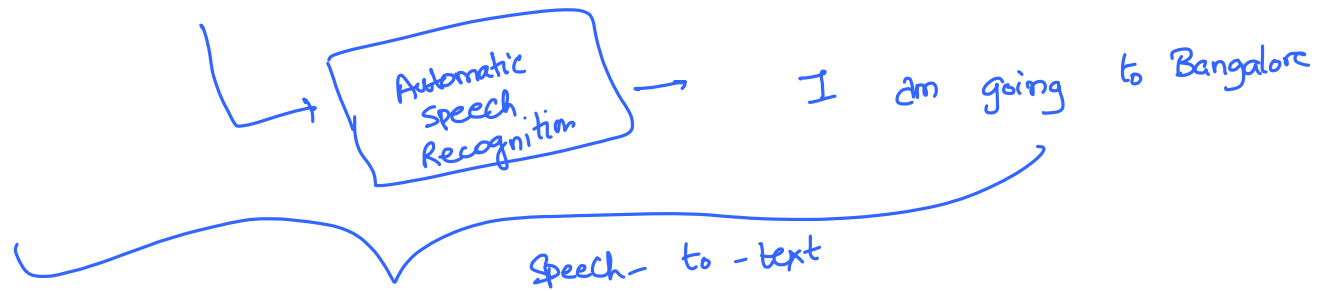


Week-2 for Speech Part → DLP- 6/11

Speaker Diarisation



Meeting

~~Handwritten scribble~~



② Speaker-Independent Mode

③ Speaker Diarization
→ Who spoken when

→ Umesh: Welcome to this class
→ Student 1: Thank you
→ Mihir: we expect you to have background in Deep Learning

① Speech-to-Text \Rightarrow Whisper ASR Model \Rightarrow 680,000 hrs
 \Rightarrow open AI



2:03

2:08

This is a Test } ✓

2:08

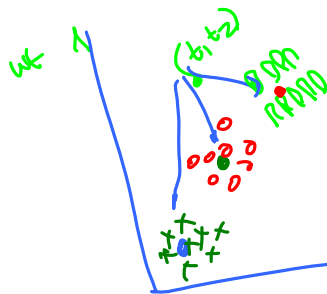
2:17

Therefore the success depends
on your input ✓

② Speaker-Identification Module
 \Rightarrow Speaker-Embedding Extractor
 "Speaker-feature" Extractor

} Speech Brain
Tool Kit
ECAPA-TDNN

| Speech Tool-kits
Kaldi \rightarrow JHU
Esnet \rightarrow JHU/CMU
Speech } \rightarrow MILA
Brain }



$$\begin{bmatrix} h_1 \\ w_1 \end{bmatrix} \leftrightarrow \begin{bmatrix} h_2 \\ w_2 \end{bmatrix} \quad \begin{bmatrix} h_3 \\ w_3 \end{bmatrix}$$

$$\begin{bmatrix} t_1 \\ t_2 \end{bmatrix}$$

"closeness"
 \rightarrow Euclidean Distance

$$\underline{x}_1 = \begin{bmatrix} x_{11} \\ x_{21} \end{bmatrix} \quad \underline{y}_1 = \begin{bmatrix} y_{11} \\ y_{21} \end{bmatrix}$$

$$\|\underline{x}_1 - \underline{y}_1\| = \sqrt{(x_{11} - y_{11})^2 + (x_{21} - y_{21})^2}$$

