

• Voice Recognition
Uses SST along with
Speaker recognition

SST (speech to text)

Falls under
audio processing

- SST (speech to text) \Rightarrow input is spoken audio \rightarrow output: text

• Goal: convert human speech into text

Steps 1) Audio capture

• Microphone records audio \rightarrow sound waves sampled (ex 16KHz)
 \rightarrow stored as a wave form

2) Feature Extraction

• Raw audio is too messy for processing directly so it's transformed into features

• Common: MFCC (Mel Frequency Cepstral Coefficients) - a way of representing sound that matches human hearing sensitivity

• think of it as: condensing sound into a compact fingerprint for each small time slice

3) Acoustic Model

• maps short audio frames \rightarrow Phonemes (basic sound units in language)

• old methods: HMMs + GMMs • Modern method: CNN, RNN, transformers

4) Language model

• helps guess the most probable sequence of words

Ex: acoustic model: "recognize speech" raw model output: "wrench a nice beach"
Then the model picks the more likely option

5) Decoding: combines acoustic and language models to produce final text

- Modern Advancements

Ex whisper \rightarrow End-to-End models: skip phoneme steps, directly map audio \rightarrow text using architectures like RNN-transducer or transformers

Multilingual: Model trained on many languages at once

Robustness: Noise suppression, speaker separation, diarization (who said what)