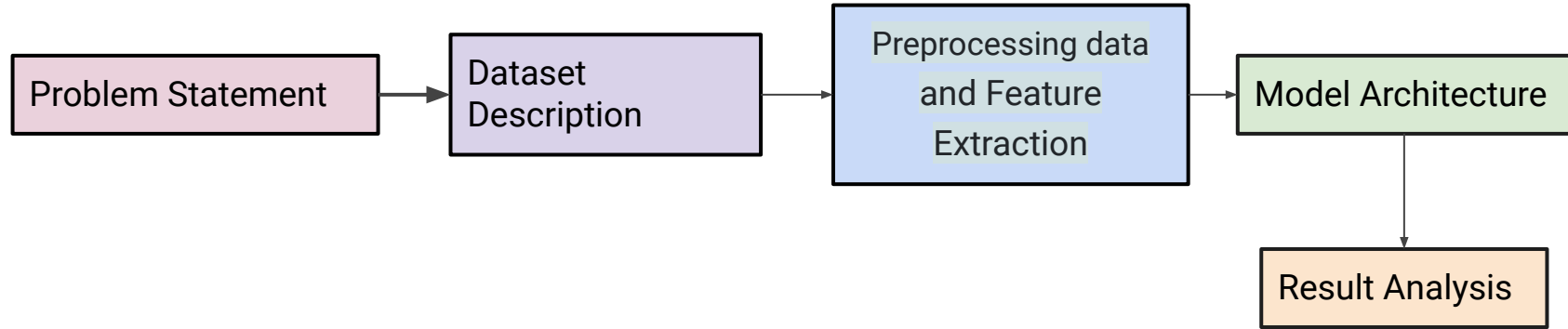# Bengali Spoken Digit Classification: A Hidden Markov Model Approach

By
Dibyendu Das
Under the Guidance Of
Dr. Sujoy Biswas
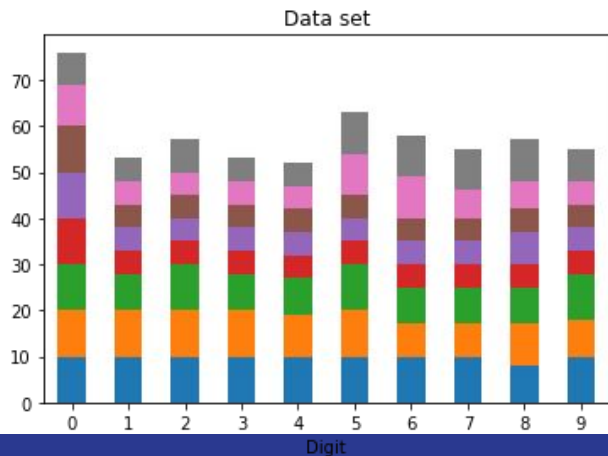
M.Sc in Big Data Analytics
Ramakrishna Mission Vivekananda Educational and Research Institute

Problem Statement → Dataset Description → Preprocessing data and Feature Extraction → Model Architecture → Result Analysis

Dataset Description

# Dataset Description

- A dataset containing 600 audio file (.wav format) was created for the experiment.
- Eight people from various parts of the State were asked to give their voice recordings Using "QuickRec" App.

| Bengali word | Bengali pronunciation | English word | English numerical |
|---|---|---|---|
| শূন্য | shun-no | zero | 0 |
| এক | a-k | one | 1 |
| দুই | du-i | two | 2 |

