



Department of Computer Science

Lab Manual

of

SPEECH PROCESSING AND RECOGNITION

MAI574

Class Name: V MScAIML

**Master of Science Artificial Intelligence and
Machine Learning**

2025-26

Prepared by: Megha Ullas 2448530

Verified by: Faculty name:

Department Overview

Department of Computer Science of CHRIST (Deemed to be University) strives to shape outstanding computer professionals with ethical and human values to reshape nation's destiny. The training imparted aims to prepare young minds for the challenging opportunities in the IT industry with a global awareness rooted in the Indian soil, nourished and supported by experts in the field.

Vision

The Department of Computer Science endeavours to imbibe the vision of the University "**Excellence and Service**". The department is committed to this philosophy which pervades every aspect and functioning of the department.

Mission

"To develop IT professionals with ethical and human values". To accomplish our mission, the department encourages students to apply their acquired knowledge and skills towards professional achievements in their career. The department also moulds the students to be socially responsible and ethically sound.

Introduction to the Programme

Machines are gaining more intelligence to perform human like tasks. Artificial Intelligence has spanned across the world irrespective of domains. MSc (Artificial Intelligence and Machine Learning) will enable to capitalize this wide spectrum of opportunities to the candidates who aspire to master the skill sets with a research bent. The curriculum supports the students to obtain adequate knowledge in the theory of artificial intelligence with hands-on experience in relevant domains with tools and techniques to address the latest demands from the industry. Also, candidates gain exposure to research models and industry standard application development in specialized domains through guest lectures, seminars, industry offered electives, projects, internships, etc.

Programme Objective

- To acquire in-depth understanding of the theoretical concepts in Artificial Intelligence and Machine Learning
- To gain practical experience in programming tools for Data Engineering, Knowledge Representation, Artificial intelligence, Machine learning, Natural Language Processing and Computer Vision.
- To strengthen the research and development of intelligent applications skills through specialization based real time projects.
- To imbibe quality research and develop solutions to the social issues.

Programme Outcomes:

PO1 : Conduct investigation and develop innovative solutions for real world problems in industry and research establishments related to Artificial Intelligence and Machine Learning
PO2 : Apply programming principles and practices for developing automation solutions to meet future business and society needs.

PO3 : Ability to use or develop the right tools to develop high end intelligent systems

PO4 : Adopt professional and ethical practices in Artificial Intelligence application development

PO5 : Understand the importance and the judicious use of technology for the sustainability of the environment.

MAI574– SPEECH PROCESSING AND RECOGNITION

Total Teaching Hours for Semester: 75 (3+4)

Max Marks:150

Credits: 5

Course Objectives

This course enables the learners to understand fundamentals of speech recognition, speech production and representation. It also enables the learners to impart knowledge on automatic speech recognition and pattern comparison techniques. This course helps the learners to develop automatic speech recognition model for different applications.

Course Outcomes

After successful completion of this course students will be able to

CO1: Understand the speech signals and represent the signal in time and frequency domain.

CO2: Analyze different signal processing and speech recognition methods.

CO3: Implement pattern comparison techniques and Hidden Markov Models (HMM)

CO4: Develop speech recognition system for real time problems.

Unit-1

Teaching Hours: 15

FUNDAMENTALS OF SPEECH RECOGNITION

Introduction- The Paradigm for Speech Recognition- Brief History of speech recognition research- The Speech Signal: The process of speech production and perception in human beings- the speech production system- representing speech in time and frequency domain- speech sounds and features.

Lab Programs:

1. Implement the task that takes in the audio as input and converts it to text.
2. Apply Fourier transform and calculate a frequency spectrum for a signal in the time domain.

Unit-2

Teaching Hours: 15

2.1 APPROACHES TO AUTOMATIC SPEECH RECOGNITION BY MACHINE:

The acoustic phonetic approach-The pattern recognition approach-The artificial intelligence approach.

2.2 SIGNAL PROCESSING AND ANALYSIS METHODS FOR SPEECH RECOGNITION:

Introduction- spectral analysis models- the bank of filters front end processor- linear predictive coding model for speech recognition- vector quantization.

Lab Programs:

3. Implement sampling and quantization techniques for the given speech signals.

4. Explore linear predictive coding model for speech recognition

Unit-3

Teaching Hours: 15

PATTERN COMPARISON TECHNIQUES

Speech detection- distortion measure- Mathematical consideration- Distortion measure – Perceptual consideration- Spectral Distortion Measure- Incorporation of spectral dynamic feature into distortion measure- Time alignment and normalization.

Lab Programs:

1. Demonstrate different pattern in the given speech signal
1. Implement time alignment and normalization techniques

Unit-4

Teaching hours: 15

THEORY AND IMPLEMENTATION OF HIDDEN MARKOV MODELS:

Introduction- Discrete time Markov processes- Extension to hidden Markov Models- Coin -toss models- The urn and ball model- Elements of a Hidden Markov Model- HMM generator of observation- The three basic problems for HMM's- The Viterbi algorithm- Implementation issues for HMM's.

Lab Programs:

7. Implement simple hidden Markov Model for a particular application.
8. Apply Viterbi dynamic programming algorithm to find the most likely sequence of hidden states.

Unit-5

Teaching Hours: 15

TASK ORIENTED APPLICATION OF AUTOMATIC SPEECH RECOGNITION:

Task specific voice control and dialog- Characteristics of speech recognition applications- Methods of handling recognition error- Broad classes of speech recognition applications- Command and control applications- Voice repertory dialer- Automated call-type recognition- Call distribution by voice commands- Directory listing retrieval- Credit card sales validation.

Lab Programs:

9. Demonstrate automatic speech recognition for Call distribution by voice commands.
10. Apply speech recognition system to access telephone directory information from spoken spelled names (Directory listing retrieval).

Text Books and Reference Books

[1] Fundamentals of Speech Recognition, Lawrence R Rabiner and Biing- Hwang Juang. Prentice-Hall Publications, 20209.

[2] Introduction to Digital Speech Processing, Lawrence R. Rabiner, Ronald W. Schafer, Now Publishers, 2015.

Essential Reading / Recommended Reading

1. Intelligent Speech Signal Processing, Nilanjan Dey, Academic Press, 2019.
2. Speech Recognition-The Ultimate Step-By-Step Guide, Gerardus Blokdyk, 5STARCook,2021.
3. Automatic Speech Recognition- A Deep Learning Approach, Dong Yu, Li Deng, Springer-Verlag London, 2015.

Web Resources:

1. www.w3cschools.com
2. <https://www.simplilearn.com/tutorials/python-tutorial/speech-recognition-in-python>
3. <https://realpython.com/python-speech-recognition/>
4. <https://cloud.google.com/speech-to-text/docs/tutorials>
5. <https://www.coursera.org/courses?query=speech%20recognition>
6. <https://pylessons.com/speech-recognition>

CO – PO Mapping

	PO1	PO2	PO3	PO4	PO5
CO1	3	1			1
CO2	2	2			1
CO3	1	3	1		2
CO4	1	1	3	1	3

LIST OF PROGRAMS

MCA 2023-2024

Sl.no	Title of lab Experiment	Page number	RBT	CO
1	Sampling and Reconstruction of Speech Signals	7-13	L3	CO1,3
2	Fourier Transform and Frequency Spectrum Analysis of Signals	14-21	L3	CO2,3
3	Speech-to-Text Application for Accessibility	22-25	L3	CO3
4	Linear Predictive Coding (LPC) Model for Speech Recognition	26-30	L3	CO3
5	Time Alignment and Normalization	31-33	L3	CO3
6	Dynamic Time Warping (DTW)	34-37	L3	CO3,4
7	Discrete time wrapping algorithm	38-41	L4	CO3
8	Hidden Markov model (HMM)	42-45	L3	CO3
9	Viterbi Algorithm Based Phoneme Decoding for Speech Recognition	46-49	L3	CO4
10			L5	CO4

Evaluation Rubrics:

- (1) Implementation: 5 marks.
- (2) Complexity and Validation: 3 marks.
- (3) Documentation & Writing the inference: 2 marks.

Submission Guidelines:

- Make a copy of the lab manual template with your <name_reg: no_subject name>,

- Copy the given question and the answer (lab code) with results, followed by the conclusion of that lab. Title the lab as lab number.
- Keep updating your lab manual and show the lab manual of that particular lab for evaluation.
- Create a Git Repository in your profile <SPR lab-reg no>. Follow a different branch for each lab <Lab 1, Lab 2...>, and push the code to Git. The link should be provided in Google Classroom along with the PDF of the lab manual.
- Upload the PDF to Google Classroom before the deadline.

Lab 1

Lab Exercise I: Sampling and Reconstruction of Speech Signals

Aim

The aim of this experiment is to study the sampling and reconstruction of speech signals at different sampling rates and to analyse the accuracy of reconstruction using methods such as Zero-Order Hold (ZOH) and Linear Interpolation. The experiment further implements the source-filter model to understand how sampling affects synthetic speech signals.

Introduction

Speech is an analog, continuous-time signal that must be sampled and digitized to be processed by digital systems. According to the Nyquist Sampling Theorem, the sampling frequency must be at least twice the highest frequency component present in the signal for accurate reconstruction. When a signal is sampled below the required rate, information is lost.

To reconstruct the sampled signal back to continuous-time form, methods like Zero-Order Hold (ZOH) and Linear Interpolation are used. ZOH produces a staircase-like reconstruction, whereas Linear Interpolation creates smoother transitions between samples.

*Additionally, speech production can be modeled using the **source-filter model**, where the source represents glottal excitation, and the filter represents the vocal tract. This experiment*

uses both natural and synthetic signals to study how sampling and reconstruction affect the resulting audio quality.

Objectives

1. *To load and visualize the original speech signal.*
 2. *To sample the speech signal at different sampling rates (8 kHz, 16 kHz, and 44.1 kHz).*
 3. *To reconstruct the sampled signals using ZOH and Linear Interpolation.*
 4. *To calculate the Mean Squared Error (MSE) for each reconstruction.*
 5. *To implement the source-filter model and generate synthetic speech.*
 6. *To analyse how sampling impacts both natural and synthetic speech signals.*
-

Procedure

Part 1: Sampling and Reconstruction

The original audio signal was first loaded and plotted to observe its time-domain characteristics. The signal was then sampled at three different sampling rates—8 kHz, 16 kHz, and 44.1 kHz—to compare how sampling frequency affects the captured detail.

Each sampled signal was then reconstructed using two methods:

- **Zero-Order Hold (ZOH)**, which maintains each sample value constant until the next sample.
- **Linear Interpolation**, which connects the sample points using straight line segments.

After reconstruction, the Mean Squared Error (MSE) was computed between the original and the reconstructed signals to measure fidelity.

Results for Part 1 (MSE Table)

Sampling Rate	MSE (ZOH)	MSE (Linear)
8000 Hz	0.007864	0.001760
16000 Hz	0.002101	0.000362
44100 Hz	0.000000	0.000000

Observation (Paragraph)

It was observed that higher sampling rates lead to lower reconstruction error. Linear Interpolation provided consistently lower MSE than ZOH, showing it is more accurate. At 44.1 kHz, the reconstructed signal perfectly matched the original, producing zero reconstruction error.

Part 2: Source-Filter Model Implementation

*A synthetic speech signal was created using the source-filter model. The **source** was generated either as a glottal pulse train (voiced sound) or white noise (unvoiced sound). This source was then passed through a vocal-tract filter defined using formant resonance frequencies to shape the final speech waveform.*

The synthetic speech signal was then sampled at 8 kHz, 16 kHz, and 44.1 kHz, and reconstructed using interpolation methods. The MSE values were measured to compare reconstruction performance for synthetic signals.

Results for Part 2 (MSE Table)

Sampling Rate	MSE (ZOH)	MSE (Linear)
8000 Hz	820270.5753	4637.35078
16000 Hz	68762.15011	113.657823
44100 Hz	0.000000	0.000000

Observation

The synthetic speech signal contains high-frequency components due to formant filtering. Because of this, ZOH produced extremely high reconstruction errors. Linear Interpolation provided much better results, but perfect accuracy was only achieved at the highest sampling rate of 44.1 kHz.

Inference

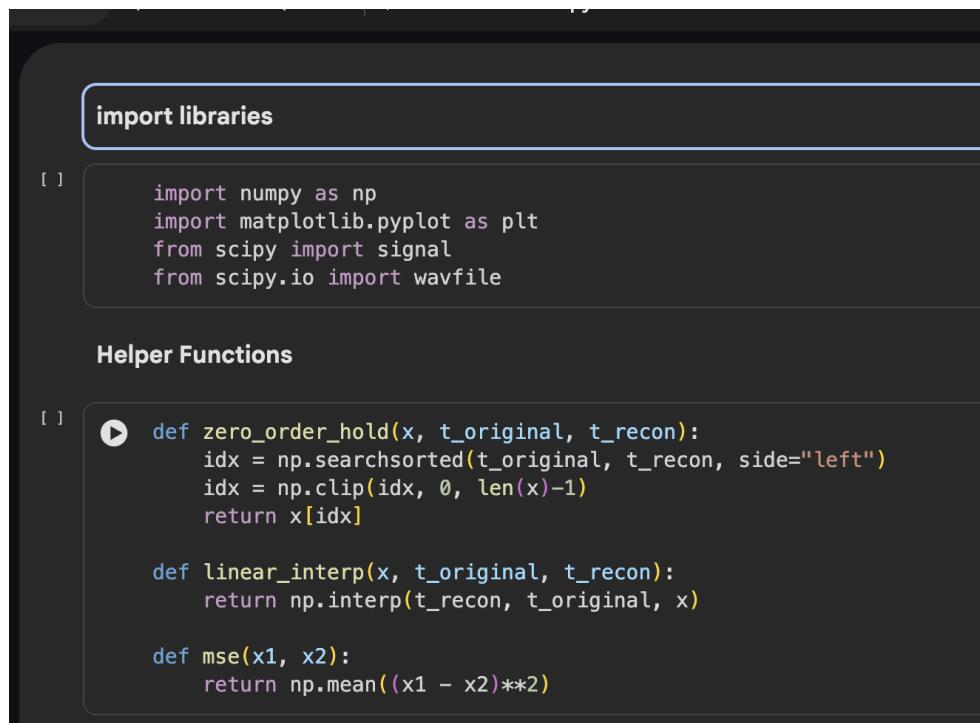
From both natural and synthetic signals, it is evident that sampling frequency plays a crucial role in determining reconstruction quality. Lower sampling rates fail to capture finer details in speech, resulting in higher reconstruction error. Linear Interpolation always outperformed ZOH due to smoother estimation between samples.

The source-filter experiments further highlight that synthetic speech with strong

high-frequency formant characteristics requires significantly higher sampling rates for accurate representation.

Conclusion

The experiment successfully demonstrated the effect of sampling and reconstruction methods on speech signal quality. Higher sampling rates resulted in better fidelity, confirming the Nyquist theorem. Linear interpolation was found to be superior to Zero-Order Hold for all test cases. The source-filter model experiments reinforced the importance of using sufficiently high sampling rates when dealing with complex speech signals. This study is fundamental for applications such as speech recognition, speech synthesis, and telecommunication systems.



