

Copy

Spectrogram

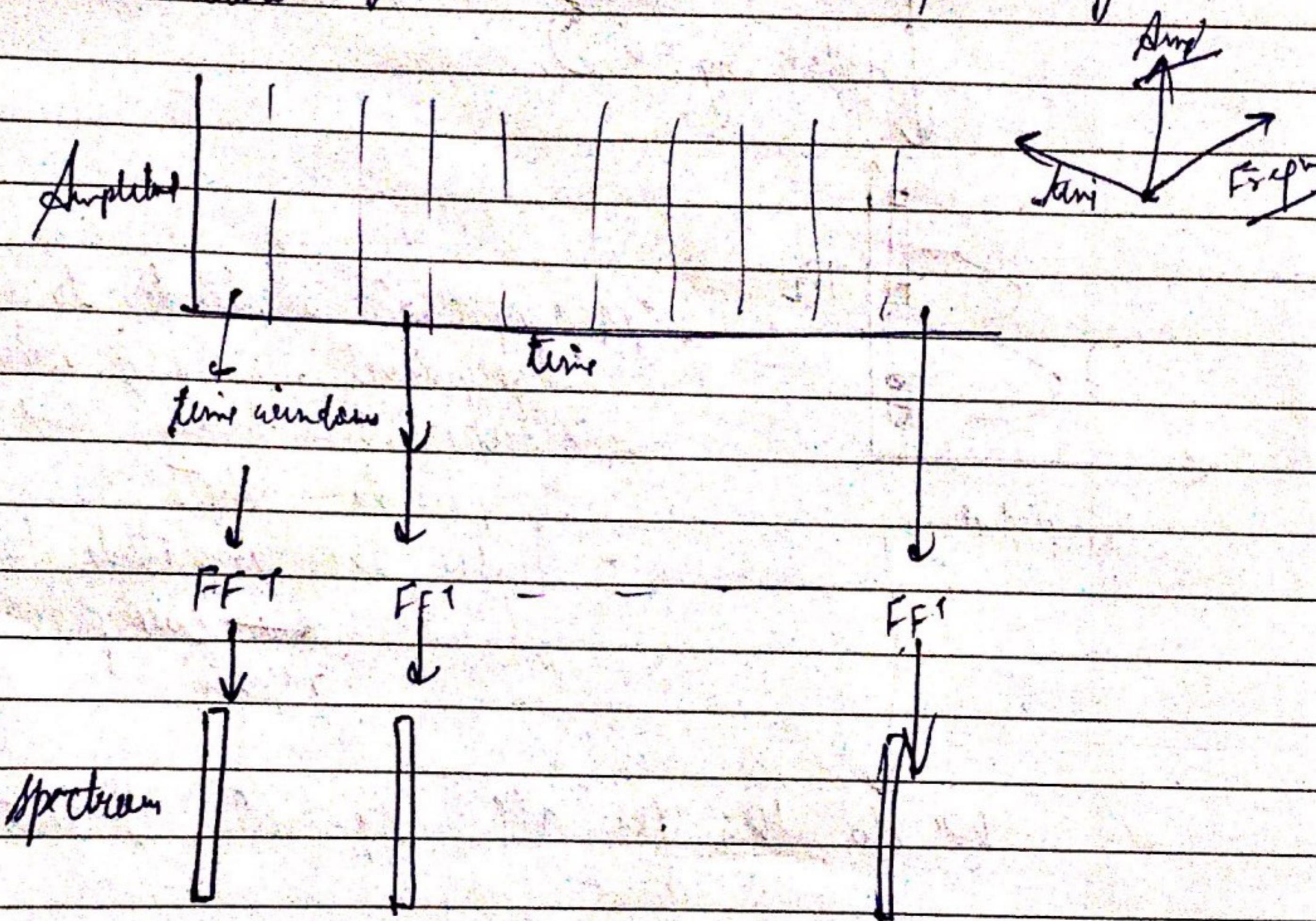
Cepstrum

Mel-frequency analysis

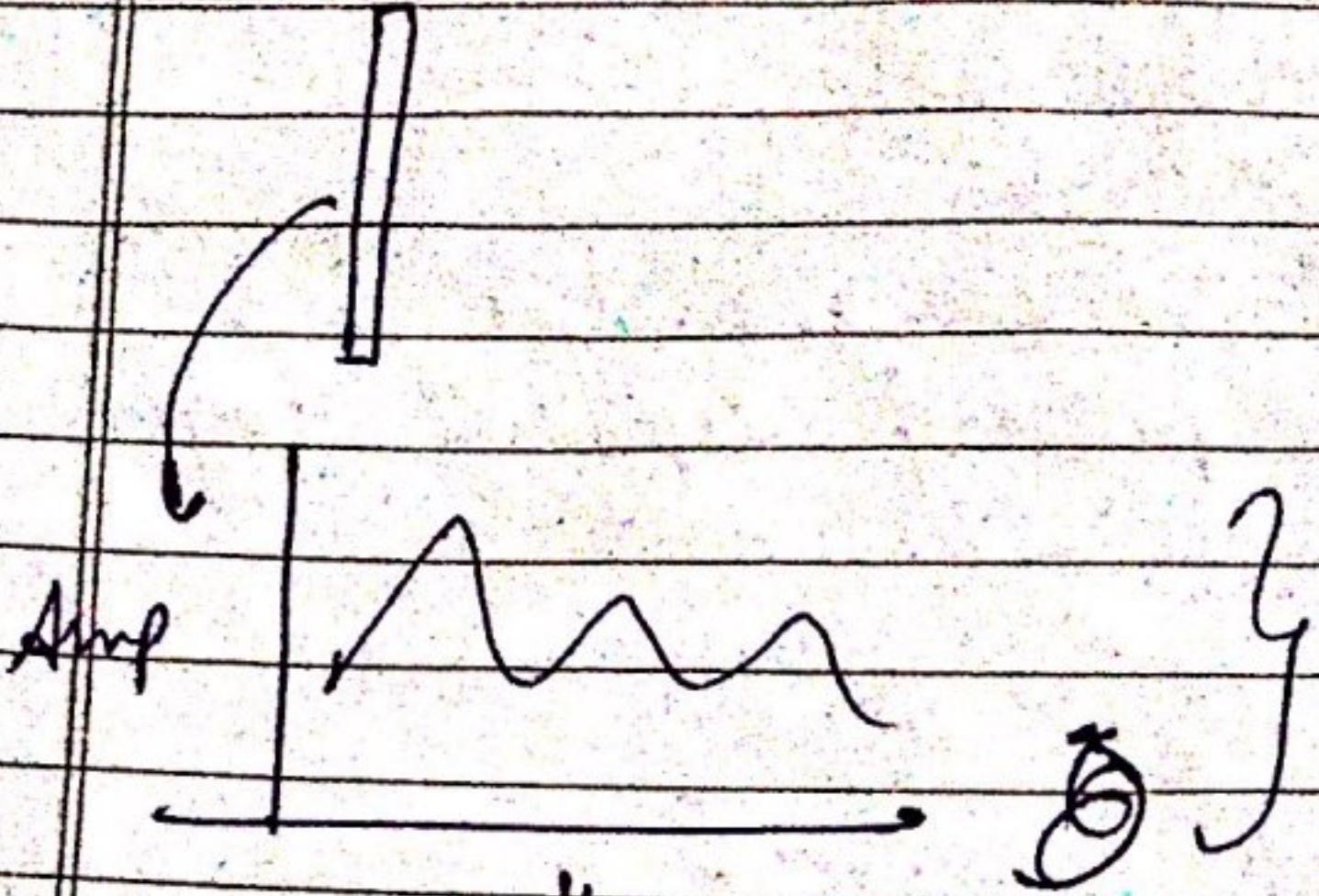
Mel-frequency cepstral coefficients

Spectrogram

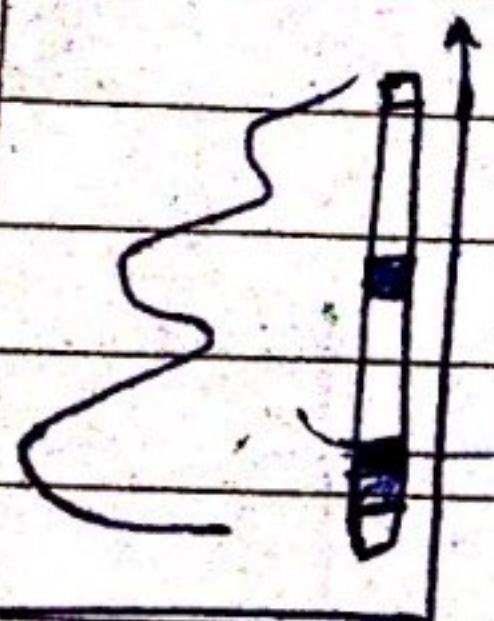
Speech signal represented as sequence of spectral vectors



→ Divide wave plot into ~~moving~~ windows and apply FFT we get spectral vectors ~~and~~ ~~feature~~