Data set

Preprocessing and Feature Extraction

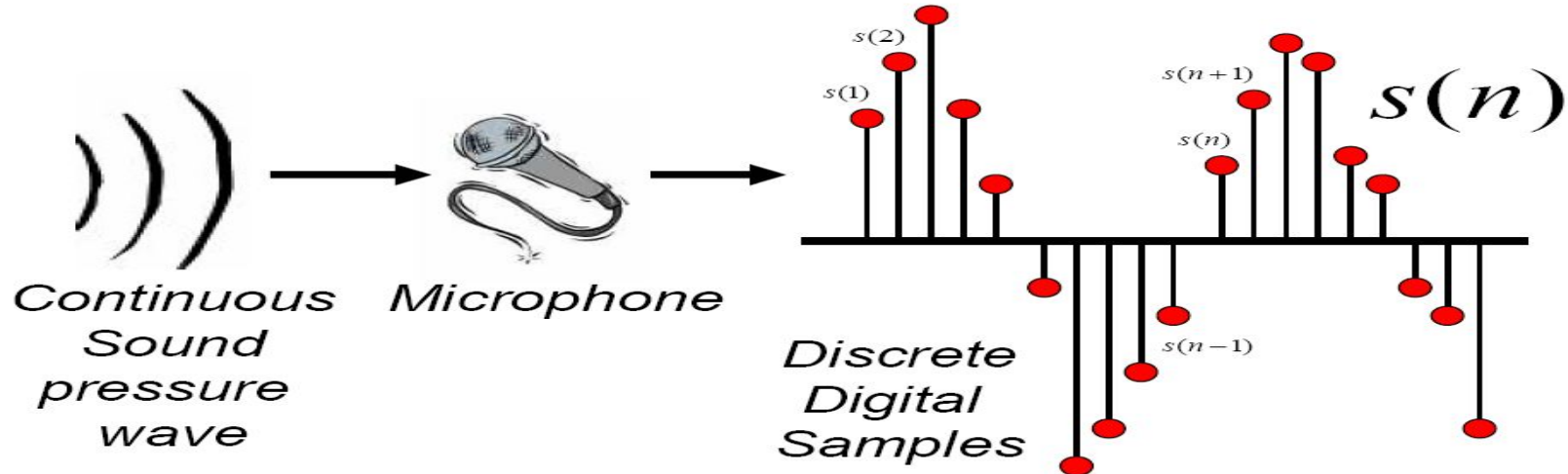# Discrete Representation of Signal :
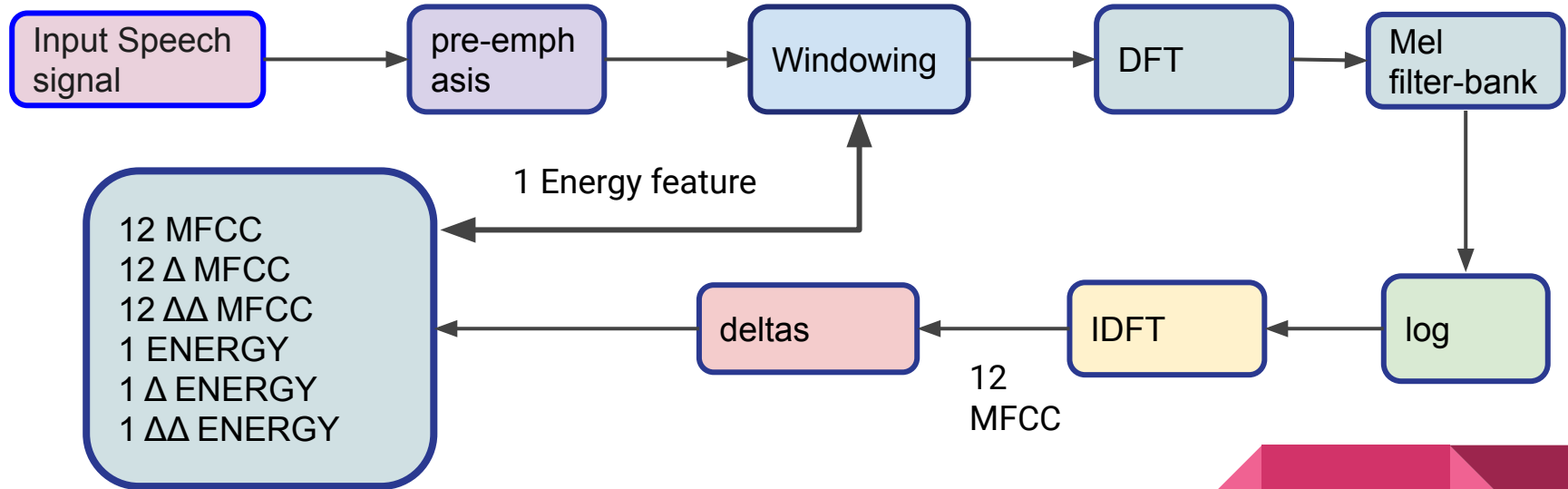
Analog-to-digital conversion has two steps :

Sampling = 2 * Nyquist frequency

**+**

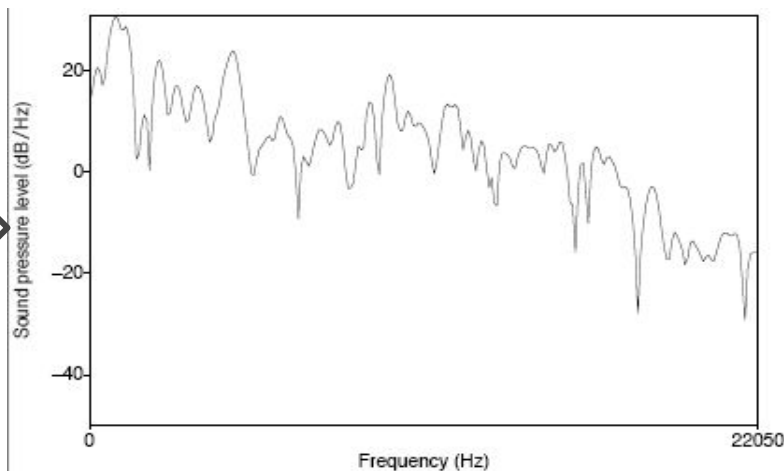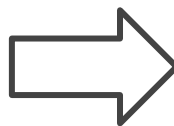Quantization = store the amplitude value as 8 bit or 16 bit



Continuous Sound pressure wave → Microphone → Discrete Digital Samples

$s(1)$  $s(2)$  $s(n)$  $s(n+1)$  $s(n-1)$  $s(n)$

# Preprocessing data and MFCC Feature Extraction :

Input Speech signal → pre-emphasis → Windowing → DFT → Mel filter-bank

1 Energy feature

12 MFCC
12 Δ MFCC
12 ΔΔ MFCC
1 ENERGY
1 Δ ENERGY
1 ΔΔ ENERGY

← deltas ← IDFT ← log ← Mel filter-bank
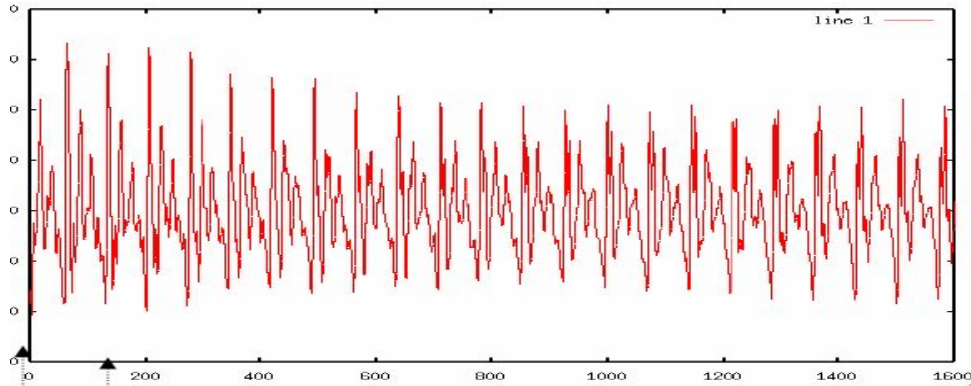
12 MFCC