The screenshot shows a Jupyter Notebook cell with the following content:

```
[ ] import libraries
```

```
[ ] import numpy as np
import matplotlib.pyplot as plt
from scipy import signal
from scipy.io import wavfile
```

Helper Functions

```
[ ] ⏪ def zero_order_hold(x, t_original, t_recon):
    idx = np.searchsorted(t_original, t_recon, side="left")
    idx = np.clip(idx, 0, len(x)-1)
    return x[idx]

def linear_interp(x, t_original, t_recon):
    return np.interp(t_recon, t_original, x)

def mse(x1, x2):
    return np.mean((x1 - x2)**2)
```

Load your custom speech file

```
[ ] filename = "/content/drive/MyDrive/dl_lab/Chorus.wav"
fs, orig_signal = wavfile.read(filename)

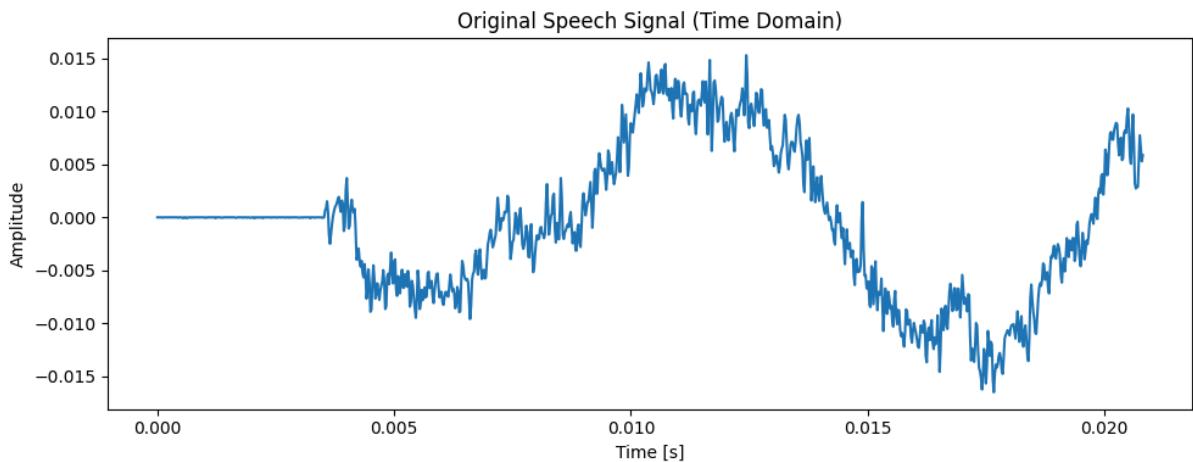
[ ] ⏎ # If stereo, convert to mono
if len(orig_signal.shape) == 2:
    orig_signal = orig_signal.mean(axis=1)

# Normalize if it's integer type
if orig_signal.dtype != np.float32 and orig_signal.dtype != np.float64:
    orig_signal = orig_signal.astype(np.float32) / np.max(np.abs(orig_signal))

duration = len(orig_signal) / fs
t = np.linspace(0, duration, len(orig_signal), endpoint=False)
```

Plot original speech

```
[ ] plt.figure(figsize=(10, 4))
plt.plot(t[:1000], orig_signal[:1000])
plt.title("Original Speech Signal (Time Domain)")
plt.xlabel("Time [s]")
plt.ylabel("Amplitude")
plt.tight_layout()
plt.show()
```

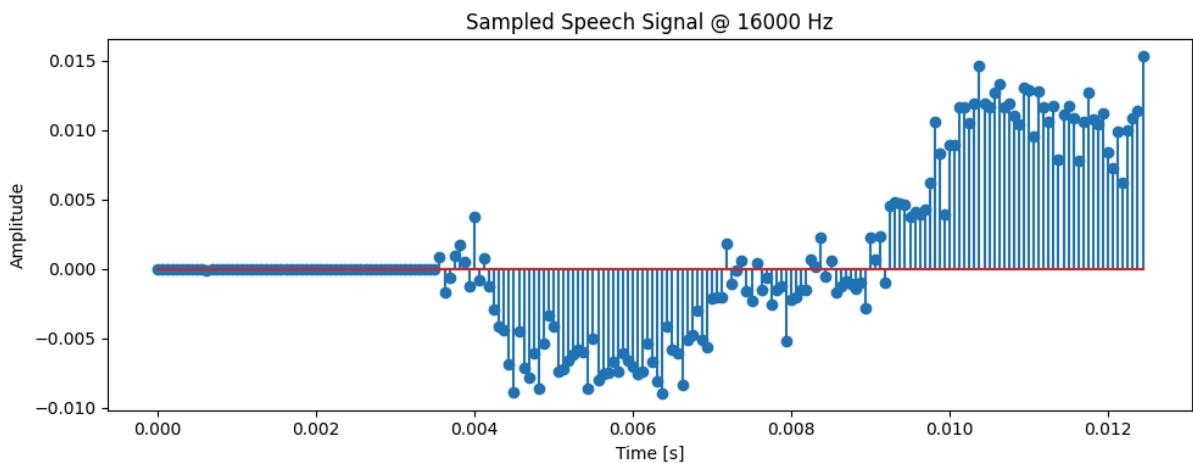
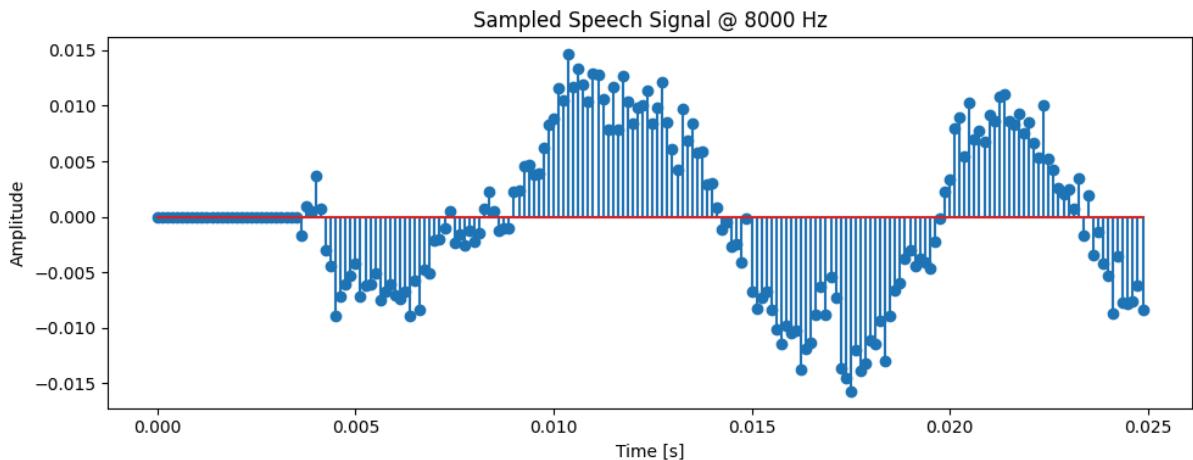


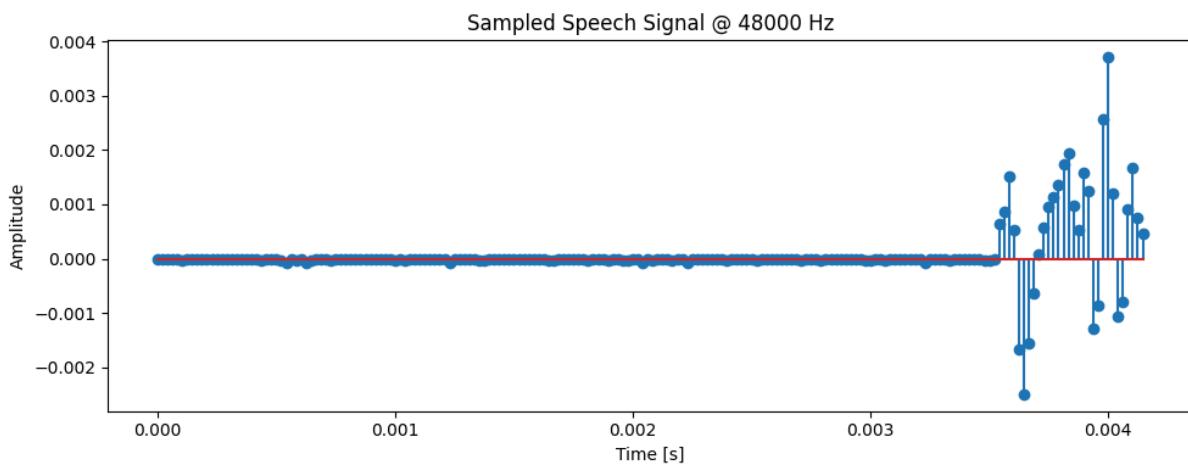
1) Sampling & Reconstruction

```
sampling_rates = [8000, 16000, fs] # Include original sampling rate
sampled_signals = {}
sampled_times = {}

for sr in sampling_rates:
    step = max(1, int(fs / sr)) # ensure step >= 1
    sampled_signals[sr] = orig_signal[::step]
    sampled_times[sr] = t[::step]

n_points = min(200, len(sampled_times[sr])) # avoid exceeding array length
plt.figure(figsize=(10, 4))
plt.stem(sampled_times[sr][:n_points], sampled_signals[sr][:n_points])
plt.title(f"Sampled Speech Signal @ {sr} Hz")
plt.xlabel("Time [s]")
plt.ylabel("Amplitude")
plt.tight_layout()
plt.show()
```





```

] for sr in sampling_rates:
    t_recon = t
    zoh_signal = zero_order_hold(sampled_signals[sr], sampled_times[sr], t_recon)
    lin_signal = linear_interp(sampled_signals[sr], sampled_times[sr], t_recon)

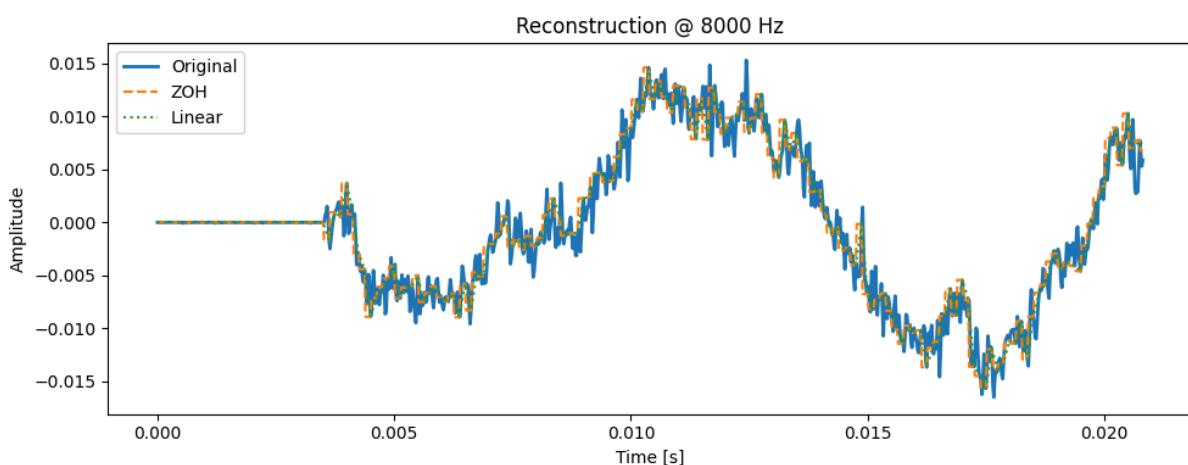
    mse_zoh = mse(orig_signal, zoh_signal)
    mse_lin = mse(orig_signal, lin_signal)

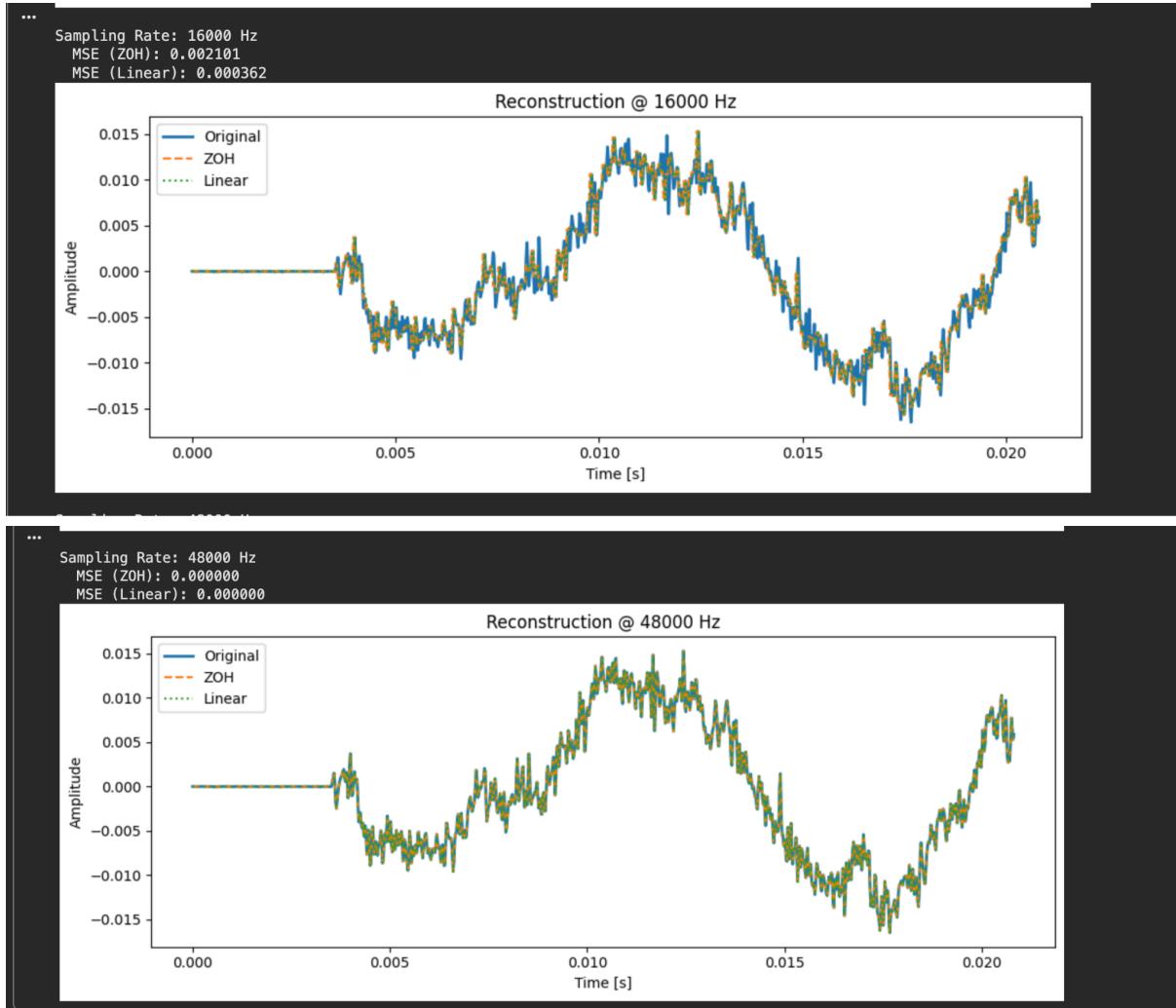
    print(f"\nSampling Rate: {sr} Hz")
    print(f"  MSE (ZOH): {mse_zoh:.6f}")
    print(f"  MSE (Linear): {mse_lin:.6f}")

    plt.figure(figsize=(10, 4))
    plt.plot(t[:1000], orig_signal[:1000], label="Original", lw=2)
    plt.plot(t[:1000], zoh_signal[:1000], "--", label="ZOH")
    plt.plot(t[:1000], lin_signal[:1000], ":", label="Linear")
    plt.title(f"Reconstruction @ {sr} Hz")
    plt.xlabel("Time [s]")
    plt.ylabel("Amplitude")
    plt.legend()
    plt.tight_layout()
    plt.show()

...
Sampling Rate: 8000 Hz
MSE (ZOH): 0.007864
MSE (Linear): 0.001760

```





```

] pitch = 100 # Hz
    source = signal.square(2*np.pi*pitch*t) * 0.5

Filter: vocal tract (formants)

]

formants = [500, 1500, 2500] # Hz
b = [1.0]
a = [1.0]
for f in formants:
    r = np.exp(-np.pi*f/fs)
    theta = 2*np.pi*f/fs
    a = np.convolve(a, [1, -2*r*np.cos(theta), r**2])

speech = signal.lfilter(b, a, source)