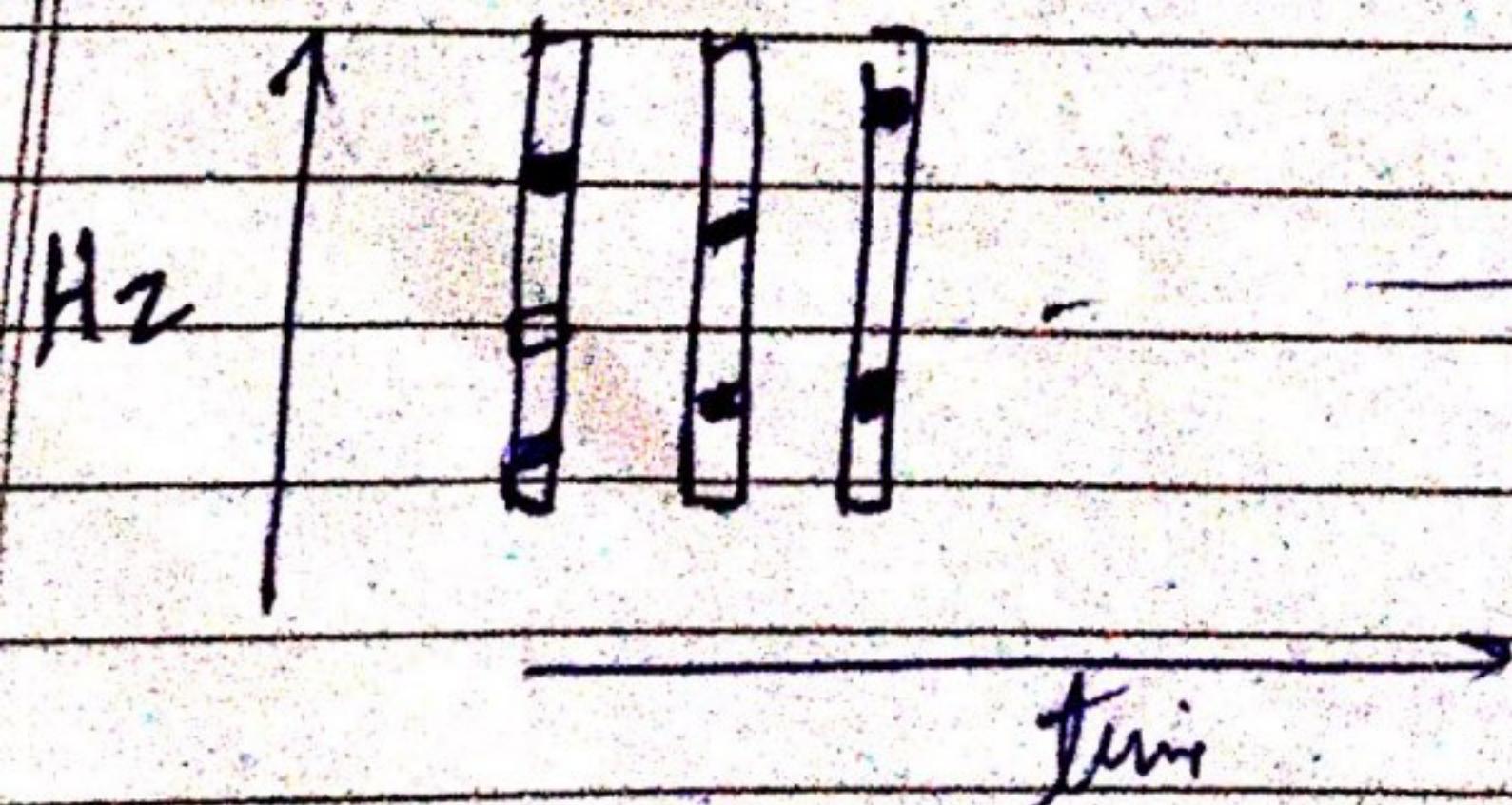
rotate by 90°



- map spectral amplitude to a grey level (0-255) value.
- 0 represents black and 255 represents white.
- Higher the amplitude darker the color corresponding region

amplitude

we do this for all spectral vectors



(Time vs Frequency representation of speech signal is referred to as spectrogram)

→

WHY SPECTROGRAMS

(speech sounds)