# Preemphasis :

- The spectrum for voiced segments has more energy at lower frequencies than higher frequencies. This is called **spectral tilt**
- Spectral tilt is caused by the nature of the glottal pulse
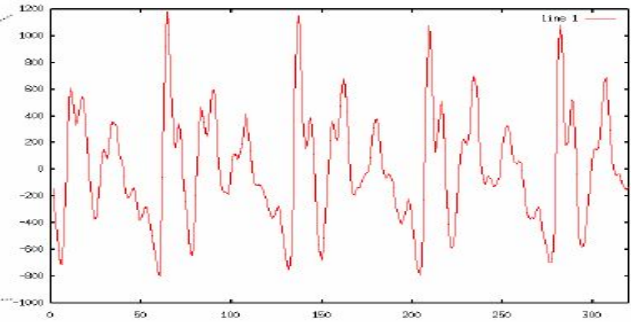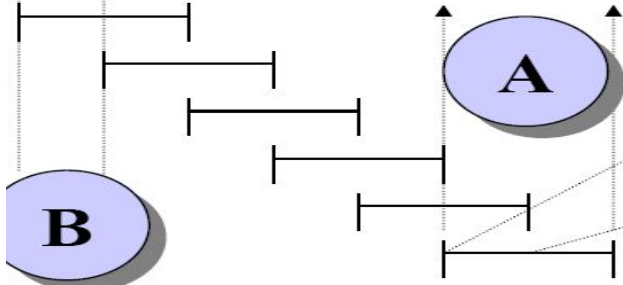
Spectral slice from the vowel [aa]
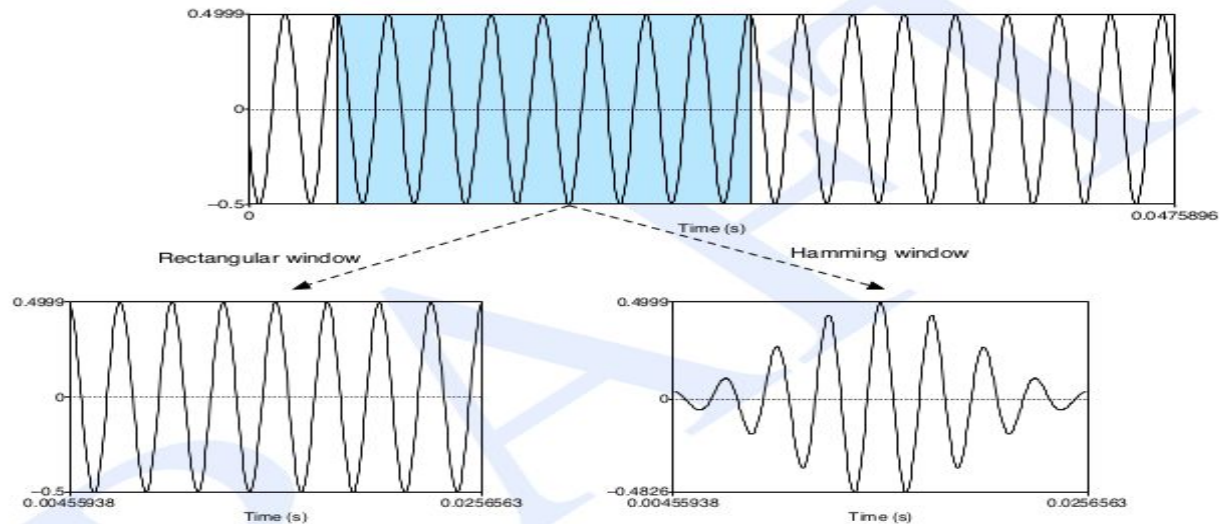
# Windowing :



$$A \sim 20 - 25 \text{ ms}$$

$$B \sim 10 \text{ ms}$$

# Common window shapes :
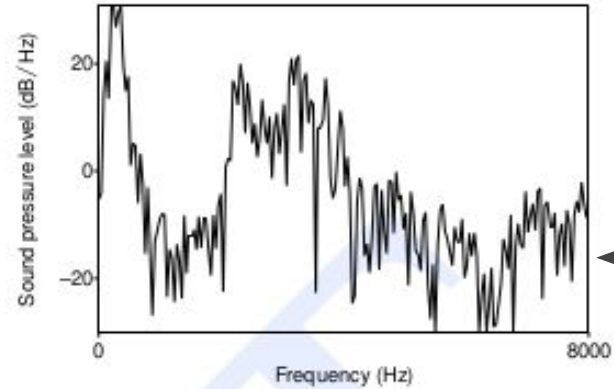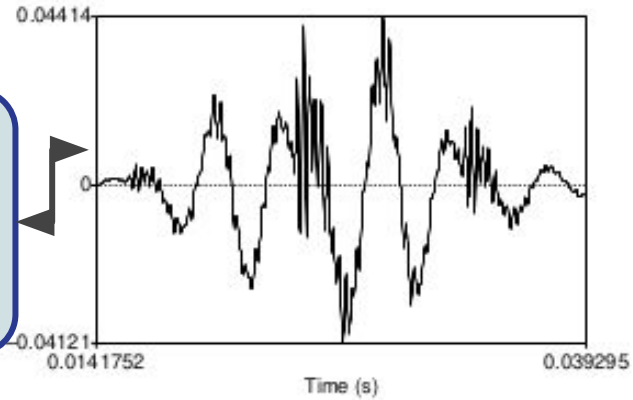
$$rectangular \quad w[n] = \begin{cases} 1 & 0 \le n \le L-1 \\ 0 & \text{otherwise} \end{cases}$$