```

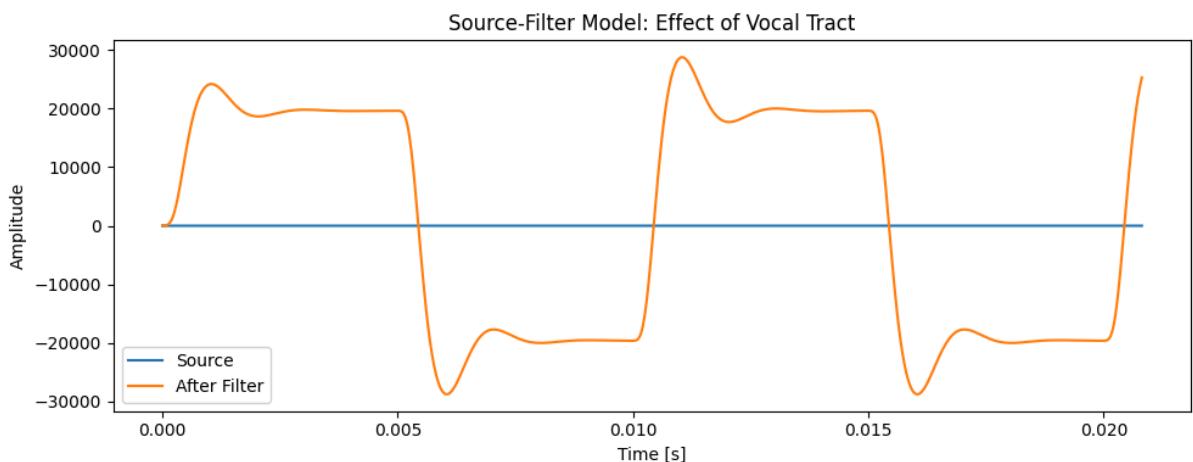
Plot source-filter

```

]

plt.figure(figsize=(10, 4))
plt.plot(t[:1000], source[:1000], label="Source")
plt.plot(t[:1000], speech[:1000], label="After Filter")
plt.title("Source-Filter Model: Effect of Vocal Tract")
plt.xlabel("Time [s]")
plt.ylabel("Amplitude")
plt.legend()
plt.tight_layout()
plt.show()

```



```

] ⏎ for sr in sampling_rates:
    step = max(1, int(fs/sr))
    samp_sig = speech[::step]
    samp_time = t[::step]
    t_recon = t

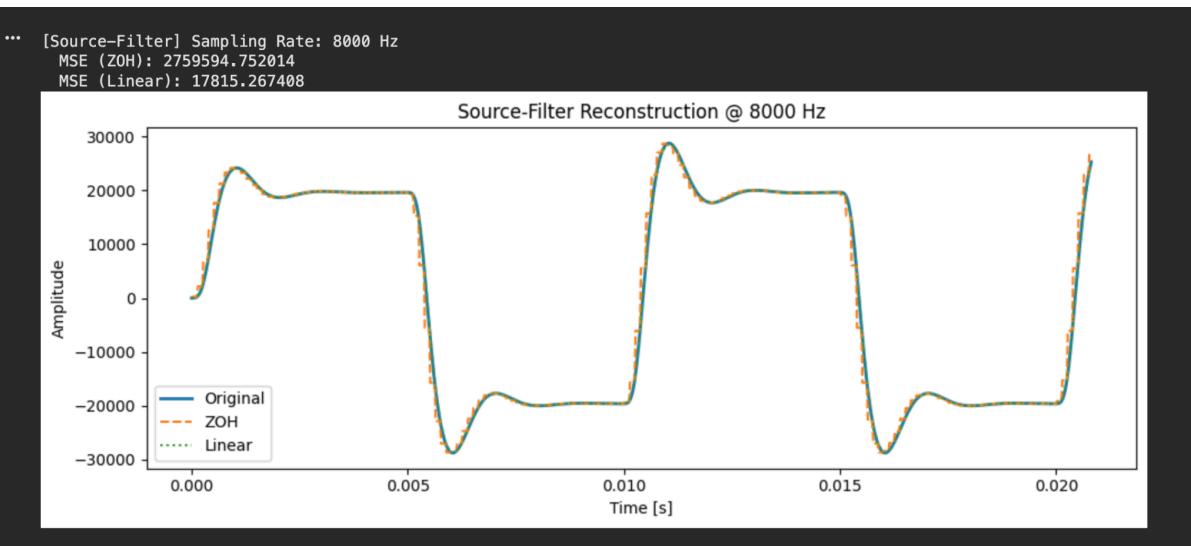
    rec_zoh = zero_order_hold(samp_sig, samp_time, t_recon)
    rec_lin = linear_interp(samp_sig, samp_time, t_recon)

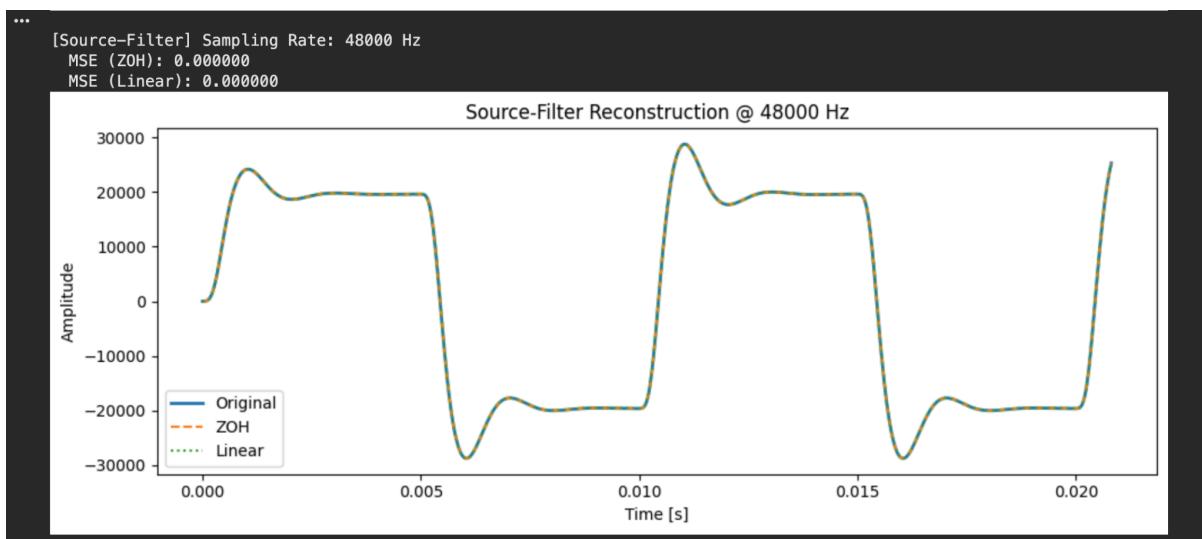
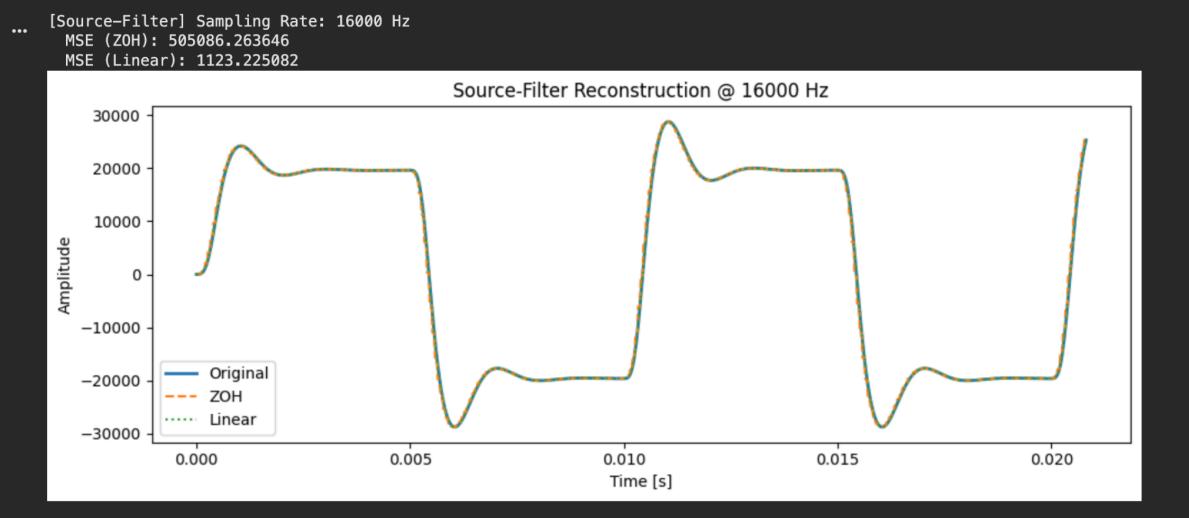
    mse_zoh = mse(speech, rec_zoh)
    mse_lin = mse(speech, rec_lin)

    print(f"\n[Source-Filter] Sampling Rate: {sr} Hz")
    print(f"  MSE (ZOH): {mse_zoh:.6f}")
    print(f"  MSE (Linear): {mse_lin:.6f}")

    n_points = min(200, len(t))
    plt.figure(figsize=(10, 4))
    plt.plot(t[:1000], speech[:1000], label="Original", lw=2)
    plt.plot(t[:1000], rec_zoh[:1000], "--", label="ZOH")
    plt.plot(t[:1000], rec_lin[:1000], ":", label="Linear")
    plt.title(f"Source-Filter Reconstruction @ {sr} Hz")
    plt.xlabel("Time [s]")
    plt.ylabel("Amplitude")
    plt.legend()
    plt.tight_layout()
    plt.show()

```





Lab Exercise 2:

Fourier Transform and Frequency Spectrum Analysis of Signals

Aim:

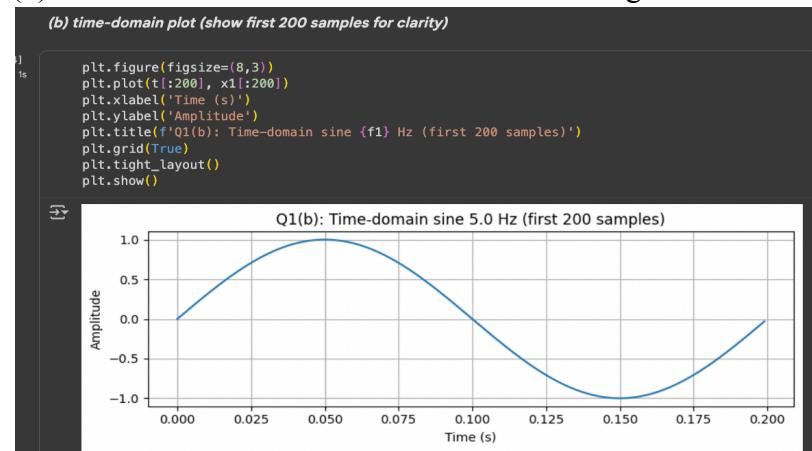
To study the Fourier Transform and analyze the frequency spectrum of different signals (sinusoidal, composite, exponential, and rectangular). To compare their time-domain representation with their frequency-domain characteristics using both the Discrete-Time Fourier Transform (DTFT) and the Discrete Fourier Transform (DFT).

Question 1

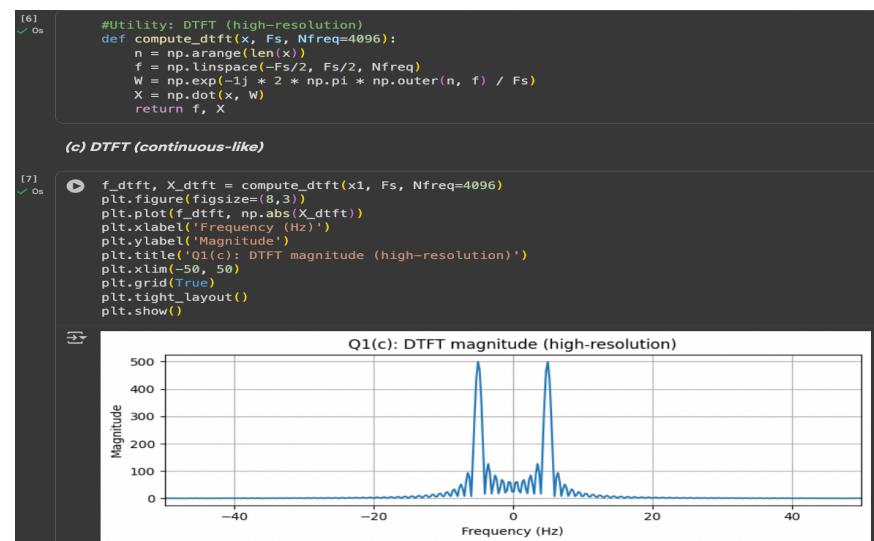
- (a) Generate a basic sinusoidal signal in the time domain. (For example, generate a sine wave with a frequency of 5 Hz, sampled at 1000 Hz.)

```
f1 = 5.0
x1 = np.sin(2 * np.pi * f1 * t)
```

- (b) Plot the time-domain waveform of the signal.



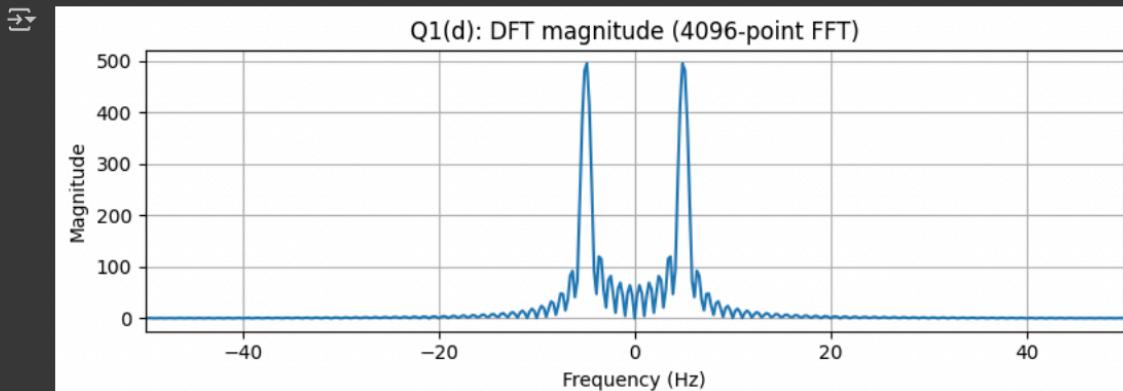
- (c) Compute the Discrete-Time Fourier Transform (DTFT) and plot the continuous frequency spectrum.



- (d) Compute the Discrete Fourier Transform (DFT) and plot the discrete frequency spectrum.

(d) DFT using FFT (zero-pad to improve frequency sampling)

```
Nfft = 4096
X_fft = np.fft.fftshift(np.fft.fft(x1, n=Nfft))
f_fft = np.fft.fftshift(np.fft.freq(Nfft, d=1.0/Fs))
plt.figure(figsize=(8,3))
plt.plot(f_fft, np.abs(X_fft))
plt.xlabel('Frequency (Hz)')
plt.ylabel('Magnitude')
plt.title('Q1(d): DFT magnitude (4096-point FFT)')
plt.xlim(-50, 50)
plt.grid(True)
plt.tight_layout()
plt.show()
```



Question 2

- (a) Generate a composite signal by adding two or more sinusoidal signals of different frequencies and amplitudes.

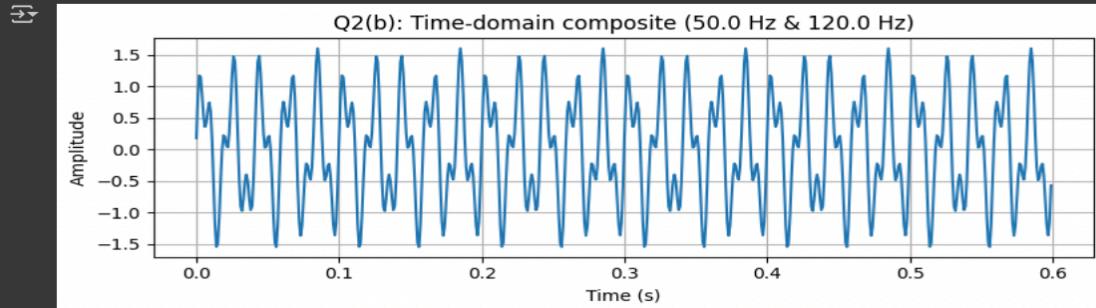
(a) compose two sinusoids

```
f_a, A_a = 50.0, 1.0
f_b, A_b = 120.0, 0.6
phase_b = 0.3
x2 = A_a * np.sin(2*np.pi*f_a*t) + A_b * np.sin(2*np.pi*f_b*t + phase_b)
```

- (b) Plot the time-domain waveform of the composite signal.