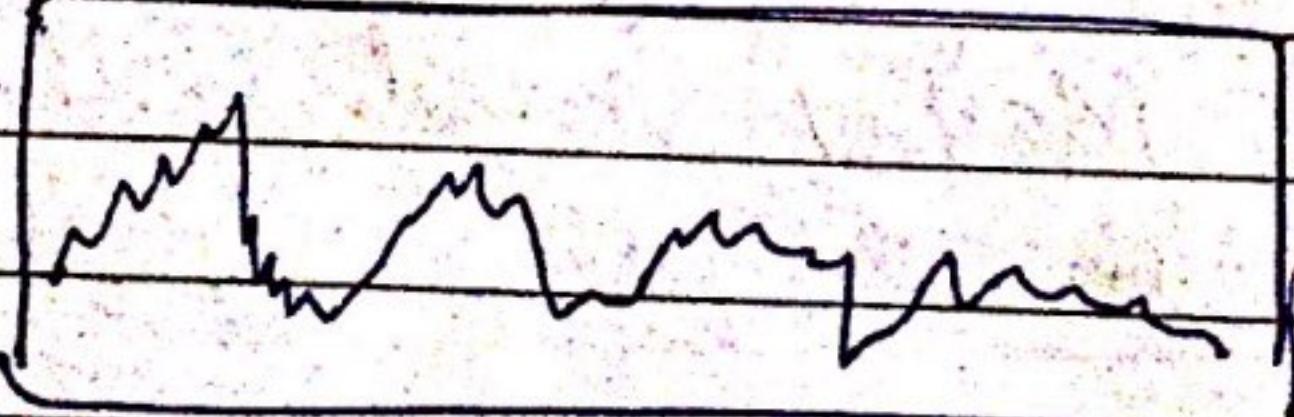
- Phones, and their prosodies are better observed in spectrograms.
- Easier detecting phone boundaries.
- Sounds can be identified by the formants and by their transitions.
→ dark region (high amplitude)
- Formants carry the identity of sound.

Note: Hidden Markov models implicitly model these spectrograms to perform speech recognition.)

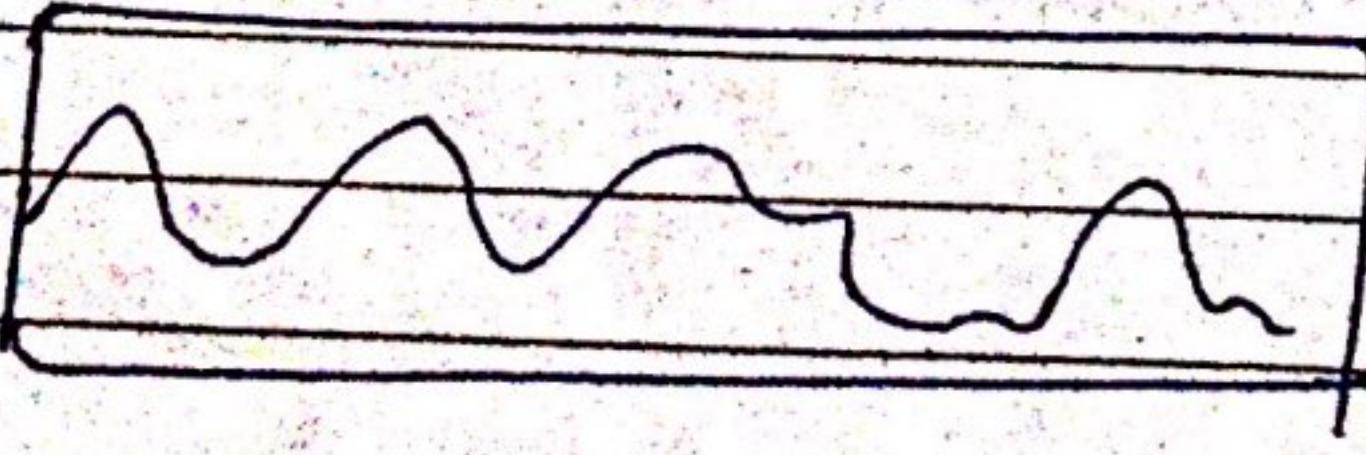
Usefulness of spectrogram

- Time-Frequency representation of the speech signal
- Spectrogram is a tool to study speech sounds (phones)
- Phones and their properties are visually studied by phoneticians
- HMM - Simplicity models spectrograms for speech to text systems
- Useful for evaluation of text to speech systems
 - A high quality text to speech system should produce synthesized speech whose spectrogram should nearly match with the natural sentences.