$$hamming \quad w[n] = \begin{cases} 0.54 - 0.46\cos\left(\frac{2\pi n}{L}\right) & 0 \le n \le L-1 \\ 0 & \text{otherwise} \end{cases}$$
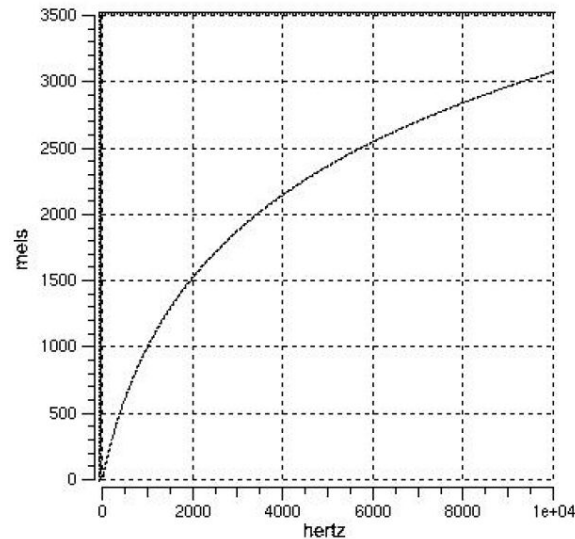
# DFT :

25 ms windowed portion

Spectrum computed by DFT

# Mel Scale :

- **Human hearing is not equally sensitive to all frequency bands**
- **Less sensitive at higher frequencies, roughly > 1000 Hz**
- **I.e. human perception of frequency is non-linear:**

A mel is a unit of pitch
**Definition:**  Pairs of sounds perceptually equidistant in pitch
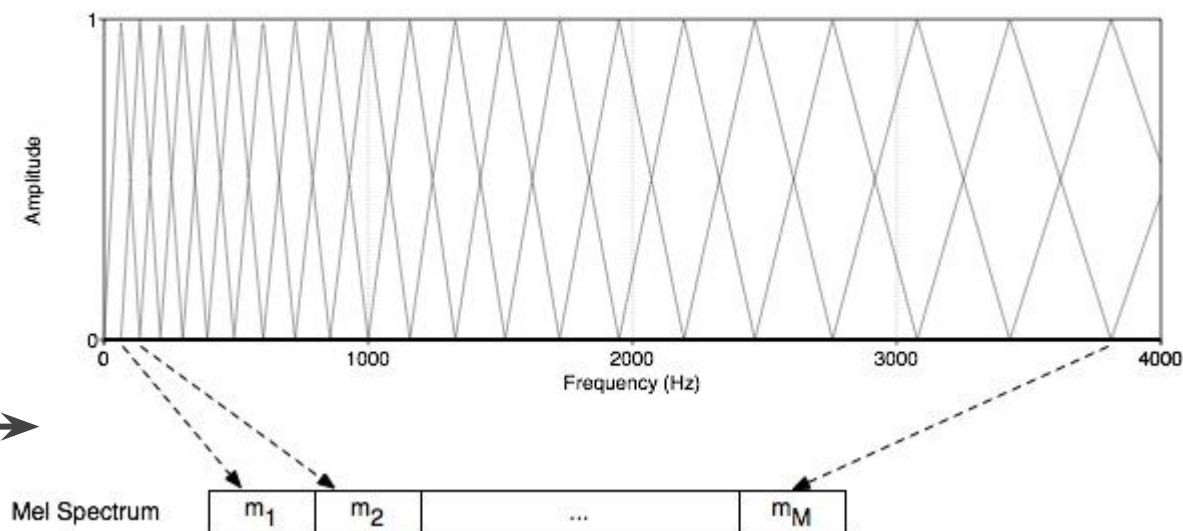Are separated by an equal number of mels