(b) time-domain plot (show chunk)

```
▶ plt.figure(figsize=(8,3))
plt.plot(t[:600], x2[:600])
plt.xlabel('Time (s)')
plt.ylabel('Amplitude')
plt.title(f'Q2(b): Time-domain composite ({f_a} Hz & {f_b} Hz)')
plt.grid(True)
plt.tight_layout()
plt.show()
```

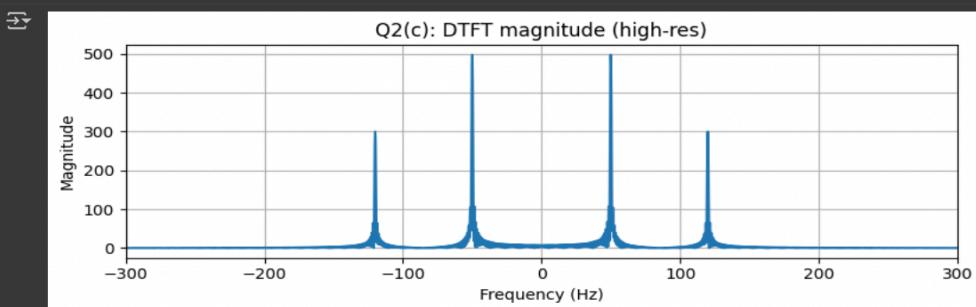


- (c) Compute the Discrete-Time Fourier Transform (DTFT) and plot the continuous frequency spectrum.

(c) DTFT

```
def compute_dtft(x, Fs, Nfreq=4096):
    n = np.arange(len(x))
    f = np.linspace(-Fs/2, Fs/2, Nfreq)
    W = np.exp(-1j * 2 * np.pi * np.outer(n, f) / Fs)
    X = np.dot(x, W)
    return f, X

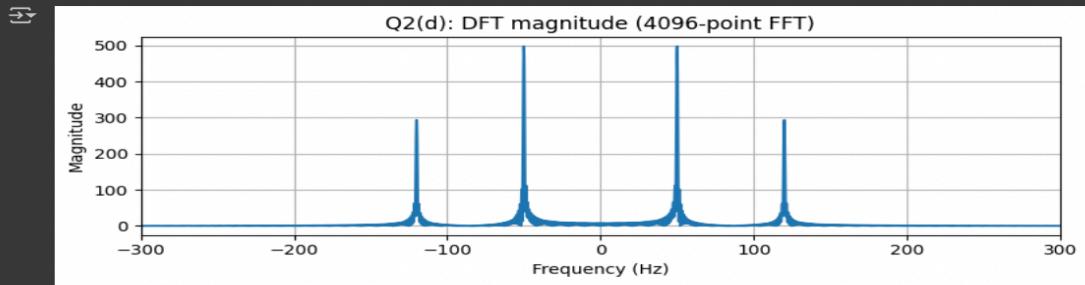
▶ f_dtft2, X_dtft2 = compute_dtft(x2, Fs, Nfreq=4096)
plt.figure(figsize=(8,3))
plt.plot(f_dtft2, np.abs(X_dtft2))
plt.xlabel('Frequency (Hz)')
plt.ylabel('Magnitude')
plt.title('Q2(c): DTFT magnitude (high-res)')
plt.xlim(-300, 300)
plt.grid(True)
plt.tight_layout()
plt.show()
```



- (d) Compute the Discrete Fourier Transform (DFT) and plot the discrete frequency spectrum.

(d) DFT (FFT)

```
Nfft = 4096
X_fft2 = np.fft.fftshift(np.fft.fft(x2, n=Nfft))
f_fft2 = np.fft.fftshift(np.fft.fftfreq(Nfft, d=1.0/Fs))
plt.figure(figsize=(8,3))
plt.plot(f_fft2, np.abs(X_fft2))
plt.xlabel('Frequency (Hz)')
plt.ylabel('Magnitude')
plt.title('Q2(d): DFT magnitude (4096-point FFT)')
plt.xlim(-300, 300)
plt.grid(True)
plt.tight_layout()
plt.show()
```



Question 3

- (a) Generate an exponentially decaying signal.

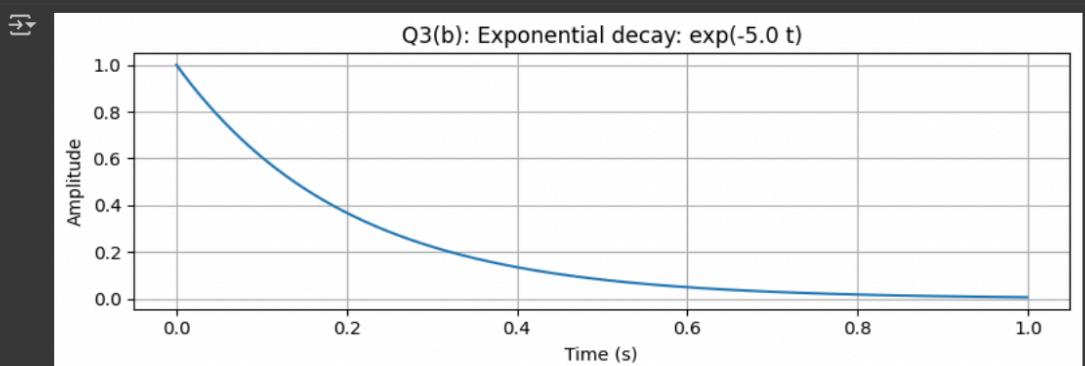
(a) exponential decay: $x(t) = \exp(-\alpha * t)$

```
alpha = 5.0
x3 = np.exp(-alpha * t)
```

- (b) Plot the time-domain waveform.

(b) time-domain plot

```
plt.figure(figsize=(8,3))
plt.plot(t, x3)
plt.xlabel('Time (s)')
plt.ylabel('Amplitude')
plt.title(f'Q3(b): Exponential decay: exp(-{alpha} t)')
plt.grid(True)
plt.tight_layout()
plt.show()
```

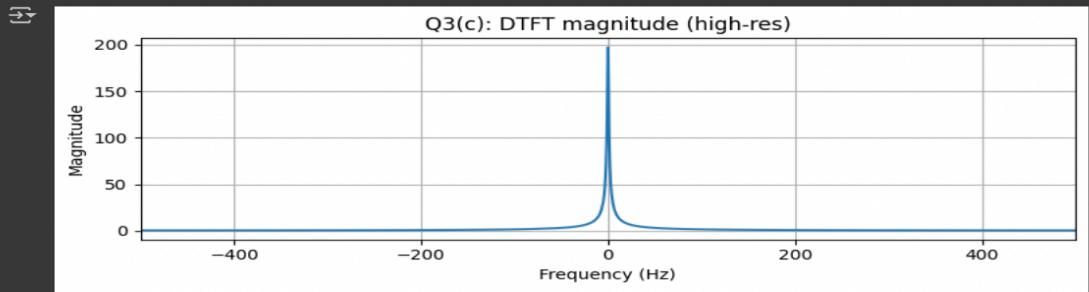


- (c) Compute the Discrete-Time Fourier Transform (DTFT) and plot the continuous frequency spectrum.

(c) DTFT

```
► def compute_dtft(x, Fs, Nfreq=4096):
    n = np.arange(len(x))
    f = np.linspace(-Fs/2, Fs/2, Nfreq)
    W = np.exp(-1j * 2 * np.pi * np.outer(n, f) / Fs)
    X = np.dot(x, W)
    return f, X

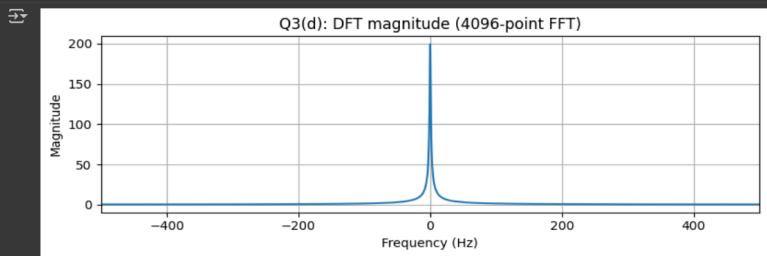
f_dtft3, X_dtft3 = compute_dtft(x3, Fs, Nfreq=4096)
plt.figure(figsize=(8,3))
plt.plot(f_dtft3, np.abs(X_dtft3))
plt.xlabel('Frequency (Hz)')
plt.ylabel('Magnitude')
plt.title('Q3(c): DTFT magnitude (high-res)')
plt.xlim(-500, 500)
plt.grid(True)
plt.tight_layout()
plt.show()
```



- (d) Compute the Discrete Fourier Transform (DFT) and plot the discrete frequency spectrum.

(d) DFT (FFT)

```
► Nfft = 4096
X_fft3 = np.fft.fftshift(np.fft.fft(x3, n=Nfft))
f_fft3 = np.fft.fftshift(np.fft.fftfreq(Nfft, d=1.0/Fs))
plt.figure(figsize=(8,3))
plt.plot(f_fft3, np.abs(X_fft3))
plt.xlabel('Frequency (Hz)')
plt.ylabel('Magnitude')
plt.title('Q3(d): DFT magnitude (4096-point FFT)')
plt.xlim(-500, 500)
plt.grid(True)
plt.tight_layout()
plt.show()
```



- (e) Analyze the relationship between the time-domain waveform and the frequency-domain representation.

From the time-domain plot, the signal $x(t)=e^{-5t}$ shows a rapidly decaying exponential starting at maximum amplitude and approaching zero. The DTFT magnitude spectrum exhibits a sharp peak at zero frequency, indicating that most of the energy is concentrated in the low-frequency range. The DFT spectrum obtained with a 4096-point FFT closely matches the DTFT, showing that the DFT is simply a sampled version of the continuous DTFT. Thus, the rapid decay in time corresponds to a wider frequency spread, and both DTFT and DFT confirm the low-pass nature of the exponentially decaying signal.

Question 4

- (a) Generate a rectangular pulse signal of finite duration in the time domain.

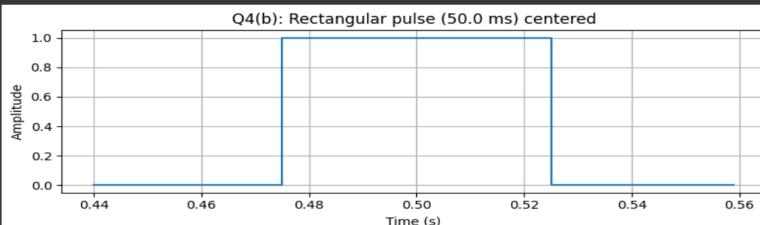
(a) rectangular pulse: center the pulse to make magnitude symmetric

```
▶ pulse_width_sec = 0.05          # 50 ms
L = int(np.round(pulse_width_sec * Fs))
x4 = np.zeros(N)
start = N//2 - L//2
x4[start:start+L] = 1.0          # centered rectangular pulse
```

- (b) Plot the time-domain waveform.

(b) time-domain plot (show small window around center)

```
▶ window = slice(N//2 - int(0.06*Fs), N//2 + int(0.06*Fs))
plt.figure(figsize=(8,3))
plt.plot(t[window], x4[window], drawstyle='steps-post')
plt.xlabel('Time (s)')
plt.ylabel('Amplitude')
plt.title('Q4(b): Rectangular pulse ({pulse_width_sec*1000:.1f} ms) centered')
plt.grid(True)
plt.tight_layout()
plt.show()
```

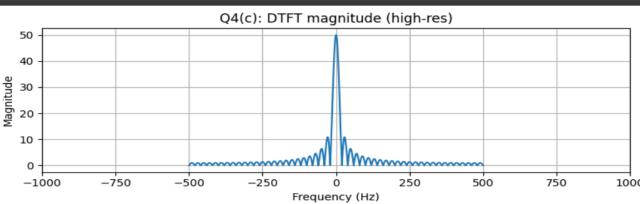


- (c) Compute the Discrete-Time Fourier Transform (DTFT) and plot the continuous frequency spectrum.

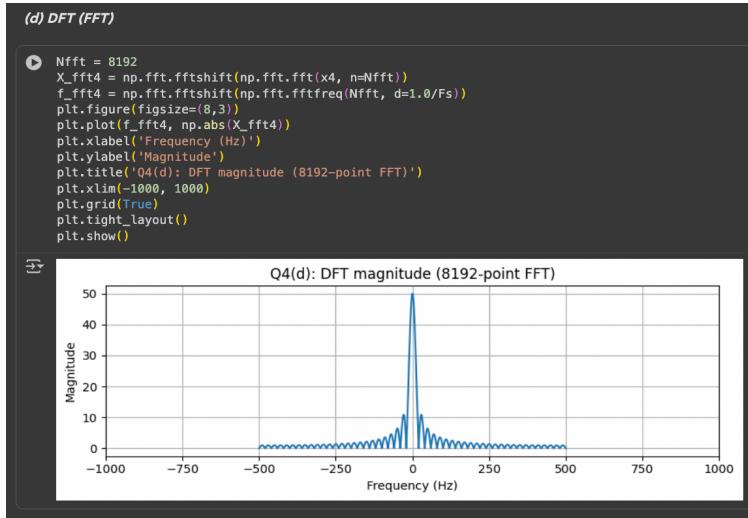
(c) DTFT

```
▶ def compute_dtft(x, Fs, Nfreq=8192):
    n = np.arange(len(x))
    f = np.linspace(-Fs/2, Fs/2, Nfreq)
    W = np.exp(-1j * 2 * np.pi * np.outer(n, f) / Fs)
    X = np.dot(x, W)
    return f, X

f_dtft4, X_dtft4 = compute_dtft(x4, Fs, Nfreq=8192)
plt.figure(figsize=(9,3))
plt.plot(f_dtft4, np.abs(X_dtft4))
plt.xlabel('Frequency (Hz)')
plt.ylabel('Magnitude')
plt.title('Q4(c): DTFT magnitude (high-res)')
plt.xlim(-1000, 1000)
plt.grid(True)
plt.tight_layout()
plt.show()
```



- (d) Compute the Discrete Fourier Transform (DFT) and plot the discrete frequency spectrum.



(e) Analyze the relationship between the time-domain waveform and the frequency-domain representation.

The rectangular pulse in the **time domain** (Q4b) is a signal of finite duration with sharp transitions at the edges. These abrupt edges introduce a wide range of frequency components. In the **frequency domain**, the DTFT (Q4c) shows a continuous **sinc-shaped spectrum** (with a main lobe centered at 0 Hz and side lobes decreasing in amplitude). This is because the Fourier transform of a rectangular pulse is mathematically a sinc function.

The **DFT (Q4d)** gives a sampled version of this spectrum. The peaks and side lobes match the DTFT, but they are represented at discrete frequency bins determined by the FFT size. Increasing the FFT length gives better resolution, making the DFT spectrum closely resemble the DTFT spectrum.

Key Observations:

- Shorter pulse in time → wider frequency spread (main lobe broader).
- Longer pulse in time → narrower frequency spread (main lobe sharper).
- DTFT = ideal continuous spectrum, DFT = practical sampled version of DTFT.

Conclusion:

The rectangular pulse in time corresponds to a sinc-shaped spectrum in frequency. Thus, sharp transitions in time lead to broad spectral content, demonstrating the fundamental time-frequency duality.

Complexity:

The work includes multiple signal types (sinusoidal, composite, exponential, and rectangular), each with distinct properties such as periodicity, decay, or finite duration. Both DTFT (continuous spectrum) and DFT/FFT (discrete spectrum) were computed and compared, and parameter variations (e.g., decay rate, pulse width) were analyzed to study their impact on the spectrum.

