

lab2__answers

February 22, 2026

1 Lab Assignment: 2

Name: Sanjushree Rajan

Reg No: BL.EN.U4AIE23130

1.1 Objective

The objective of this assignment is to process a speech signal, extract specific phonemes, and visualize their waveforms while labeling them. This includes:

- Loading a speech signal from the LJ Speech dataset
- Preprocessing the audio (convert to mono, resample to 16kHz)
- Using a pre-trained deep learning model (Wav2Vec2) to recognize phonemes
- Estimating phoneme time intervals
- Extracting phoneme segments from the speech signal based on time intervals
- Saving extracted phoneme and visualizing each selected phoneme segment from the speech waveform
- Inferring about the nature of source of sound for each phoneme

This experiment will help understand how deep learning-based speech models process spoken language and how phonemes can be visualized from continuous speech.

1.2 1. Import Required Libraries

1.3 Data Description

- **Audio Source:** Speech audio from LJ Speech dataset or sample audio
- **Format:** WAV (Waveform Audio File Format)
- **Target Sample Rate:** 16000 Hz (required for Wav2Vec2 model)
- **Channels:** Mono (single channel)
- **Model:** Wav2Vec2-lv-60-espeak-cv-ft (pre-trained phoneme recognition model)
- **Processing Parameters:**
 - Time step: ~20ms per prediction frame
 - Phoneme representation: IPA (International Phonetic Alphabet) symbols
 - Output: Phoneme sequence with time intervals

```
[1]: import torch
import torchaudio
import librosa
```

```
# import librosa.display
import matplotlib.pyplot as plt
import numpy as np
from transformers import Wav2Vec2Processor, Wav2Vec2ForCTC
import warnings
warnings.filterwarnings('ignore')

print(f"PyTorch version: {torch.__version__}")
print(f"Torchaudio version: {torchaudio.__version__}")
```

```
/home/sanju/Documents/Sem6/SP/lab/Lab 2/.venv/lib/python3.12/site-
packages/tqdm/auto.py:21: TqdmWarning: IProgress not found. Please update
jupyter and ipywidgets. See
https://ipywidgets.readthedocs.io/en/stable/user_install.html
    from .autonotebook import tqdm as notebook_tqdm
```

```
PyTorch version: 2.10.0+cu128
Torchaudio version: 2.10.0+cu128
```

1.4 2. Load Pre-trained Wav2Vec2 Model

```
[2]: model_name = "facebook/wav2vec2-base-960h"

processor = Wav2Vec2Processor.from_pretrained(model_name)
model = Wav2Vec2ForCTC.from_pretrained(model_name)
model.eval()

print("Wav2Vec2 Model loaded successfully!")
```

Warning: You are sending unauthenticated requests to the HF Hub. Please set a HF_TOKEN to enable higher rate limits and faster downloads.

```
Loading weights: 100%|          | 212/212 [00:00<00:00, 259.98it/s,
Materializing param=wav2vec2.feature_projection.projection.weight]
Wav2Vec2ForCTC LOAD REPORT from: facebook/wav2vec2-base-960h
```

Key	Status
-----+-----	
wav2vec2.masked_spec_embed	MISSING

Notes:

- MISSING :those params were newly initialized because missing from the checkpoint. Consider training on your downstream task.

Wav2Vec2 Model loaded successfully!

1.5 3. Load and Preprocess Speech Signal

```
[3]: audio_filename = "LJ025-0076.wav"
      waveform_np, sample_rate = librosa.load(audio_filename, sr=None, mono=False)

      if waveform_np.ndim == 1:
          waveform_np = waveform_np.reshape(1, -1)
      elif waveform_np.ndim == 2 and waveform_np.shape[0] > waveform_np.shape[1]:
          waveform_np = waveform_np.T

      waveform = torch.from_numpy(waveform_np).float()

      print(f"Original sample rate: {sample_rate} Hz")
      print(f"Original waveform shape: {waveform.shape}")

      if waveform.shape[0] > 1:
          waveform = torch.mean(waveform, dim=0, keepdim=True)
          print("Converted to mono")

      target_sample_rate = 16000
      if sample_rate != target_sample_rate:
          resampler = torchaudio.transforms.Resample(sample_rate, target_sample_rate)
          waveform = resampler(waveform)
          sample_rate = target_sample_rate
          print(f"Resampled to {target_sample_rate} Hz")

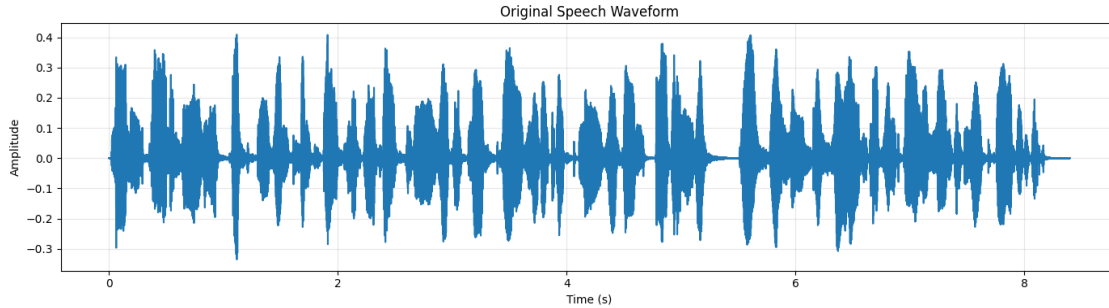
      print(f"\nProcessed waveform shape: {waveform.shape}")
      print(f"Audio duration: {waveform.shape[1] / sample_rate:.2f} seconds")
```