$$\text{Mel } (f) = 1127 \ln \left[ 1+(f/100) \right]$$

Mel Filter bank
Uniformly spaced before 1 kHz
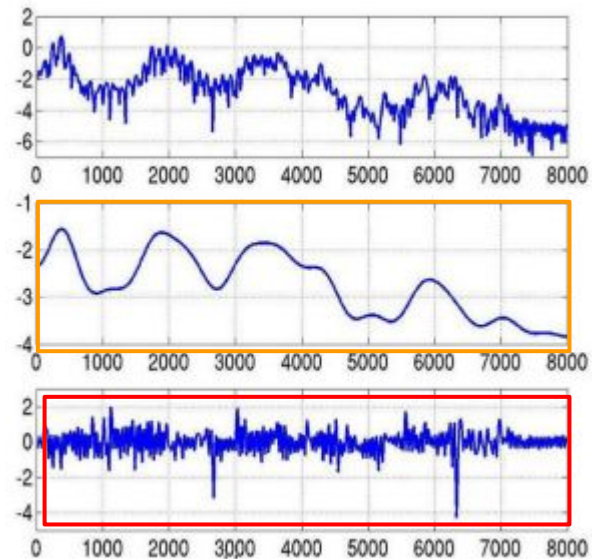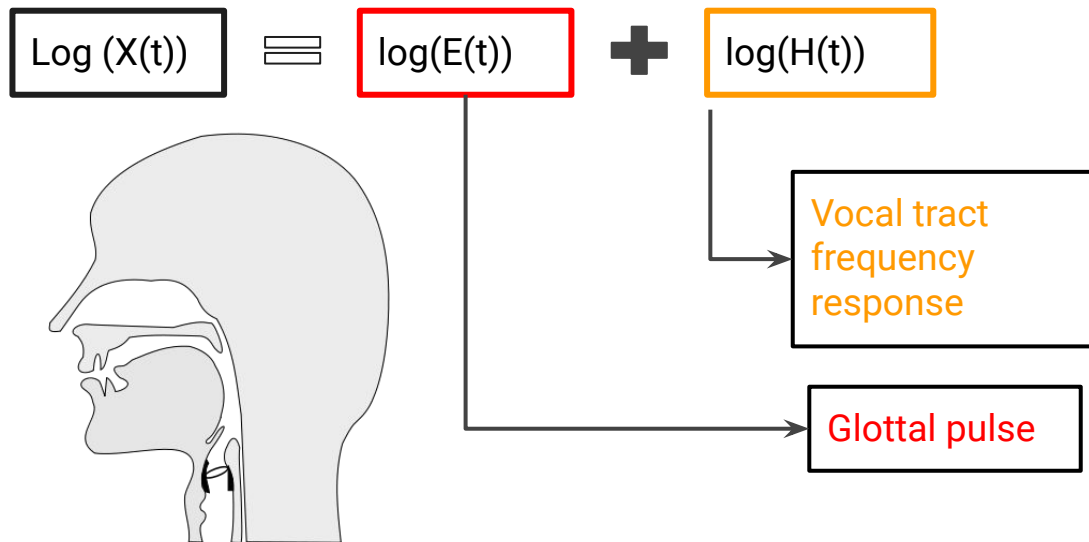logarithmic scale after 1 kHz

Each triangular filter collects energy from given frequency range

# The Cepstrum : IDFT

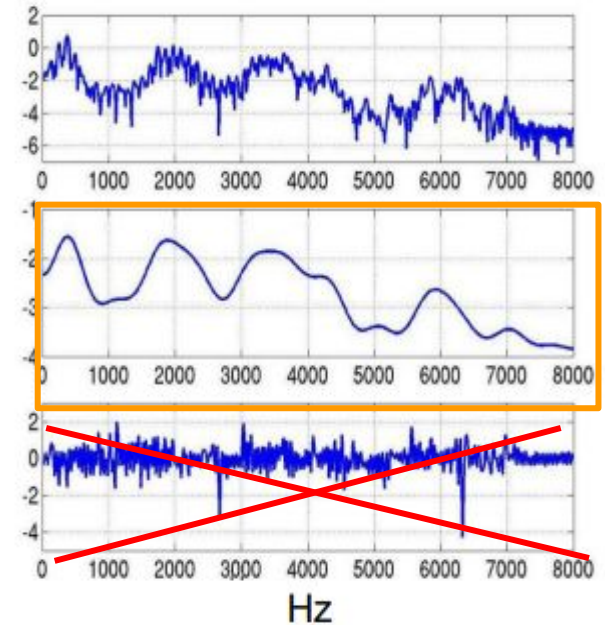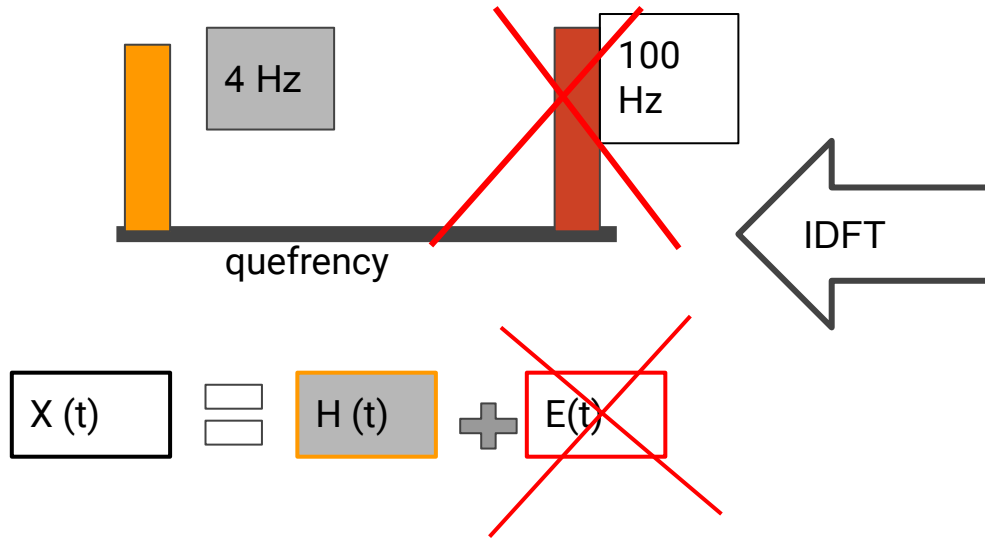Speech = Convolution of vocal tract frequency response with glottal pulse