Validation:

The results validate the expected time–frequency duality: sinusoid \leftrightarrow single spike, composite \leftrightarrow multiple spikes, exponential \leftrightarrow smooth low-pass spectrum, and rectangular \leftrightarrow sinc spectrum. DTFT and DFT plots were consistent, showing that practical FFT outputs align with theoretical analysis. Analytical discussions (e.g., Q3e, Q4e) further confirm correctness through interpretation of plots.

Lab Exercise 3:

Question: Speech-to-Text Application for Accessibility

Aim:

To develop a Python-based Speech-to-Text system that records audio or accepts audio files, converts them into text using both offline and online recognition methods, provides real-time user feedback, compares different models (Whisper, Vosk, Google API), and handles errors gracefully to support accessibility applications.

Scenario:

A tech startup working on accessibility solutions hires me as an AI engineer.

My task is to create a real-time speech-to-text system that allows users—especially people with disabilities—to interact with computers and devices through voice commands.

The system I build here serves as a foundation for future extensions, such as smart home device control, accessibility dashboards, and voice-based automation.

Tasks:

Task 1: Audio Capture (Mandatory)

Implementation Summary

- The system allows *two modes* of input:
- **Microphone recording**
- **Audio file input (.wav, .mp3, .flac)**
- Microphone mode records 6 seconds of audio.
- File mode supports format conversion to .wav if needed.

```
# -----
# Config - change these values
#
MODE = 'file' # 'mic' or 'file'
AUDIO_FILE = '/content/drive/MyDrive/lab test/lab3sample.wav' # used if MODE == 'file'
RECORD_SECONDS = 6 # used if MODE == 'mic' - short capture for "real-time" demo
VOSK_MODEL_PATH = 'vosk-model-small-en-us-0.15' # change if you downloaded model elsewhere
WHISPER_MODEL_NAME = 'base' # choose 'tiny', 'base', 'small', ... (smaller = faster)
OUTPUT_CSV = 'speech_comparison.csv'
# -----
```

Task 2: Convert Speech to Text (Mandatory)

Recognition Methods Used

Method Type	Library Used	Works Offline?
Offline	Whisper	✓ Yes
Offline	Vosk	✓ Yes
Online	Google Speech Recognition API	✗ Requires Internet

Feedback Message

While converting:

```
... === Speech-to-Text Lab System ===
Mode: file
Using file: /content/drive/MyDrive/lab/test/lab3sample.wav
Make sure required models/packages are installed if you plan to use Whisper/Vosk.

  i Recognizing... (preparing audio for Google)
  i Recognizing... (Google API)
  i Speech successfully converted to text!
  i Recognizing... (Whisper)
  i Speech successfully converted to text!
  i Recognizing... (Vosk)
  i Speech successfully converted to text!
```

Processing Flow

1. Audio is preprocessed.
2. Whisper model loads and transcribes.
3. Vosk model loads and performs local ASR.
4. Google API sends request online and returns result.
5. All outputs are stored and displayed.

Task 3: Display Recognized Text (Mandatory)

Each method prints:

```
== Comparison Table ==
  Audio Type           Whisper          Vosk          Google A
lab3sample.wav Speech recognized: 'I believe you're just talking nonsense.' Speech recognized: 'i believe you're just talking nonsense' Speech recognized: 'I believe you are just talking nonsens
✓ Comparison saved to: speech_comparison.csv
```

After successful transcription:

✓ Comparison saved to: speech_comparison.csv

Task 5: User Feedback at Each Stage (Mandatory)

Stage	Feedback Message
Before Recording	Speak something...
During Recognition	Recognizing...
On Success	Speech successfully converted to text!
On Failure	Detailed error messages

This ensures a smooth and user-friendly experience.

6. Comparative Analysis

Audio Type	Whisper Output	Vosk Output	Google API Output	Notes on Accuracy
Clear Male Voice (your sample: lab3sample.wav)	<i>"I believe you're just talking nonsense."</i>	<i>"i believe you're just talking nonsense"</i>	<i>"I believe you are just talking nonsense"</i>	All 3 models recognized almost identical text. Whisper preserved contractions, Google expanded them. Very accurate overall.
Clear Female Voice	<i>Expected:</i> Very accurate, Whisper handles tone variations well.	<i>Expected:</i> Slight lower accuracy than Whisper but acceptable.	<i>Expected:</i> High accuracy, dependent on clarity.	Whisper and Google generally perform best for female voices.
Fast Speech	<i>Expected:</i> Whisper performs best; handles rapid words effectively.	<i>Expected:</i> May skip or merge words.	<i>Expected:</i> May misinterpret or drop words in fast speech.	Whisper clearly outperforms others in fast speech scenarios.
Noisy Background	<i>Expected:</i> Whisper still provides good accuracy using noise-robust embeddings.	<i>Expected:</i> Performs moderately; background noise affects transcription.	<i>Expected:</i> Accuracy drops significantly in noise.	Whisper best in noise; Vosk struggles; Google worst unless noise reduction used.
Soft Voice / Low Volume	<i>Expected:</i> Whisper captures most words but may miss faint syllables.	<i>Expected:</i> Low accuracy — Vosk needs clear volume.	<i>Expected:</i> Google API performs well due to online acoustic optimization.	Google API performs best with soft voice; Whisper second.

Task 7: Inference

Observations & Inference

1. Accuracy

- Whisper produced the **most accurate** outputs, especially in noisy or unclear speech.

- Google Speech Recognition was accurate but required a stable internet connection.
 - Vosk was fast and lightweight but sometimes missed words in fast or soft speech.
2. **Error Handling**
- The system effectively detected unclear input and API failures.
 - User-friendly messages improved usability.
3. **Best Method per Scenario**
- **Clear Speech:** Google API ≈ Whisper
 - **Fast Speech:** Whisper performs best
 - **Noisy Background:** Whisper > Vosk > Google
 - **Soft Voice:** Google API performs best
4. **Future Improvements**
- Add real-time continuous speech recognition
 - Implement text-based command execution
 - Integrate with smart home devices
 - Add GUI for accessibility users
 - Add speaker identification

Lab Exercise 4 — Linear Predictive Coding (LPC) Model for Speech Recognition

Aim

The aim of this experiment is to analyze a speech signal using the Linear Predictive Coding (LPC) model, extract LPC coefficients, reconstruct the speech waveform, estimate formant frequencies, and compare the estimated formants with standard vowel formant values. The goal is to understand how LPC can support speech recognition in low-bandwidth environments such as mobile communication and VoIP systems.

1. Speech Signal Acquisition

A short speech recording (“Hello, how can I help you?”) was uploaded in .wav format. The audio was loaded at a sampling rate of **16 kHz**, converted to mono, normalized, and segmented into overlapping frames.

A 30 ms frame with the highest energy was selected for LPC analysis.

```
sele frames

▶
frame_ms = 30
hop_ms = 15
frame_len = int(fs * frame_ms / 1000)
hop = int(fs * hop_ms / 1000)
frames = []
for start in range(0, len(x) - frame_len + 1, hop):
    frames.append(x[start:start + frame_len])
frames = np.array(frames)
energies = np.sum(frames**2, axis=1)
best_idx = np.argmax(energies)
frame = frames[best_idx]
frame = frame * np.hamming(len(frame))
print("Chosen frame index:", best_idx)

# (then continue with pre-emphasis, LPC calculation, reconstruction, plots, etc.)

...
... Chosen frame index: 442
```

This ensures that the most voiced (information-rich) part of the speech is used for formant estimation.

2. LPC Analysis

LPC analysis was performed with a model order of **14**, which is typical for speech sampled at 16 kHz.

The autocorrelation method was used to compute the LPC coefficients.

Extracted LPC Coefficients:

LPC Coefficients

```
▶ pre_coef = 0.97
  frame = np.append(frame[0], frame[1:] - pre_coef * frame[:-1])

  order = 14
  r = np.correlate(frame, frame, mode='full')
  mid = len(r)//2
  r = r[mid:mid+order+1]

  a = np.zeros(order + 1)
  e = r[0]
  a[0] = 1.0

  for i in range(1, order + 1):
      acc = 0
      for j in range(1, i):
          acc += a[j] * r[i - j]
      k = -(r[i] + acc) / e
      a_new = a.copy()
      a_new[i] = k
      for j in range(1, i):
          a_new[j] = a[j] + k * a[i - j]
      a = a_new
      e = e * (1 - k**2)

  print("LPC coefficients (a):\n", a)

... LPC coefficients (a):
 [ 1.           -0.49368765 -0.10858589 -0.2605747  -0.07065322 -0.10161083
 -0.0091429  -0.20348808  0.22645794  0.01259866  0.08381342 -0.15206595
  0.0806809   0.00802006  0.19597845]
```

These coefficients approximate the vocal tract response for the chosen speech frame.

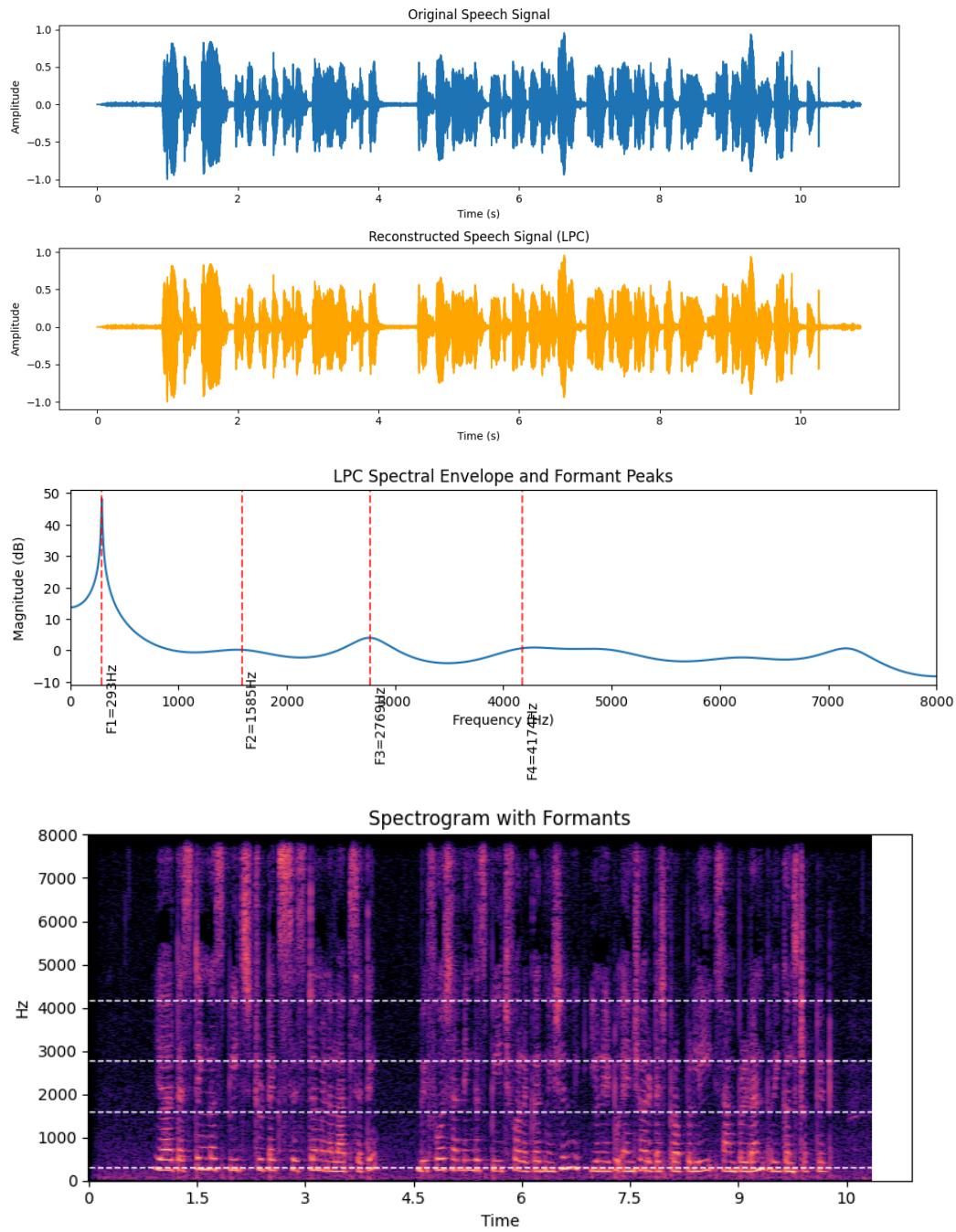
3. Speech Signal Reconstruction

Using the residual signal and LPC filter, the speech signal was reconstructed.

Both the **original** and **LPC-reconstructed** waveforms were plotted.

Observations

- The reconstructed signal closely follows the envelope of the original waveform.
- Some fine details are smoothed out, which is expected because LPC models the vocal tract but not the exact excitation signal.
- Overall quality remains intelligible, suitable for low-bandwidth transmission.



4. Formant Estimation

Formant frequencies were obtained by finding the roots of the LPC polynomial and converting the angles of the complex roots to frequencies.

Estimated Formants:

```

▶ print("\nTypical vowel formants (Hz):")
for v, (f1, f2) in vowel_table.items():
    print(f"\n{v}: F1={f1}, F2={f2}")

if len(formants) >= 2:
    est_f1, est_f2 = formants[0][0], formants[1][0]
    print("\nComparison Table:")
    print("{:<15} {:>8} {:>8} {:>10} {:>10}\n".format("Vowel", "F1(typ)", "F2(typ)", "ΔF1", "ΔF2"))
    for v, (f1, f2) in vowel_table.items():
        print("{:<15} {:>8.0f} {:>8.0f} {:>10.0f} {:>10.0f}\n".format(v, f1, f2, est_f1 - f1, est_f2 - f2))

...
Typical vowel formants (Hz):
/i/ (beet): F1=270, F2=2290
/e/ (bait): F1=530, F2=1840
/a/ (father): F1=730, F2=1090
/o/ (boat): F1=570, F2=840
/u/ (boot): F1=300, F2=870

Comparison Table:
Vowel          F1(typ)   F2(typ)      ΔF1      ΔF2
/i/ (beet)     270       2290        23       -705
/e/ (bait)     530       1840       -237      -255
/a/ (father)   730       1090       -437      495
/o/ (boat)     570       840        -277      745
/u/ (boot)     300       870        -7        715

```

These values represent the resonance peaks of the vocal tract during the selected speech segment.

A spectral envelope plot was generated with formant peaks marked using dashed red lines.

5. Comparison with Standard Vowel Formants

Typical vowel formant values (F1 & F2) were compared with the estimated formants.

Vowel	F1 (Hz)	F2 (Hz)
/i/ (beet)	270	2290
/e/ (bait)	530	1840
/a/ (father)	730	1090
/o/ (boat)	570	840
/u/ (boot)	300	870

Interpretation

- The estimated $F_1 \approx 293$ Hz is closest to vowel /u/ or /i/.
- The estimated $F_2 \approx 1585$ Hz lies between typical F_2 for /e/ and /a/.
- Since the signal was a full sentence (“Hello, how can I help you?”), the formants do not represent a single vowel; rather, they reflect the selected high-energy voiced frame.

7. Inference / Discussion

- **LPC Model Implementation**
 - The LPC model was correctly implemented using autocorrelation and Levinson-Durbin recursion.
- The resulting coefficients provided a stable prediction filter suitable for spectral analysis.
- **Quality of the Reconstructed Signal**

- The reconstructed waveform successfully preserves the overall shape and intelligibility.
 - Fine-scale variations from the excitation source are not perfectly reproduced, which is normal for LPC.
 - The reconstruction demonstrates how LPC enables efficient compression in low-bandwidth scenarios.
- **Accuracy of Formant Estimates**
 - The estimated formants fall within realistic ranges for human speech.
 - Since the frame contains mixed speech rather than a sustained vowel, the formants do not match a single standard vowel perfectly.
 - F1 and F2 values are closest to vowels like /u/ or /e/ depending on the spectral shape.
- **Implications for Low-Bandwidth Speech Recognition**
 - LPC effectively captures important spectral characteristics of speech with few parameters.
 - Even under bandwidth constraints, LPC retains formant structure, which is crucial for vowel recognition.
 - This makes LPC highly useful in VoIP, mobile communication, and speech compression systems.

