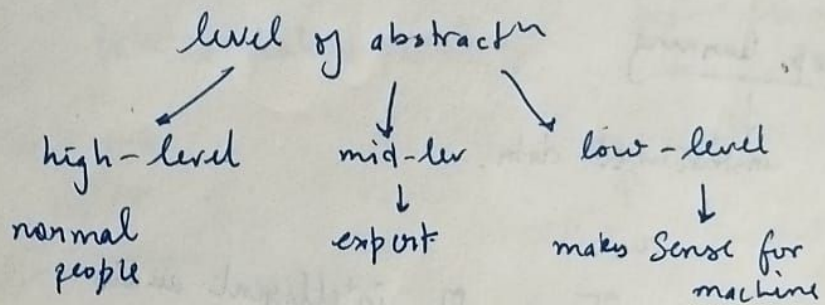- Spectral centroid
- Spectral flux

extracted from wave form

Computing MFCC

Waveform
↓
DFT
↓
log-Amp spectrum
↓
Mel-scaling
↓
Discrete cosine transform
↓
MFCCs

Why Discrete cosine transform.

- Simplified version of FT
- Get real valued coeff.

⇓

Gives the _coeff._

- De correlate energy in differen mel band
- Reduces # dimension to represent spectrum

↓

Tells us how well that frequency fits in the signal

coeff ↑ ⟹ good fitting.

How many coeff?

- Traditionally 1st 12-13 coeff.
- Use Δ & ΔΔ MFCCs          (derivatives)   change in coff
- Total 39 coeff per frame.                  in subsequent frame

# Mel Spectrograms.

Humans perceive frequency logarithmically.

Mel-scale

$$1000 \, mel = 1000 Hz$$

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

## Mel-Spectrograms

1) STFT

2) Amp $\longrightarrow$ DBs

3) Convert frequency to Mel scale

$\downarrow$

1. Choose number of mel bands
2. Construct mel filter banks
3. Apply mel filter banks to spectrograms

Cepstrum     Quefrency     Liftering     Rhamonic
$\downarrow$
Spectrum     Frequency     filtering     harmonic

$$C(x(t)) = F^{-1}\left[\log\left(\underbrace{F(x(t))}_{\text{DFT}}\right)\right]$$

log spectrum.

spectrum

$C(x(t))$
$\downarrow$
time-domain signal

$\underbrace{\qquad\qquad\qquad\qquad}_{\text{Cepstrum}}$

Speech = Convolution of vocal tract frequency response with glottal pulse.

Formalising speech

$$x(t) = e(t) \cdot h(t).$$

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

$$\log(X(t)) = \log(E(t)) + \log(H(t))$$

$\log(X(t))$
$\downarrow$
Speech

$\log(E(t))$
$\downarrow$
Glott
Glottal

$\log(H(t))$
$\searrow$
Vocal tract frequency Rsponse

$\Downarrow$ IDFT

$$X(t) = \cancel{E(t)} + H(t)$$
$\downarrow$
not interested    (low pass lifter)

M-L approach → traditional
→ Deep learning.

Deep learning
↓
unstructured data.

Types of intelligent audio
systems.

Traditional ML          Deep learn
↓                       ↓
feature engineering     automatic feature
                        extraction

Extract audio features.              Feature value / vector
↓                                         ↑
Time - domain feature pipeline       aggregation
↓                                         ↑
Analogue Sound
↓ ADC                                Feature computat
digital signal.          ←            ↑
↓
Framing (making frames of samples)

→ Typical value → (256 — 8192)
→ powers of 2
Frames are constructed to stack samples so that the
duration is humanly percievable.