$x(t) = e(t) * h(t)$  ⟶  $X(t) = E(t) * H(t)$  ⟶  $\log(X(t)) = \log(E(t)) + \log(H(t))$

Log (X(t))  $\equiv$  log(E(t))  ✚  log(H(t))

Vocal tract frequency response

Glottal pulse

HZ

# The Cepstrum :



4 Hz

100 Hz

quefrency

IDFT

X (t) = H (t) + E(t)

Hz

# Feature :

- The cepstral coefficients do not capture energy
- So we add an energy feature
- Also, we know that speech signal is not constant
- So we want to add the changes in features (the slopes).
- We call these delta features

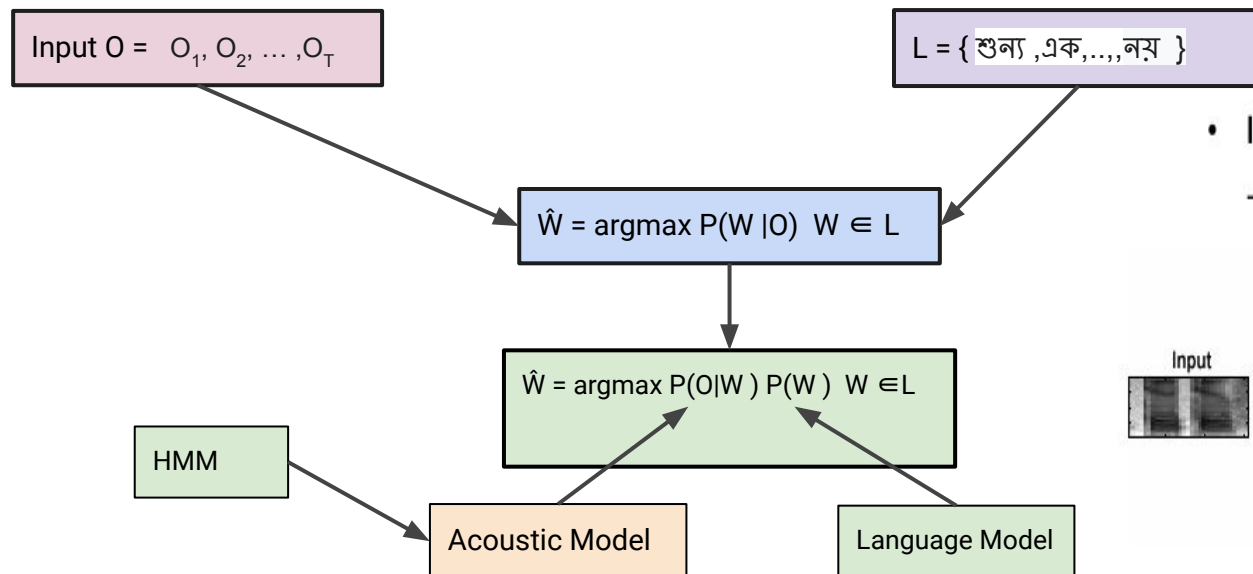$$d(t) = \frac{c(t+1) - c(t-1)}{2}$$

c(t) = cepstral value at time t

12 MFCC
12 Δ MFCC
12 ΔΔ MFCC
1 ENERGY
1 Δ ENERGY
1 ΔΔ ENERGY

- We also add double-delta acceleration feature

HMM Architecture

# Our main goal :

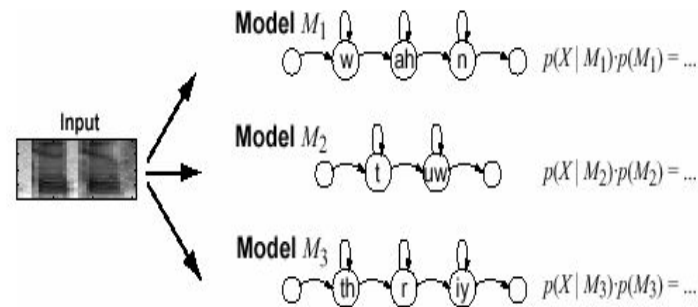**"What is the most likely word out of all words in the language L given some acoustic input O?"**

Input O = $O_1, O_2, \ldots, O_T$

L = { শূন্য ,এক,..,,নয় }

$\hat{W}$ = argmax P(W |O)  W ∈ L

$\hat{W}$ = argmax P(O|W ) P(W )  W ∈ L

HMM

Acoustic Model

Language Model

- **Isolated word**

  - choose best $p(M|X) \propto p(X|M)p(M)$

Model $M_1$  $p(X \mid M_1) \cdot p(M_1) = \ldots$

  w  ah  n

Model $M_2$  $p(X \mid M_2) \cdot p(M_2) = \ldots$

  t  uw

Input

Model $M_3$  $p(X \mid M_3) \cdot p(M_3) = \ldots$

  th  r  iy

# Overall Architecture :
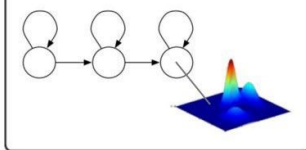
Forward Algorithm

Likelihood : Given an HMM λ = ( A, B ) and a observation sequence O , determine the likelihood P( O | λ ) .