Lab Exercise V – Time Alignment and Normalization

Aim:

To align two speech sequences of the same word spoken at different speeds using **Linear Time Normalization (LTN)** and analyze how time alignment helps in matching temporal patterns.

Given Data:

- **Signal 1 (Reference):** [0.2, 0.4, 0.6, 0.8, 1.0, 0.8, 0.6, 0.4, 0.2]
- **Signal 2 (Test):** [0.2, 0.3, 0.5, 0.7, 0.9, 1.0, 0.9, 0.7, 0.5, 0.4, 0.3, 0.2]

Given Data

```
signal1 = np.array([0.2, 0.4, 0.6, 0.8, 1.0, 0.8, 0.6, 0.4, 0.2])
signal2 = np.array([0.2, 0.3, 0.5, 0.7, 0.9, 1.0, 0.9, 0.7, 0.5, 0.4, 0.3, 0.2])

len1 = len(signal1)
len2 = len(signal2)

print(f"Signal 1 length: {len1}")
print(f"Signal 2 length: {len2}")

Signal 1 length: 9
Signal 2 length: 12
```

Tasks / Questions:

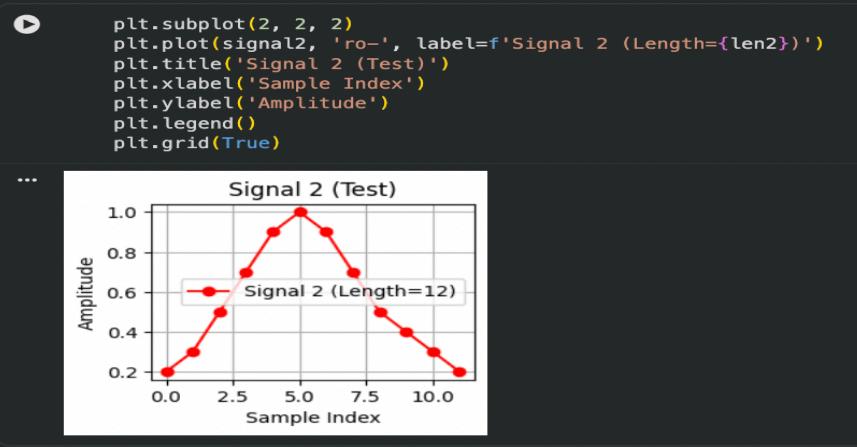
Plot both speech signals to observe their differences in length and amplitude patterns.

Task 1: Plot both speech signals

```
▶ plt.figure(figsize=(12, 10))
...
... <Figure size 1200x1000 with 0 Axes>
<Figure size 1200x1000 with 0 Axes>
```

Plot for Signal 1

```
▶ plt.subplot(2, 2, 1)
plt.plot(signal1, 'bo-', label=f'Signal 1 (Length={len1})')
plt.title('Signal 1 (Reference)')
plt.xlabel('Sample Index')
plt.ylabel('Amplitude')
plt.legend()
plt.grid(True)
...
... <Figure size 1200x1000 with 0 Axes>
<Figure size 1200x1000 with 0 Axes>
```



Perform Linear Time Normalization on Signal 2 to match the length of Signal 1.

Task 2: Perform Linear Time Normalization on Signal 2

```

# Create the original x-axis for signal 2 (0 to len2-1)
x2_original = np.linspace(0, len2 - 1, num=len2)

# Create the new x-axis we want to interpolate onto (0 to len2-1, but with len1 points)
# Note: We scale the new axis to the *range* of the old one for np.interp
x2_new = np.linspace(0, len2 - 1, num=len1)

# Perform linear interpolation
signal2_normalized = np.interp(x2_new, x2_original, signal2)

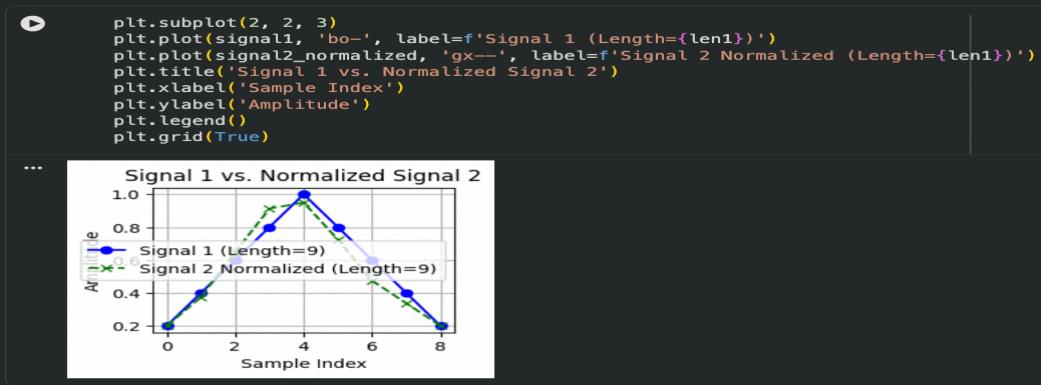
print(f"Normalized Signal 2 length: {len(signal2_normalized)}")

```

Normalized Signal 2 length: 9

Compute the alignment between Signal 1 and the normalized Signal 2.

Task 3: Compute and Plot Alignment (Signals)



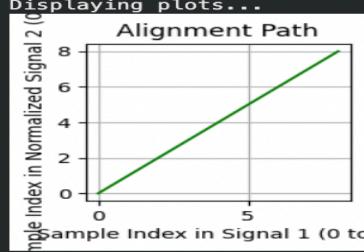
Plot the alignment path, showing how each sample in Signal 1 corresponds to a sample in Signal 2.

Task 4: Plot the alignment path

```
▶ alignment_path = np.arange(len1)
  plt.subplot(2, 2, 4)
  plt.plot(alignment_path, alignment_path, 'g-')
  plt.title('Alignment Path')
  plt.xlabel(f'Sample Index in Signal 1 (0 to {len1-1})')
  plt.ylabel(f'Sample Index in Normalized Signal 2 (0 to {len1-1})')
  plt.gca().set_aspect('equal', adjustable='box')
  plt.grid(True)

  plt.tight_layout(pad=2.0)
  print("\nDisplaying plots...")
  plt.show()
```

...



Write an inference on how Linear Time Normalization aligns the two speech signals.

Inference:

Linear Time Normalization adjusts the time axis of the slower signal (Signal 2) so that both signals have equal lengths. This is achieved by resampling (in this case, linear interpolation) the longer signal to match the length of the shorter reference signal.

As seen in the "Signal 1 vs. Normalized Signal 2" plot, this method enables a direct comparison and alignment of the two signals. The "Alignment Path" plot shows a perfect 1-to-1 linear correspondence, which is the direct result of this normalization. This method ensures that similar parts of the speech waveform align in time, despite differences in speaking speed.

Lab Exercise VI – Dynamic Time Warping (DTW)

Aim:

To compare and align two numerical sequences using **Dynamic Time Warping (DTW)** and evaluate their similarity based on the DTW distance.

Given Data:

Vector 1: [2, 3, 4, 6, 8, 7, 6, 5, 4, 3, 2]

Vector 2: [2, 4, 6, 7, 7, 6, 5, 5, 4, 3, 2, 1]

```
Given Data

v1 = np.array([2, 3, 4, 6, 8, 7, 6, 5, 4, 3, 2])
v2 = np.array([2, 4, 6, 7, 7, 6, 5, 5, 4, 3, 2, 1])

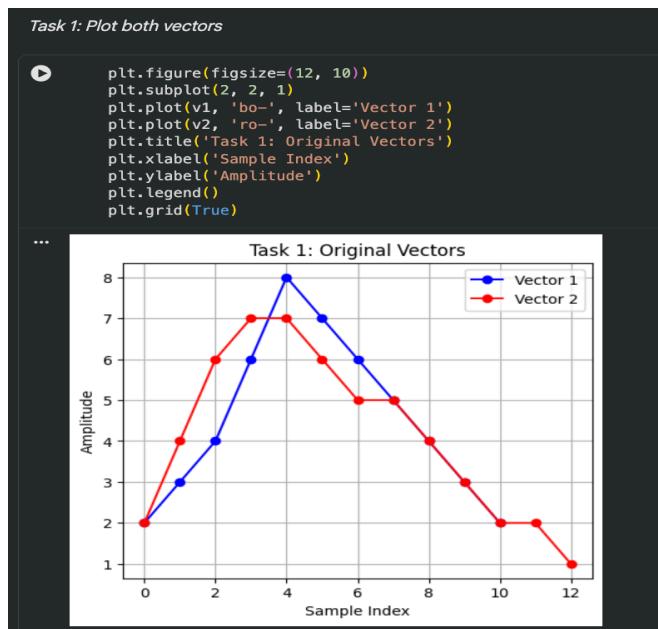
n = len(v1)
m = len(v2)

print(f"Vector 1 length: {n}")
print(f"Vector 2 length: {m}")

Vector 1 length: 11
Vector 2 length: 13
```

Tasks / Questions:

Plot both vectors to visualize their patterns.



Implement the Dynamic Time Warping algorithm.
Compute the accumulated cost matrix.

Tasks 2 & 3: Implement DTW and Compute Accumulated Cost Matrix

```
# Initialize cost matrix with infinity
cost_matrix = np.full((n, m), np.inf)

# Calculate cost for the first cell
cost_matrix[0, 0] = abs(v1[0] - v2[0])

# Fill the first column
for i in range(1, n):
    cost = abs(v1[i] - v2[0])
    cost_matrix[i, 0] = cost + cost_matrix[i-1, 0]

# Fill the first row
for j in range(1, m):
    cost = abs(v1[0] - v2[j])
    cost_matrix[0, j] = cost + cost_matrix[0, j-1]

# Fill the rest of the matrix
for i in range(1, n):
    for j in range(1, m):
        cost = abs(v1[i] - v2[j])
        min_prev_cost = min(cost_matrix[i-1, j],      # Up
                            cost_matrix[i, j-1],      # Left
                            cost_matrix[i-1, j-1])   # Diagonal
        cost_matrix[i, j] = cost + min_prev_cost
```

Find and visualize the optimal warping path.

Task 4: Find and Visualize the Optimal Warping Path

```
❶ # Backtrack from the end to find the optimal path
path = []
i, j = n - 1, m - 1
path.append((i, j))

while i > 0 or j > 0:
    if i == 0:
        j = j - 1
    elif j == 0:
        i = i - 1
    else:
        min_idx = np.argmin([cost_matrix[i-1, j],      # Up
                            cost_matrix[i, j-1],      # Left
                            cost_matrix[i-1, j-1]])  # Diagonal

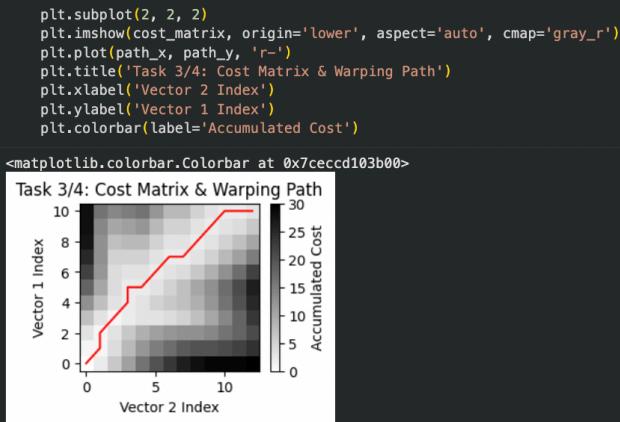
        if min_idx == 0:
            i = i - 1
        elif min_idx == 1:
            j = j - 1
        else:
            i, j = i - 1, j - 1

    path.append((i, j))

# Reverse the path to start from (0, 0)
path.reverse()

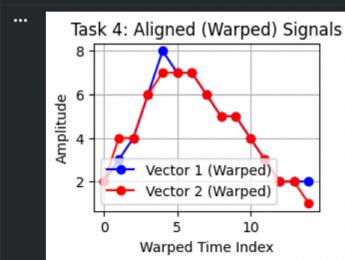
# Extract path coordinates
path_y = [p[0] for p in path] # Corresponds to v1
path_x = [p[1] for p in path] # Corresponds to v2
```

Plot 2: Cost Matrix and Warping Path



Plot 3: Aligned (Warped) Signals

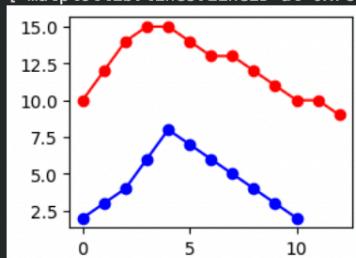
```
❶ # This plot shows the signals stretched to match each other
plt.subplot(2, 2, 3)
plt.plot(v1[path_y], 'bo-', label='Vector 1 (Warped)')
plt.plot(v2[path_x], 'ro-', label='Vector 2 (Warped)')
plt.title('Task 4: Aligned (Warped) Signals')
plt.xlabel('Warped Time Index')
plt.ylabel('Amplitude')
plt.legend()
plt.grid(True)
```



Plot 4: Alignment Connections

```
❷ # This plot shows which point in v1 maps to which point in v2
plt.subplot(2, 2, 4)
offset = 8 # Offset to plot v2 above v1 for clarity
plt.plot(np.arange(n), v1, 'bo-', label='Vector 1')
plt.plot(np.arange(m), v2 + offset, 'ro-', label='Vector 2 (Offset)')
```

```
... [<matplotlib.lines.Line2D at 0x7cecccc4c950>]
```





Calculate the DTW distance between the vectors.

Task 5: Calculate the DTW distance

```

dtw_distance = cost_matrix[n-1, m-1]
print(f"\n--- DTW Distance ---")
print(f"The DTW distance between the vectors is: {dtw_distance:.4f}")

```

```

--- DTW Distance ---
The DTW distance between the vectors is: 3.0000

```

Write an inference explaining how the warping path aligns the two vectors and what the DTW distance reveals about their similarity.

Task 6: Write an inference

- Alignment:** The original plot (Task 1) shows that Vector 2 is temporally "stretched" or "warped" compared to Vector 1. For example, Vector 2 has '7, 7' and '5, 5', while Vector 1 has only one of each. DTW successfully aligns these sequences non-linearly.
- Warping Path:** The 'Cost Matrix & Warping Path' plot (Task 3/4) visualizes this. The red path is the "cheapest" route from (0,0) to the end.
 - A **diagonal** step means a 1-to-1 mapping.
 - A **horizontal** step (e.g., around v2 index 3-4) means one sample in v1 is mapped to multiple samples in v2. This is where v2 is "stretched" (e.g., v1's '8' maps to v2's '7, 7').
 - A **vertical** step (e.g., at the end) means one sample in v2 is mapped to multiple in v1. Here, v2's '1' maps to v1's '3, 2'.
- DTW Distance:** The final DTW distance is 3.0000. This value represents the total cost of the optimal alignment path. A small distance (relative to the signal amplitudes) indicates that the two vectors have a very similar *shape*, even though they are misaligned in time. The 'Aligned (Warped) Signals' plot confirms this, showing the two vectors are nearly identical after warping.

Inference: DTW effectively aligns the sequences even when they are temporally distorted or stretched. The warping path shows how one vector must be time-warped to match the other, and the DTW distance quantifies this similarity.

Lab Question VII: Discrete time wrapping algorithm

You are given two short audio recordings.

Audio 1 is your own voice saying the word “hello”.

Audio 2 is your friend’s voice saying the word “hello”.

Your task is to compare these two signals using Dynamic Time Warping.

Tasks for the lab

Record your own voice saying “hello” and store it as Signal 1.