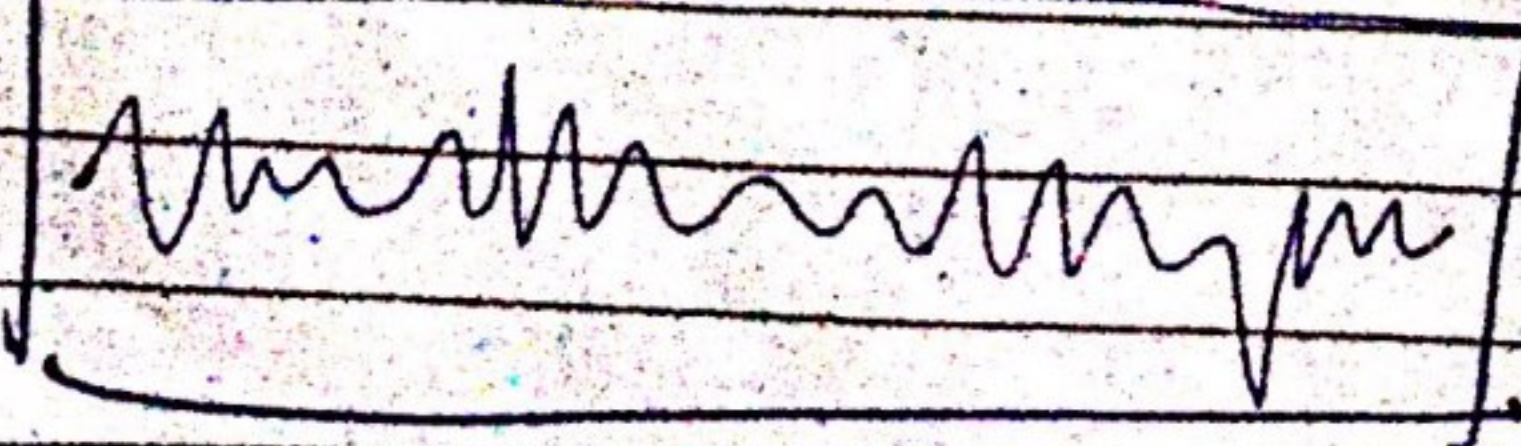
Spectrum



Spectral Envelope db



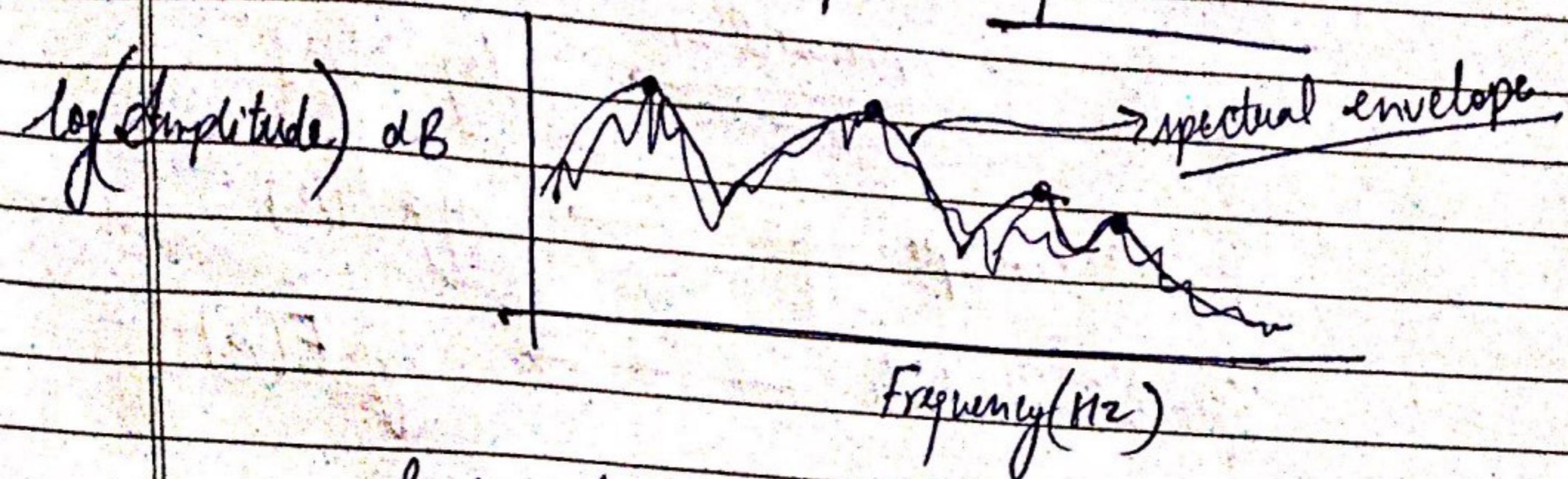
Spectral details



Hz

Cepstral Analysis

A sample speech spectrum



- Peaks denote dominant frequency components in speech signals
- Peaks are referred to as formants
- Formants carry identity of sound

What we want to Extract

Spectral Envelope

- Formants and a smooth curves connecting them
- This smooth curve is referred to as spectral envelope

Find peaks and curves connecting them



After Extracting spectral envelope from spectrum what is left there is spectral details

$$\text{Spectrum} \xrightarrow{\text{represented as}} \log X[k]$$

$$\text{Spectral envelope} \xrightarrow{} \log H[k]$$

$$\text{Spectral details} \xrightarrow{} \log E[k]$$

$$\log X[k] = \log H[k] + \log E[k]$$

- 1) Goal :- separate spectral envelope and spectral details from the spectrum
- 2) Given $\log X[k]$ separate $\log H[k]$ and $\log E[k]$

How to achieve separation

As we know spectrum is already a FFT of wave plot. We have converted amplitude into log/amplitude).