Baum-Welch Algorithm

Learning : Given an observation sequence O and the set of states in the HMM , Learn HMM parameter A and B
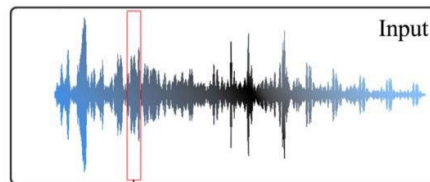
Acoustic model

Lexicon

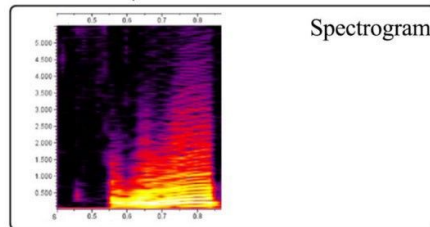smile • $\underset{s}{\bullet}$ $\underset{m}{\bullet}$ $\underset{ay}{\bullet}$ $\underset{l}{\bullet}$

Input

Spectrogram

39 features     Features extraction

Decoding search

Language model

she — 0.5 — is

0.2     0.3

$w_3$     $w_4$     $w_5$

$p(w_t \mid w_{t-1})$

Word sequence

$$\mathbf{W}^* = \arg\max_{\mathbf{W}} P(\mathbf{W} \mid \mathbf{X})$$

$$\mathbf{W}^* = \arg\max_{\mathbf{W}} \; p(\mathbf{X}|\mathbf{W}) \; P(\mathbf{W})$$
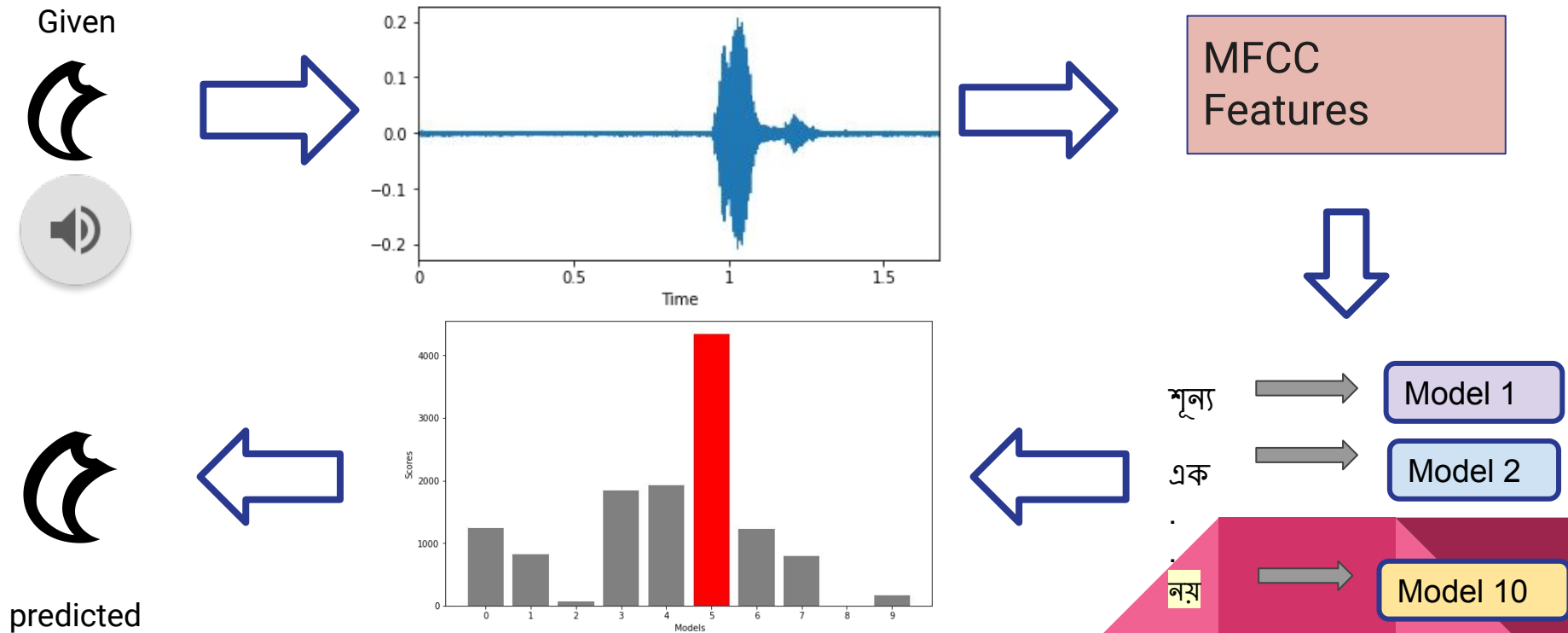
word sequence     acoustic model   language model

Viterbi Algorithm:
 Given an observation sequence O and an
 HMM  λ = ( A, B ) , discover the best hidden sequence Q

# Example :

Given



MFCC Features

শূন্য → Model 1

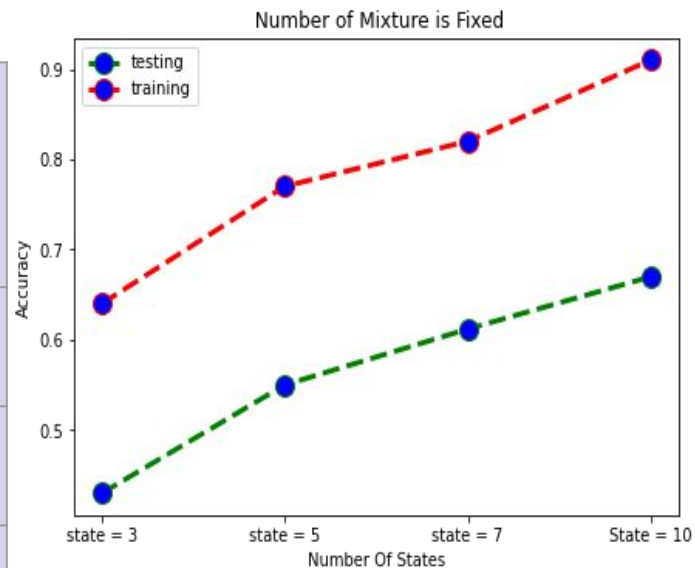এক → Model 2

.
.