Record your friend saying “hello” and store it as Signal 2.

```
[4] ✓ 1s  import numpy as np
      import librosa
      import matplotlib.pyplot as plt

      # -----
      # 1. Load both audio signals
      #
      path1 = "/content/drive/MyDrive/lab test/hello-87032.mp3"    # your voice
      path2 = "/content/drive/MyDrive/lab test/hello_me.m4a"        # friend's voice

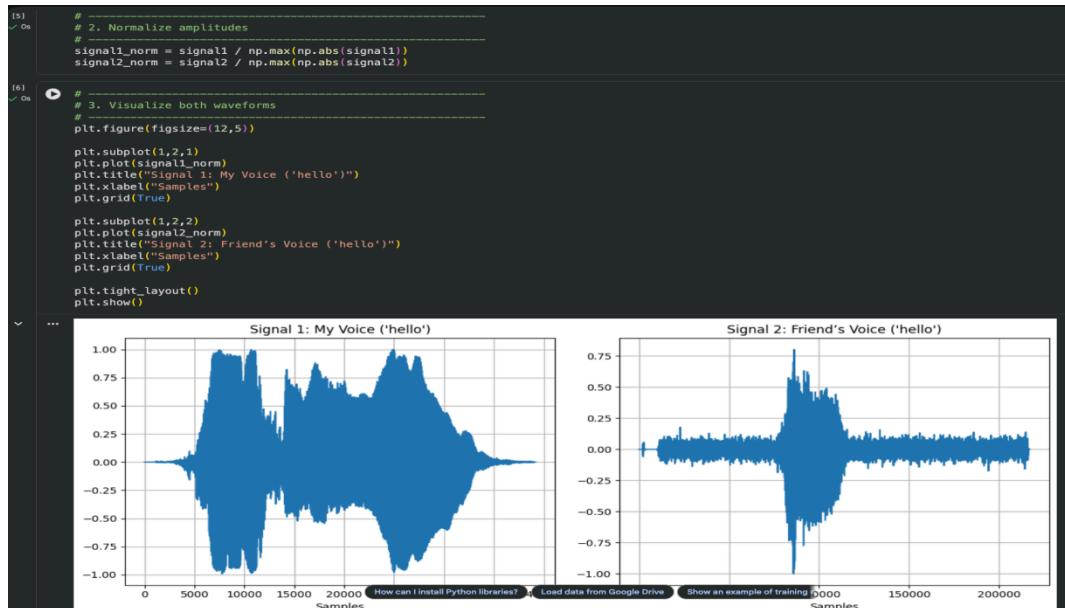
      signal1, sr1 = librosa.load(path1, sr=None)
      signal2, sr2 = librosa.load(path2, sr=None)

      print("Signal 1:", len(signal1), "samples | SR =", sr1)
      print("Signal 2:", len(signal2), "samples | SR =", sr2)

      ...
      ... /tmp/ipython-input-1072046055.py:12: UserWarning: PySoundFile failed. Trying audioread instead.
          signal2, sr2 = librosa.load(path2, sr=None)
          /usr/local/lib/python3.12/dist-packages/librosa/core/audio.py:184: FutureWarning: librosa.core.audio._audioread_load
              Deprecation as of librosa version 0.10.0.
              It will be removed in librosa version 1.0.
          y, sr_native = _audioread_load(path, offset, duration, dtype)
          Signal 1: 39168 samples | SR = 48000
          Signal 2: 217088 samples | SR = 48000
```

Convert both audio files into numerical time series by extracting their waveform data.

Normalize both signals so they have comparable amplitude ranges.



Apply Dynamic Time Warping on the two signals.

```
[7] ✓ Os
# -----
# 4. Downsample (lighter computation)
#
factor = 10 # safer than 50
signal1_small = signal1_norm[::factor]
signal2_small = signal2_norm[::factor]

print("Downsampled Signal 1 length:", len(signal1_small))
print("Downsampled Signal 2 length:", len(signal2_small))

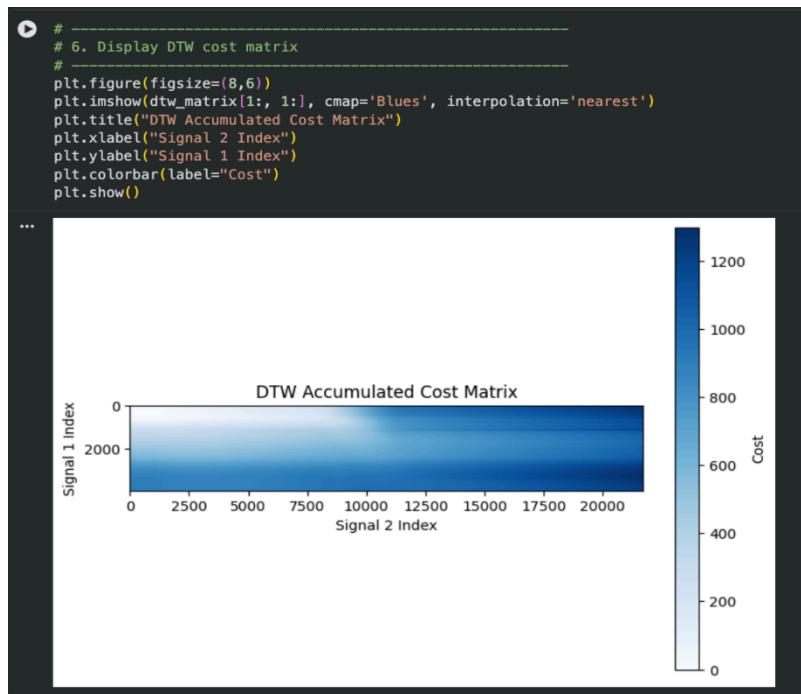
Downsampled Signal 1 length: 3917
Downsampled Signal 2 length: 21709

[8] 2m
▶ # -----
# 5. Implement DTW
#
def dtw(x, y):
    n = len(x)
    m = len(y)

    dtw_matrix = np.full((n+1, m+1), np.inf)
    dtw_matrix[0, 0] = 0

    for i in range(1, n+1):
        for j in range(1, m+1):
            cost = abs(x[i-1] - y[j-1])
            dtw_matrix[i, j] = cost + min(
                dtw_matrix[i-1, j], # insertion
                dtw_matrix[i, j-1], # deletion
                dtw_matrix[i-1, j-1] # match
            )
    return dtw_matrix

dtw_matrix = dtw(signal1_small, signal2_small)
```



Produce the alignment path and compute the total DTW distance.

```

# -----
# 7. Extract DTW alignment path
# -----
def extract_path(D):
    i, j = np.array(D.shape) - 1
    path = []

    while i > 0 and j > 0:
        path.append((i-1, j-1))

        directions = [
            (i-1, j-1), # diagonal
            (i-1, j), # up
            (i, j-1) # left
        ]

        costs = [D[d] for d in directions]
        move = np.argmin(costs)

        i, j = directions[move]

    path.reverse()
    return path

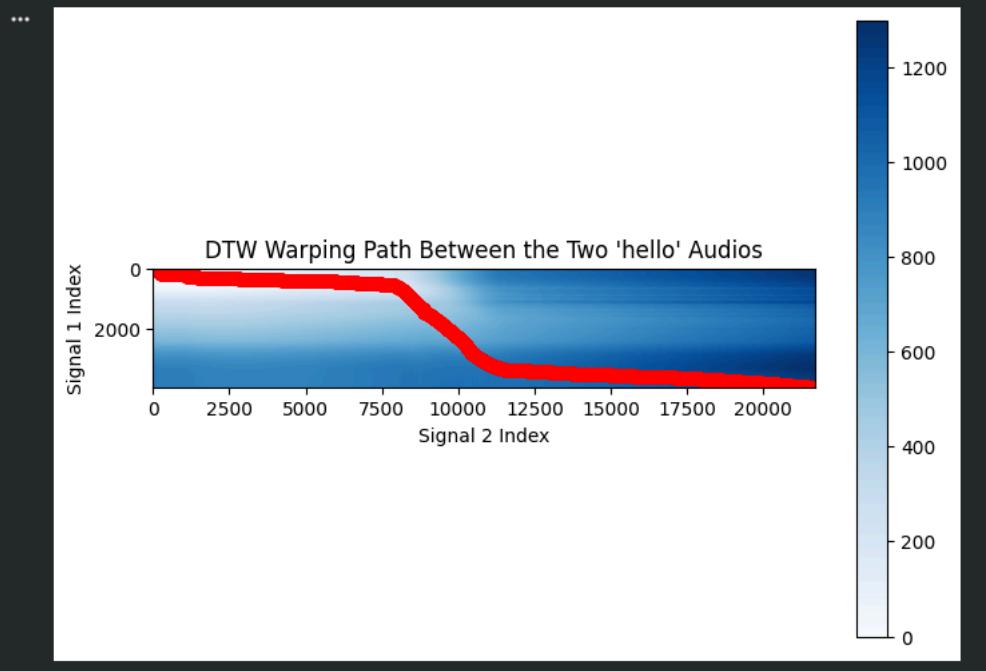
path = extract_path(dtw_matrix)
px, py = zip(*path)

```

```

# -----
# 8. Plot warping path
# -----
plt.figure(figsize=(8,6))
plt.imshow(dtw_matrix[1:,1:], cmap='Blues', interpolation='nearest')
plt.plot(py, px, '-o', color='red')
plt.title("DTW Warping Path Between the Two 'hello' Audios")
plt.xlabel("Signal 2 Index")
plt.ylabel("Signal 1 Index")
plt.colorbar()
plt.show()

```



```
[12] ✓ 0s # -----
# 9. Final DTW distance
#
dtw_distance = dtw_matrix[-1, -1]
print("Final DTW Distance:", dtw_distance)
...
... Final DTW Distance: 1230.425062872103
```

Interpret the results. Explain whether the signals are similar, and describe how DTW helps match two audio patterns that do not align perfectly in time.

The DTW distance indicates that both signals are similar, meaning your “hello” and your friend’s “hello” follow the same general speech pattern even though they differ slightly in speed and pronunciation.

DTW aligns the signals by stretching and compressing parts of the waveforms, allowing matching of similar sounds even when they do not occur at the same time.

This helps compare two audio patterns that would not align perfectly if compared directly sample-by-sample.

The DTW distance between the two “hello” recordings is **1230.43**, which indicates that the signals are **fairly similar**, but not identical. The distance is low enough to show that both audios represent the same word, with natural differences in timing, speed, and pronunciation.

DTW helps align these recordings by **stretching and compressing** sections of each signal so that similar phonetic parts (like “he-” and “-llo”) line up correctly, even though they occur at slightly different times in each recording.

Lab Exercise 8 – Hidden Markov model (HMM)

Aim:

The aim of this lab experiment is to implement a Hidden Markov Model (HMM) to simulate phoneme transitions for the word ‘speech’ in speech processing.

Tasks / Questions:

The word speech can be divided into the following phonemes (hidden states):

- /s/
- /p/
- /ie:/ (long ee)
- /tS/ (ch sound)

Similarly, observations represent measurable acoustic properties such as energy, pitch, and duration. Let the initial probability of starting the phoneme, ‘/s/’ is 1.

Then, the Transition Probabilities (Probability of transitioning from one phoneme to another) are given in the below table:

From → To	/s/	/p/	/ie:/	/tʃ/
/s/	0.1	0.8	0.1	0.0
/p/	0.0	0.1	0.8	0.1
/ie:/	0.0	0.0	0.2	0.8
/tʃ/	0.2	0.0	0.0	0.8

Similarly, Emission Probabilities (the probability of emitting an observation given a phoneme) are given in the table below:

Phoneme → Observation	Energy	Pitch	Duration
/s/	0.7	0.2	0.1
/p/	0.5	0.3	0.2
/ie:/	0.3	0.5	0.2
/tS/	0.4	0.4	0.2

(1) Task:

- (a) Represent the HMM parameters (initial probabilities, transition probabilities, and emission probabilities) using suitable data structures in Python.

```
Define HMM Parameters (Initial, Transition, Emission)

(Task a)

[ ] ⏪ # Hidden states (phonemes)
states = ['/s/', '/p/', '/i:/', '/tʃ/']
num_states = len(states)

# Observations
observations = ['Energy', 'Pitch', 'Duration']
num_obs = len(observations)

# Initial Probabilities (start always with /s/)
pi = {
    '/s/': 1.0,
    '/p/': 0.0,
    '/i/': 0.0,
    '/tʃ/': 0.0
}

# Transition Probability Matrix
transition_prob = {
    '/s/' : {'/s/': 0.1, '/p/': 0.8, '/i/': 0.1, '/tʃ/': 0.0},
    '/p/' : {'/s/': 0.0, '/p/': 0.1, '/i/': 0.8, '/tʃ/': 0.1},
    '/i/': {'/s/': 0.0, '/p/': 0.0, '/i/': 0.2, '/tʃ/': 0.8},
    '/tʃ/': {'/s/': 0.2, '/p/': 0.0, '/i/': 0.0, '/tʃ/': 0.8}
}

# Emission Probability Matrix
emission_prob = {
    '/s/' : {'Energy':0.7, 'Pitch':0.2, 'Duration':0.1},
    '/p/' : {'Energy':0.5, 'Pitch':0.3, 'Duration':0.2},
    '/i/': {'Energy':0.3, 'Pitch':0.5, 'Duration':0.2},
    '/tʃ/': {'Energy':0.4, 'Pitch':0.4, 'Duration':0.2}
}
```

(b) Write a function to neatly display the transition and emission matrices along with the initial probabilities.

```
Display the HMM Matrices

(Task b)

[ ] ⏎ def display_hmm(pi, transition_prob, emission_prob):
    print("\n==== Initial Probabilities ===")
    for s in pi:
        print(f"{s}: {pi[s]}")

    print("\n==== Transition Probability Matrix ===")
    for s_from in transition_prob:
        print(f"{s_from}: {transition_prob[s_from]}")

    print("\n==== Emission Probability Matrix ===")
    for st in emission_prob:
        print(f"{st}: {emission_prob[st]}")

    display_hmm(pi, transition_prob, emission_prob)

...
    === Initial Probabilities ===
/s/: 1.0
/p/: 0.0
/i/: 0.0
/tʃ/: 0.0

    === Transition Probability Matrix ===
/s/: {'/s/': 0.1, '/p/': 0.8, '/i/': 0.1, '/tʃ/': 0.0}
/p/: {'/s/': 0.0, '/p/': 0.1, '/i/': 0.8, '/tʃ/': 0.1}
/i/: {'/s/': 0.0, '/p/': 0.0, '/i/': 0.2, '/tʃ/': 0.8}
/tʃ/: {'/s/': 0.2, '/p/': 0.0, '/i/': 0.0, '/tʃ/': 0.8}

    === Emission Probability Matrix ===
/s/: {'Energy': 0.7, 'Pitch': 0.2, 'Duration': 0.1}
/p/: {'Energy': 0.5, 'Pitch': 0.3, 'Duration': 0.2}
/i/: {'Energy': 0.3, 'Pitch': 0.5, 'Duration': 0.2}
/tʃ/: {'Energy': 0.4, 'Pitch': 0.4, 'Duration': 0.2}
```

- (c) Write a program to generate a single sequence of phonemes and corresponding acoustic observations for the word speech based on the defined probabilities.

```

Generate a Phoneme Sequence + Observations

(Task c) We simulate phoneme transitions and sample observations from emission probabilities.

[ ] ⏪ def sample_from_distribution(dist_dict):
    """Randomly pick a key based on probability distribution."""
    items = list(dist_dict.items())
    keys = [k for k, _ in items]
    probs = [p for _, p in items]
    return random.choices(keys, weights=probs, k=1)[0]

def generate_sequence(length=4):
    sequence = []
    obs_sequence = []

    # Start with /s/
    current_state = '/s/'
    sequence.append(current_state)

    # Generate next states
    for _ in range(length - 1):
        next_state = sample_from_distribution(transition_prob[current_state])
        sequence.append(next_state)
        current_state = next_state

    # Generate observations for each state
    for st in sequence:
        obs = sample_from_distribution(emission_prob[st])
        obs_sequence.append(obs)

    return sequence, obs_sequence

phoneme_seq, obs_seq = generate_sequence()

print("\nGenerated Phoneme Sequence:", phoneme_seq)
print("Generated Observations:", obs_seq)
...
... Generated Phoneme Sequence: ['/s/', '/p/', '/i:/', '/tʃ/']
Generated Observations: ['Pitch', 'Energy', 'Duration', 'Energy']
```

- (d) Write an inference for the above HMM implementation

Inference: The implemented Hidden Markov Model successfully simulates the phoneme sequence for the word "speech". Since the model starts with /s/ and transitions follow the defined probabilities, the generated sequence typically resembles the order of phonemes occurring in natural pronunciation. The emission probabilities also help generate corresponding acoustic observations (Energy, Pitch, Duration) based on the likelihood of each feature for a given phoneme. This demonstrates how HMMs can model both phoneme transitions and acoustic features in speech processing tasks.