# Fourier transform.

- Compare signal with sinusoides of various frequencies
- for each frequency we get a magnitude & a phase
- high mag indicates high similarity btw signal & a sinusoid.
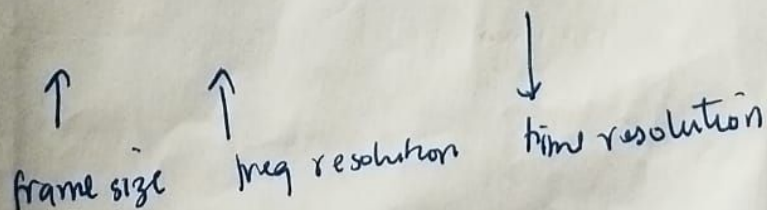
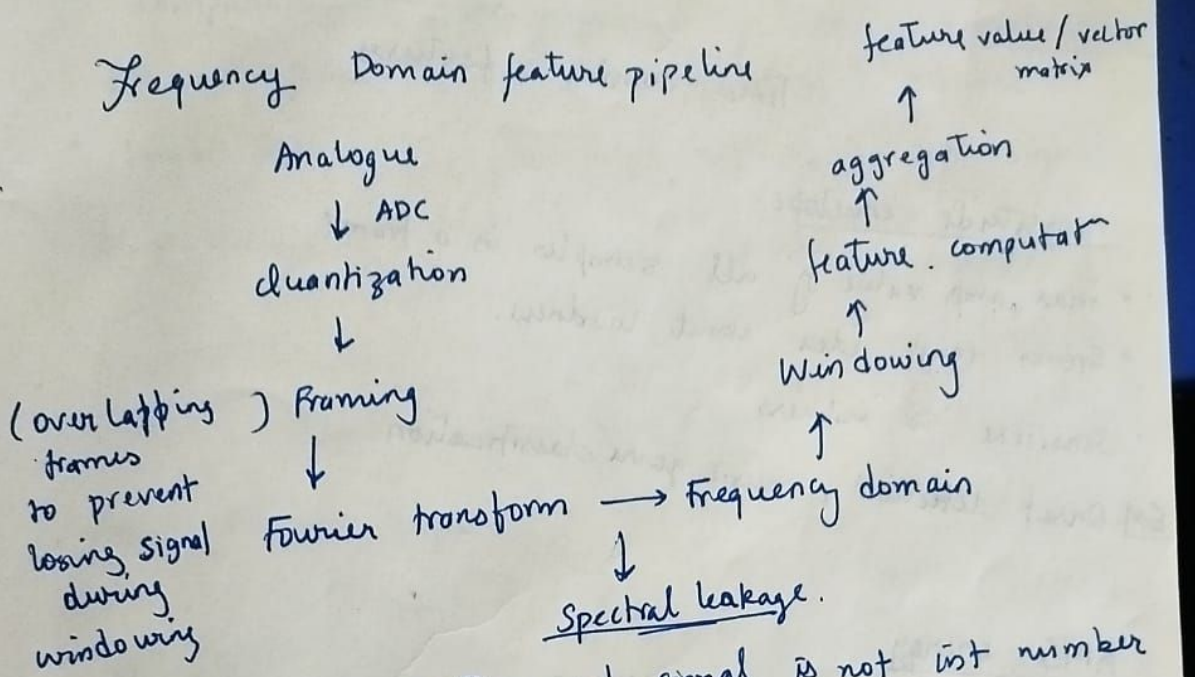## Short Fourier Transform.

- Apply FT for frames.

Outputs.

- DFT
  - Spectral vector ( # frequency bins )
  - N complex Fourier coeff

- STFT
  - Spectral matrix ( # freq bins , # frames )
  - Complex Fourier coeff.

Time / frequency trade off.

↑              ↑                          ↓

frame size    freq resolution    time resolution

Durat$^n$ of
1 sat sample  $<<$  human resolut$^n$

Frequency Domain feature pipeline    feature value / vector
                                                matrix
                                                  ↑
          Analogue                           aggregation
           ↓ ADC                                  ↑
          Quantization                      feature . computat$^n$
             ↓                                    ↑
                                             Windowing
(over lapping ) Framing                           ↑
   frames           ↓
to prevent    Fourier transform → Frequency domain
losing signal                          ↓
during                           Spectral leakage.
windowing
                    ⎧  o Processed signal is not ist number
                    ⎪     of periods.
                    ⎨
                    ⎪  o End points are discontinuous.
                    ⎩
                       o Discontinuities appear as high-freq
                         components not present in the org signal

Windowing: ⎧ Applying windows function to each frame
           ⎪ eliminates samples at both ends of a frame
           ⎨ Generates a periodic signal.         ↓
           ⎩                                    To make
                         ↓                      it$^k$ perodic.
                    minimizes spectral
                         leakage.

popular funt$^n$: Honn window.

hop length : over lap betw fff win frames.

---

## Time domain features

### Amplitude envelope.
- max amp value of all samples in a frame.
- Gives rough idea about loudness.
- Sensitive to outliers

Ex) Onset detection, music genre classification

### RMS - energy.
- RMS of all samples in a frame.
- Indicator of loudness
- less sensitive to outliers than AE

Ex) Audio segmentatⁿ, music genre classification

### Zero crossing Rate
- number of times a signal crosses $x = 0$
- Ex   Recognizing ~~percu~~ percussive vs pitched
  sounds
  Speech
  Recognition

Fourier
- Compa
- for each
- high
  N a s

- Apply

Output

- DFT
  o Spe
  o N c

- STFT
  o Sp
  o Cor

Time /

↑
frame