নয় → Model 10

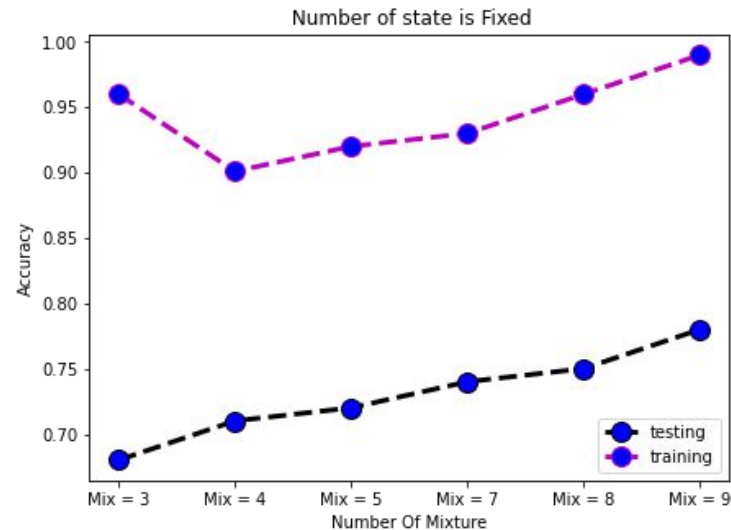predicted

# Result Analysis

Total Data = 600
Train Data = 478
Test Data = 122

| Number Of States | Number Of MIxture | Number Of Iteration | Training Accuracy | Testing Accuracy |
|---|---|---|---|---|
| 3 | 2 | 100 | 64% | 43% |
| 5 | 2 | 100 | 77% | 55% |
| 7 | 2 | 100 | 82% | 61% |
| 10 | 2 | 100 | 91% | 67% |



Number of Mixture is Fixed

# Changing Parameters :

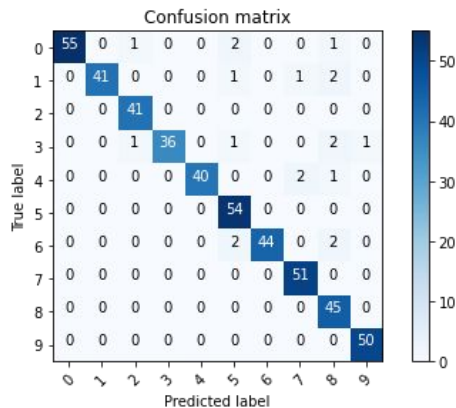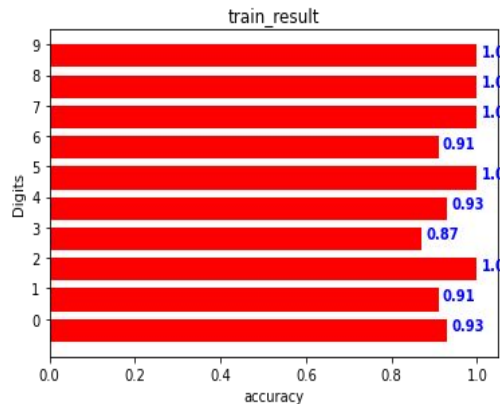| Number Of States | Number Of MIxture | Number Of Iteration | Training Accuracy | Testing Accuracy |
|---|---|---|---|---|
| 10 | 3 | 100 | 96% | 68% |
| 10 | 4 | 100 | 90% | 71% |
| 10 | 8 | 100 | 96% | 74% |
| 10 | 9 | 400 | 96% | 76% |



Number of state is Fixed

Data = 478

Training Analysis

Accuracy = 96%

Confusion Matrix for training set

Digit wise accuracy for training set

Data = 122
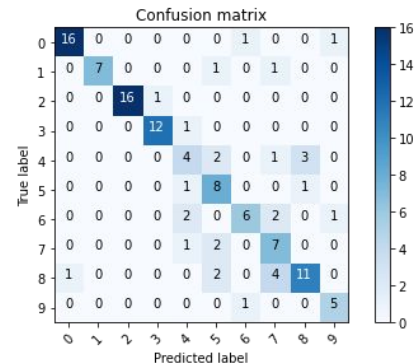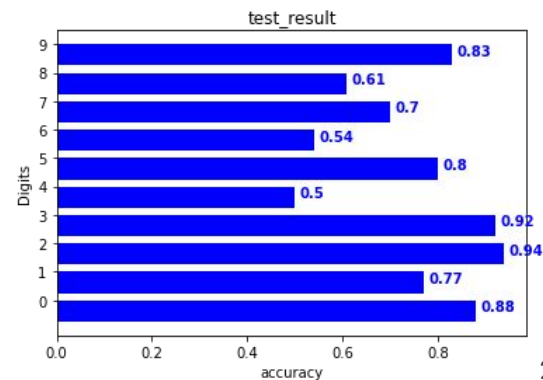
Testing Analysis

Accuracy = 76%

Confusion Matrix for test set

Digit wise accuracy for test set

24

Any Questions **?**

Thank You