An FFT on spectrum referred to as IFFT (Inverse FFT)

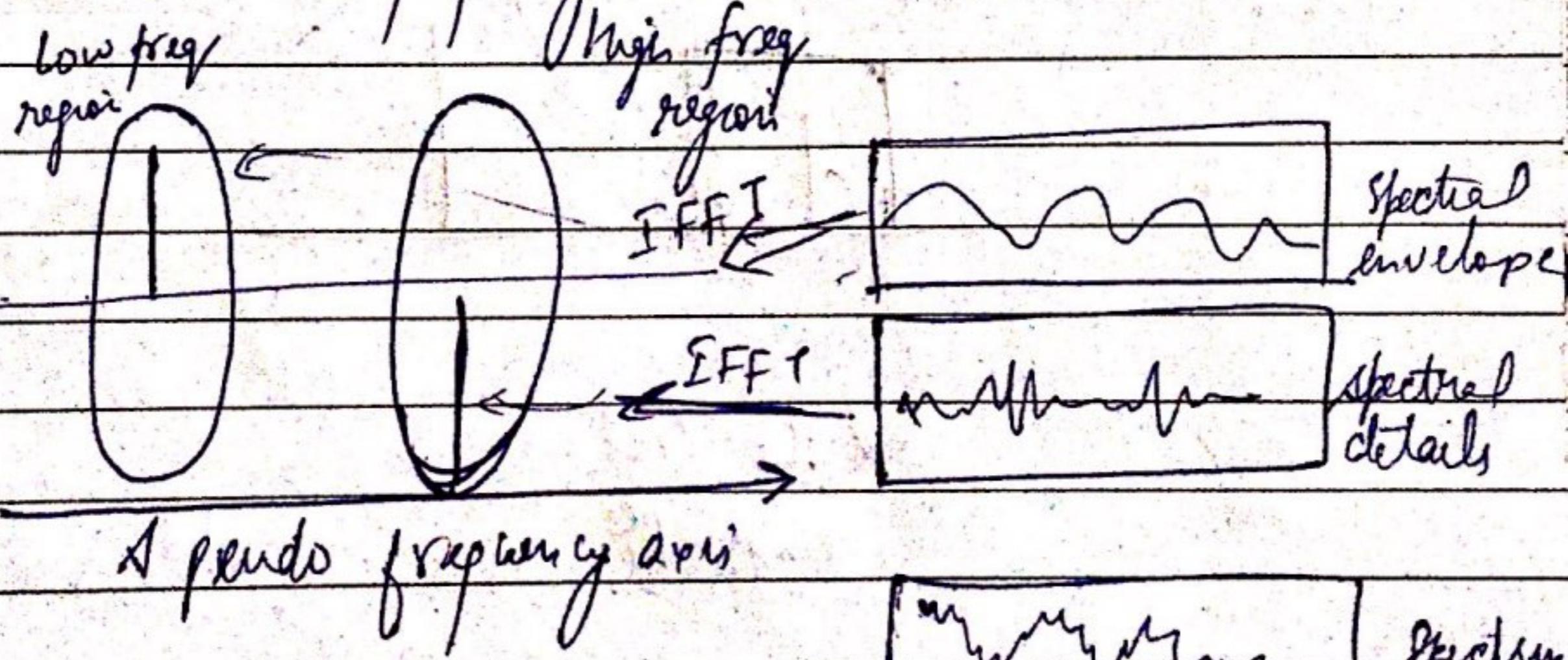
$$\text{FFT}(\text{FFT}(\text{waveplot})) = \text{waveplot}$$

But after taking $\text{FFT}(\text{waveplot})$ we always get y-axis (amplitude) so $\text{IFFT}(\text{spectrum})$ no more return frequency.

• IFFT of log spectrum would represent the signal in pseudo-frequency axis

Bottom-up approach

assume spectral envelope as a sine wave with four cycles per sec (assume x-axis as time). If we take its FFT then it gives a peak at 4Hz in frequency axis.