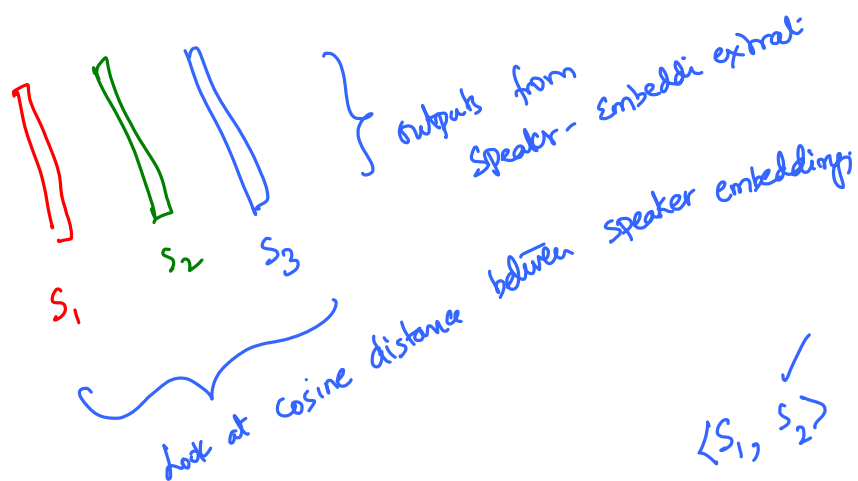
Cosine distance

$$\cos \theta = \frac{\underline{x}_1^T \underline{y}_1}{\|\underline{x}_1\| \|\underline{y}_1\|}$$

$$\underline{x}_1$$

$$\underline{y}_1$$

$$\underline{x}_1$$



$$\langle s_1, s_2 \rangle = \frac{s_1 s_2}{\|s_1\| \|s_2\|}$$

$\langle s_1, s_3 \rangle$

$\langle s_2, s_3 \rangle$

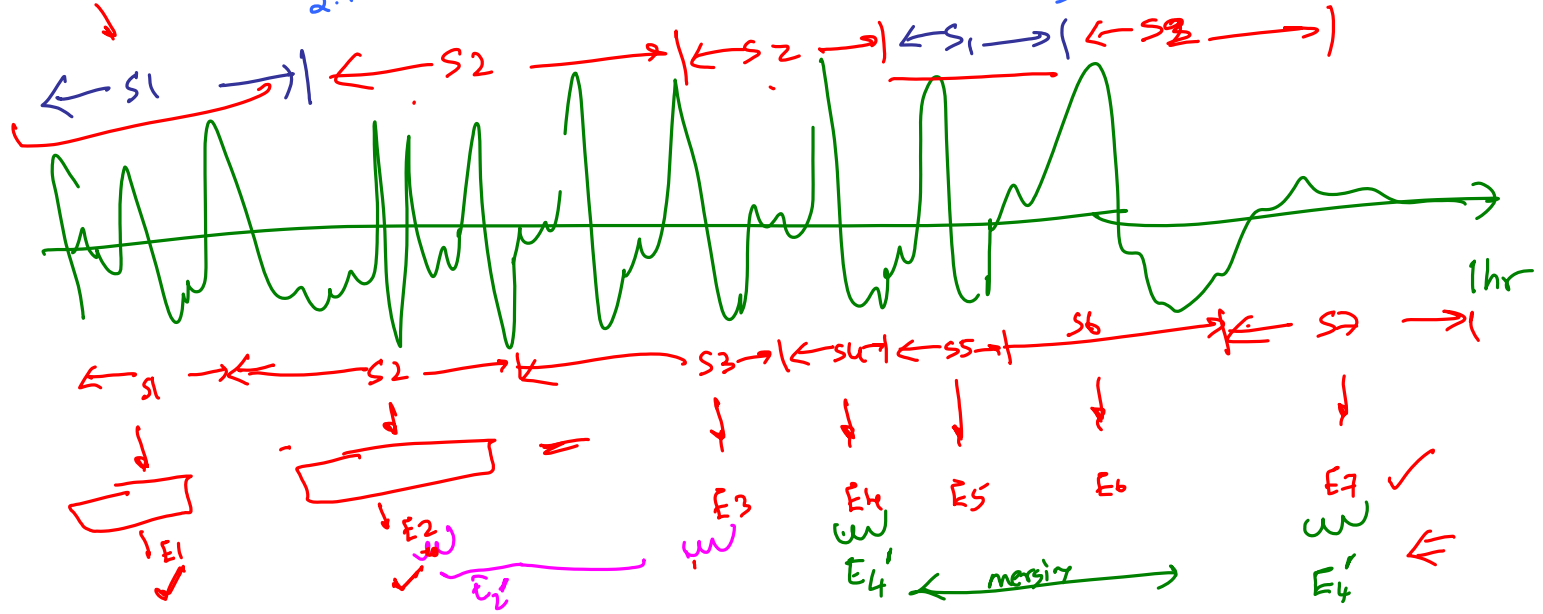
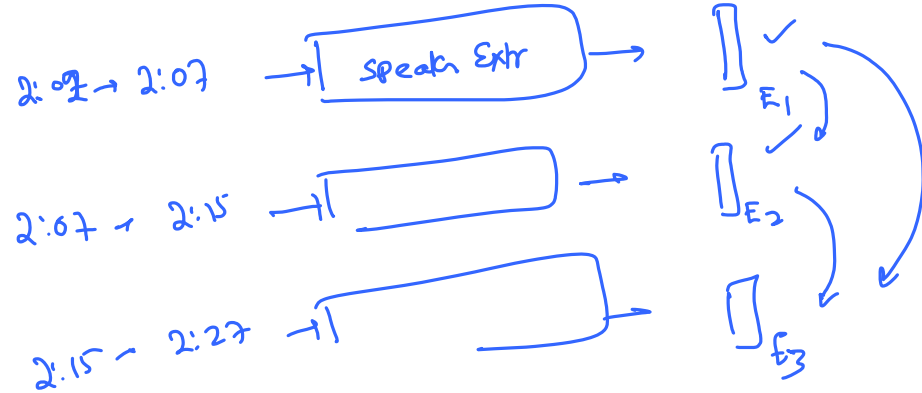
Recall: Speaker Diarisation: "Who spoken When"