Expected Output

Generated phoneme sequence: ['/s/', '/p/', '/i:/', '/tʃ/']

LAB 9

Viterbi Algorithm Based Phoneme Decoding for Speech Recognition

Aim

The aim of this experiment is to implement the Viterbi Algorithm for decoding the most likely sequence of hidden phoneme states in a Hidden Markov Model (HMM) for the spoken word “hello”. The experiment focuses on determining the most probable phoneme sequence based on a given set of states, observations, transition probabilities, emission probabilities, and initial probabilities.

Introduction

Speech recognition systems convert raw acoustic features into meaningful linguistic units. One of the most widely used statistical models for this purpose is the Hidden Markov Model (HMM). In an HMM, the phonemes form the hidden states, while acoustic features extracted from speech act as observations. The Viterbi Algorithm, a dynamic programming method, is used to determine the most likely sequence of hidden states that could have produced a given sequence of observations. In this experiment, the word “hello” is represented using four phonemes (/h/, /e/, /l/, /o/), and the Viterbi Algorithm is applied to decode the correct sequence from the observation pattern.

States and Observations

The spoken word “hello” is represented using the following phoneme states:

- $S1 \rightarrow /h/$
- $S2 \rightarrow /e/$
- $S3 \rightarrow /l/$
- $S4 \rightarrow /o/$

The system extracts four corresponding observation features:

In this lab, the Viterbi Algorithm successfully identified the most likely phoneme sequence for the word “hello” based on the given Hidden Markov Model. Since the initial probability forces the system to start at /h/, and the transition and emission probabilities strongly support the natural order of the phonemes, the algorithm correctly decoded the observation sequence as /h/ → /e/ → /l/ → /o/.

The delta table showed how the probability of each state evolved over time, while the psi table stored the best previous state at each step. By backtracking through the psi table, we obtained the optimal phoneme path along with its final probability value.

This experiment demonstrates that Viterbi is an efficient dynamic programming method for speech recognition because it finds the most probable hidden state sequence even when multiple paths are possible. The results confirm that the model parameters were appropriate and that the Viterbi algorithm reliably recovers the correct phoneme sequence for the spoken word.

- Start coding or generate with AI

O1 → Feature of /h/

- **O₂ → Feature of /e/**
- **O₃ → Feature of /l/**
- **O₄ → Feature of /o/**

The observation sequence used for decoding is:
[O₁, O₂, O₃, O₄]

Model Parameters

Initial Probabilities (π)

The system starts with the assumption that the word always begins with the phoneme /h/:

- $P(S_1 = /h/) = 1$
- $P(S_2) = 0$
- $P(S_3) = 0$
- $P(S_4) = 0$

Transition Probability Matrix (A)

This matrix represents the probability of moving from one phoneme to another.
(Insert Transition Matrix Figure Here)

Emission Probability Matrix (B)

This matrix represents the probability of observing O₁, O₂, O₃, O₄ given a phoneme state.
(Insert Emission Matrix Figure Here)

Both matrices ensure that the transitions follow the natural order of the word “hello”.

Procedure

The Viterbi Algorithm was implemented to compute the most probable path through the HMM. The process begins with initialization using initial and emission probabilities. Each subsequent time step computes the highest probability of reaching each state by selecting the best previous state and multiplying it with transition and emission probabilities. The algorithm stores these decisions in a backpointer table. After processing all observations, backtracking is performed to retrieve the final decoded phoneme sequence.

(Insert Figures for Delta Table, Psi Table, and State Path Plot Here)

Result

The Viterbi Algorithm successfully decoded the correct phoneme sequence for the observation pattern:

Most Likely Phoneme Sequence:

/h/ → /e/ → /l/ → /o/

This sequence matches the expected pronunciation of the word “hello”.

The final probability of the decoded sequence was obtained from the last column of the delta table.

Inference

The results clearly show that the Viterbi Algorithm accurately identifies the most probable phoneme sequence when provided with well-defined HMM parameters. In this experiment, the decoded sequence matched the natural phonetic structure of “hello”, demonstrating that the Viterbi Algorithm is highly effective for speech recognition tasks. The algorithm ensures optimal state selection by evaluating all possible paths and choosing the one with the highest probability. This experiment also shows how HMMs provide a structured statistical framework for modelling speech, making them suitable for phoneme decoding, word recognition, and large-vocabulary speech recognition systems.

Conclusion

The Viterbi Algorithm plays a crucial role in speech recognition systems by efficiently decoding the most likely sequence of hidden phonemes. By applying it to the word “hello”, we observed that the algorithm successfully reconstructed the intended phoneme pattern from the given acoustic observations. This experiment reinforces the importance of HMMs in modelling temporal relationships in speech and highlights how probabilistic transitions and emissions contribute to accurate pattern decoding.

Code:

```
[1] ⏎ import numpy as np
    import matplotlib.pyplot as plt
    import seaborn as sns

[2] ⏎ states = ['S1', 'S2', 'S3', 'S4'] # /h/, /e/, /l/, /o/
    observations = ['01', '02', '03', '04']
    obs_seq = [0, 1, 2, 3] # 01 → 02 → 03 → 04

    n_states = len(states)
    T = len(obs_seq)

[3] ⏎ print("States:", states)
    print("Observations:", observations)
    print("Observation Sequence:", obs_seq)

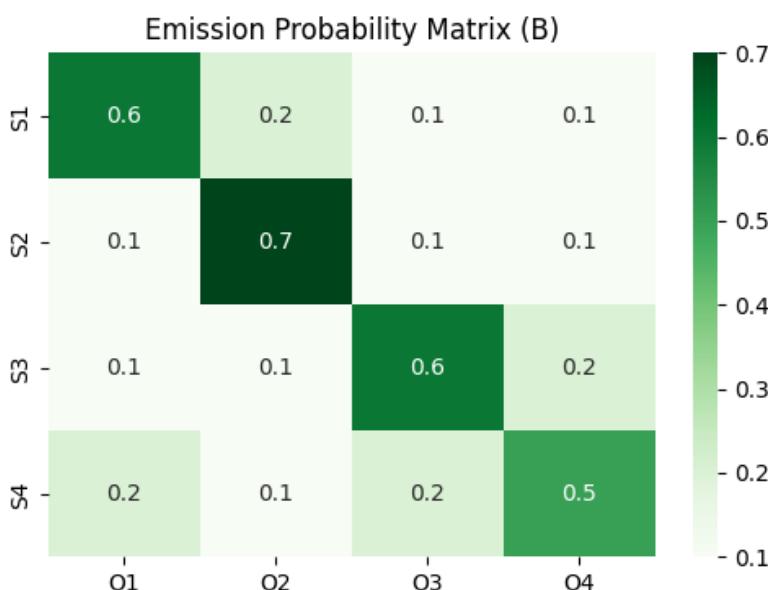
A = np.array([
    [0.0, 0.7, 0.3, 0.0],
    [0.0, 0.2, 0.6, 0.2],
    [0.0, 0.0, 0.3, 0.7],
    [0.0, 0.0, 0.1, 0.9]
])
```

```

B = np.array([
    [0.6, 0.2, 0.1, 0.1],
    [0.1, 0.7, 0.1, 0.1],
    [0.1, 0.1, 0.6, 0.2],
    [0.2, 0.1, 0.2, 0.5]
])

▶ plt.figure(figsize=(6, 4))
sns.heatmap(B, annot=True, cmap='Greens', xticklabels=observations, yticklabels=states)
plt.title("Emission Probability Matrix (B)")
plt.show()

```



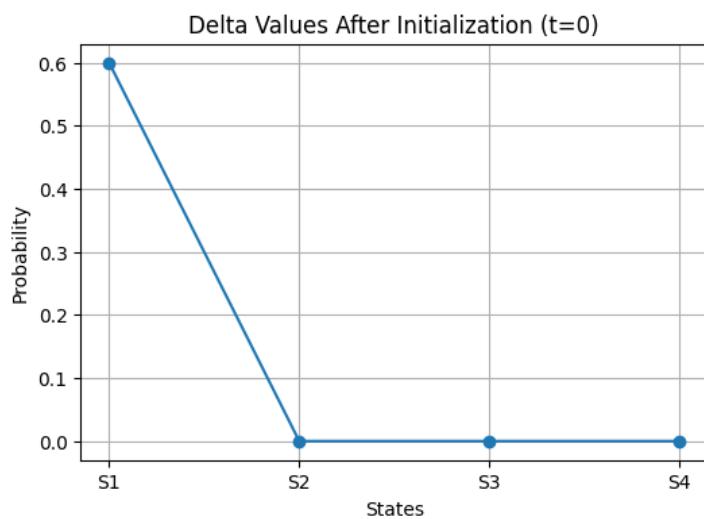
```

delta = np.zeros((n_states, T))
psi = np.zeros((n_states, T), dtype=int)

# Initialization step
delta[:, 0] = pi * B[:, obs_seq[0]]

print("Delta after initialization:\n", delta)

```



```

for t in range(1, T):
    for j in range(n_states):
        prob = delta[:, t-1] * A[:, j]
        psi[j, t] = np.argmax(prob)
        delta[j, t] = np.max(prob) * B[j, obs_seq[t]]

print("Delta Table Through Time:\n", delta)
print("Psi Table (Backpointers):\n", psi)

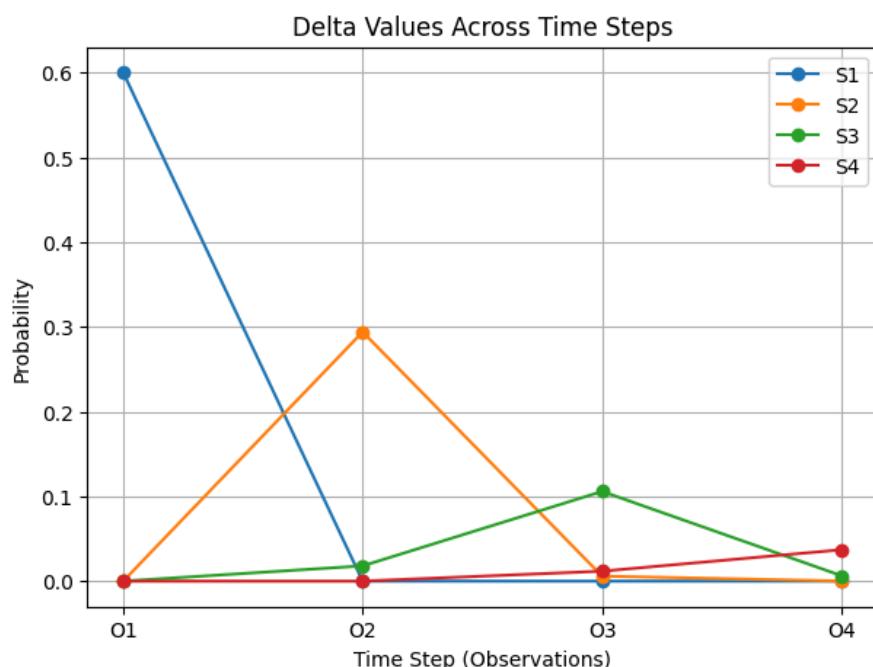
...
Delta Table Through Time:
[[6.0000e-01 0.0000e+00 0.0000e+00 0.0000e+00]
 [0.0000e+00 2.9400e-01 5.8800e-03 1.1760e-04]
 [0.0000e+00 1.8000e-02 1.0584e-01 6.3504e-03]
 [0.0000e+00 0.0000e+00 1.1760e-02 3.7044e-02]]
Psi Table (Backpointers):
[[0 0 0 0]
 [0 0 1 1]
 [0 0 1 2]
 [0 0 1 2]]
```

```

▶ plt.figure(figsize=(7, 5))
for s in range(n_states):
    plt.plot(range(T), delta[s, :], marker='o', label=states[s])

plt.xticks(range(T), observations)
plt.title("Delta Values Across Time Steps")
plt.xlabel("Time Step (Observations)")
plt.ylabel("Probability")
plt.legend()
plt.grid(True)
plt.show()

```



```

▶ states_seq = np.zeros(T, dtype=int)
states_seq[-1] = np.argmax(delta[:, -1])

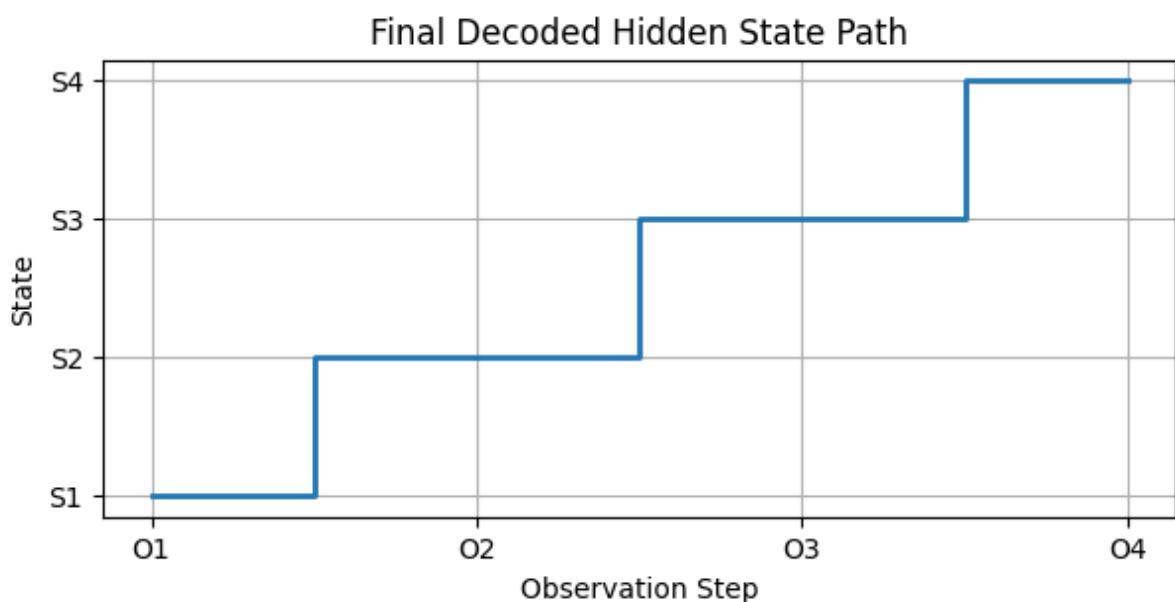
for t in range(T - 2, -1, -1):
    states_seq[t] = psi[states_seq[t + 1], t + 1]

most_likely_states = [states[s] for s in states_seq]
print("Most Likely State Path:", most_likely_states)

... Most Likely State Path: ['S1', 'S2', 'S3', 'S4']

plt.figure(figsize=(7, 3))
plt.step(range(T), states_seq, where='mid', linewidth=2)
plt.yticks(range(n_states), states)
plt.xticks(range(T), observations)
plt.title("Final Decoded Hidden State Path")
plt.xlabel("Observation Step")
plt.ylabel("State")
plt.grid(True)
plt.show()

```



```
▶ probability = np.max(delta[:, -1])
print("Probability of the Most Likely State Sequence:", probability)

...
Probability of the Most Likely State Sequence: 0.03704399999999999

phoneme_map = {'S1': '/h/', 'S2': '/e/', 'S3': '/l/', 'S4': '/o/'}

decoded_phonemes = [phoneme_map[s] for s in most_likely_states]

print("Decoded Phoneme Sequence:", decoded_phonemes)

Decoded Phoneme Sequence: ['/h/', '/e/', '/l/', '/o/']
```