Similarly assume spectral details as sine wave with 100 cycles per second. If we take IFFT we get peak at 100Hz

$$x[k] = w[k] + e[k]$$

IFFT

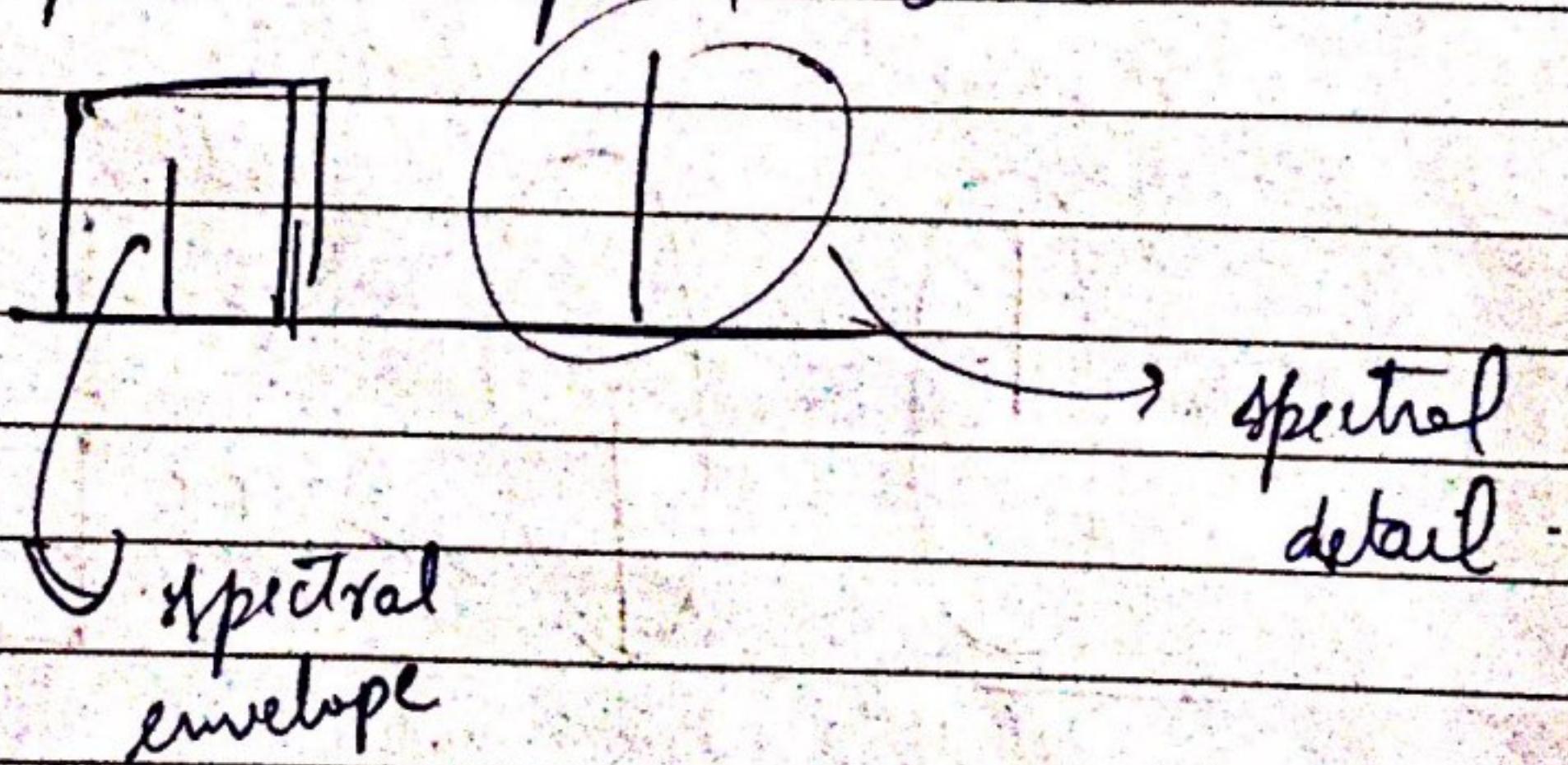
$\log H[k]$ | spectral envelope

spectral detail $\log E[k]$

A pseudo-freq. axis

2 But we have access to only $\log X[k]$ and hence you can obtain $x[k]$

and it have a lower freq and higher freq component that correspond to spectral envelope and spectral detail.



3 If you know $x[k]$ filter the low frequency regions to get $h[k]$.

$x[k]$ is referred to as CEPSTRUM

REP h[k] represents the spectral envelope and is widely used as feature for speech recognition.