① → ~~AAAA~~ speech - sign → Whisper

s_1	2:02	2:07] This is a test	
s_1	2:07	2:15		} Good Test
s_2	2:15 - 2:27			

How are you

②



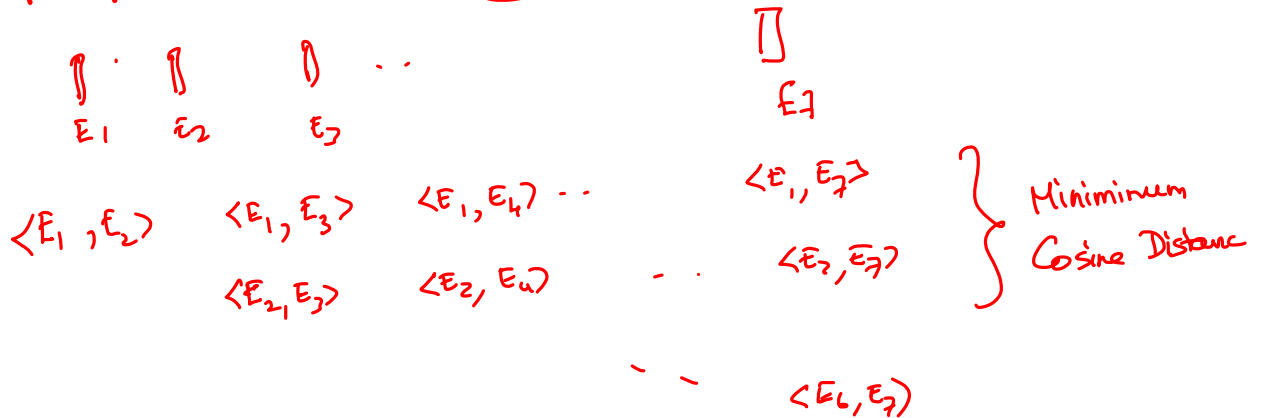
I know only "2 speakers"

AGGLOMERATIVE CLUSTERING

⇒ "clustering" ⇒ and keep merging clusters
 ⇒ until I get "desire number of clusters"

Initial

7 segments → 7 speakers



let us say E_4, E_7 has min. cosine distance