Mel-Frequency Analysis

Speech recognition
behaves like a human
ear

Review : What we did

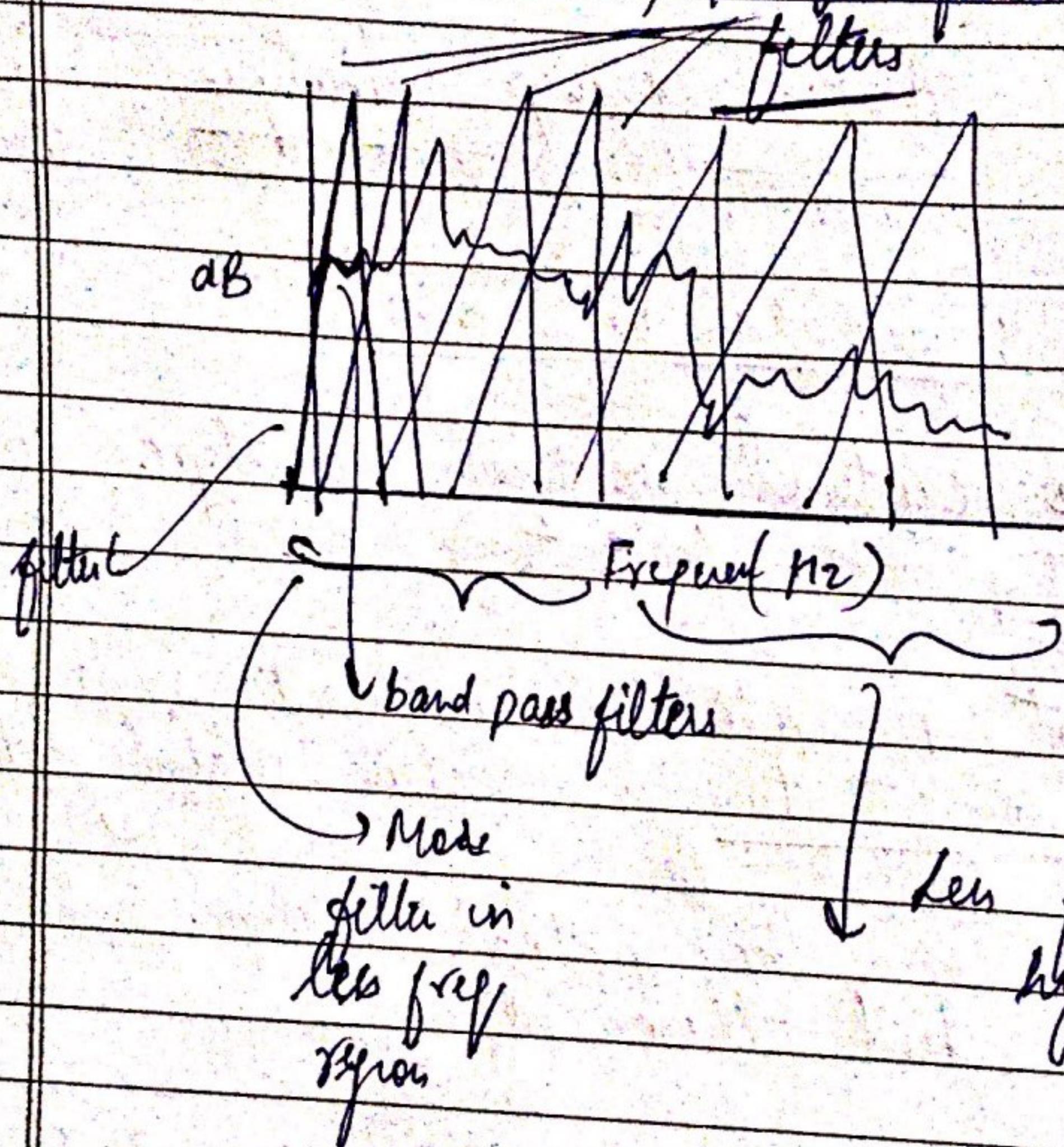
- We captured spectral envelope (curve connecting all formants)
- But perceptual experiment says human ear concentrates on certain regions rather than whole of the spectral envelope. So how to find those regions.

(Concentration points also known as speech perception)

- #
- ⇒ Mel-Frequency analysis of speech is based on human perceptual experiments
 - ⇒ It is observed that human ear acts as a filter and concentrates on only certain frequency components.
 - ⇒ These filters are non uniformly spaced on frequency axis i.e. more filters at low freq

regions. and less filters in high frequency regions.

MEL-frequency filter



-1 filter only pass freq ~~inside~~ not fall wind

-1 41 filter are found to be effective

There 41 filters set have particular position to be placed based on Mel-Scale

Mel-Frequency Cepstral Coefficients (MFCC)

wavelet \xrightarrow{FFT} Spectrum $\xrightarrow{M.F. \text{ filter}}$ MEL spectrum ~~- inverse~~

$$\log X[k] = \log (\text{MEL spectrum})$$

$$\log X[k] = \log N[k] + \log E[k]$$

Taking IFFT (we get low freq and high freq region)

$$x[k] = h[k] + e[k]$$

Cepstral coefficient $h[k]$ obtained for
MEL spectrum are called MFCC

2) Apply this to all spectrum we
get MFCCs for all time windows called
cepstral vectors.