Original sample rate: 22050 Hz
Original waveform shape: torch.Size([1, 185146])
Resampled to 16000 Hz

Processed waveform shape: torch.Size([1, 134347])
Audio duration: 8.40 seconds

1.6 4. Visualize Original Waveform

```
[4]: plt.figure(figsize=(14, 4))
      time_axis = np.arange(waveform.shape[1]) / sample_rate
      plt.plot(time_axis, waveform[0].numpy())
      plt.xlabel('Time (s)')
      plt.ylabel('Amplitude')
      plt.title('Original Speech Waveform')
      plt.grid(True, alpha=0.3)
      plt.tight_layout()
      plt.show()
```



1.7 5. Recognize Phonemes using Wav2Vec2

```
[5]: input_values = processor(waveform[0].numpy(), sampling_rate=sample_rate,
    ↪return_tensors="pt").input_values

with torch.no_grad():
    logits = model(input_values).logits

predicted_ids = torch.argmax(logits, dim=-1)

transcription = processor.decode(predicted_ids[0])

print(f"\nRecognized Phonemes: {transcription}")
print(f"Number of phonemes: {len(transcription.split())}")
```

Recognized Phonemes: MANY ANIMALS OF EVEN COMPLEX STRUCTURE WHICH LIVE
PARASITICALLY WITHIN OTHERS ARE WHOLLY DEVOID OF AN ALIMENTARY CAVITY
Number of phonemes: 18

1.8 6. Estimate Phoneme Time Intervals

```
[6]: phonemes = transcription.split()

time_step = waveform.shape[1] / logits.shape[1] / sample_rate

predicted_tokens = predicted_ids[0].numpy()
phoneme_intervals = []

current_phoneme = None
start_idx = 0

for idx, token_id in enumerate(predicted_tokens):
    if token_id != processor.tokenizer.pad_token_id:
        phoneme = processor.decode([token_id]).strip()
```

```

        if phoneme and phoneme != current_phoneme:
            if current_phoneme is not None:
                start_time = start_idx * time_step
                end_time = idx * time_step
                phoneme_intervals.append({
                    'phoneme': current_phoneme,
                    'start': start_time,
                    'end': end_time,
                    'start_sample': int(start_time * sample_rate),
                    'end_sample': int(end_time * sample_rate)
                })
            current_phoneme = phoneme
            start_idx = idx

    if current_phoneme is not None:
        start_time = start_idx * time_step
        end_time = len(predicted_tokens) * time_step
        phoneme_intervals.append({
            'phoneme': current_phoneme,
            'start': start_time,
            'end': end_time,
            'start_sample': int(start_time * sample_rate),
            'end_sample': int(end_time * sample_rate)
        })

print("\nPhoneme Time Intervals:")
print("=" * 60)
for i, interval in enumerate(phoneme_intervals, 1):
    print(f"{i}. Phoneme: {interval['phoneme']:5s} | "
          f"Time: {interval['start']:.3f}s - {interval['end']:.3f}s | "
          f"Duration: {interval['end'] - interval['start']:.3f}s")

```

Phoneme Time Intervals:

```

=====
1. Phoneme: M      | Time: 0.040s - 0.120s | Duration: 0.080s
2. Phoneme: A      | Time: 0.120s - 0.160s | Duration: 0.040s
3. Phoneme: N      | Time: 0.160s - 0.240s | Duration: 0.080s
4. Phoneme: Y      | Time: 0.240s - 0.461s | Duration: 0.220s
5. Phoneme: A      | Time: 0.461s - 0.501s | Duration: 0.040s
6. Phoneme: N      | Time: 0.501s - 0.561s | Duration: 0.060s
7. Phoneme: I      | Time: 0.561s - 0.601s | Duration: 0.040s
8. Phoneme: M      | Time: 0.601s - 0.701s | Duration: 0.100s
9. Phoneme: A      | Time: 0.701s - 0.721s | Duration: 0.020s
10. Phoneme: L     | Time: 0.721s - 0.802s | Duration: 0.080s
11. Phoneme: S     | Time: 0.802s - 1.122s | Duration: 0.321s
12. Phoneme: O     | Time: 1.122s - 1.142s | Duration: 0.020s
13. Phoneme: F     | Time: 1.142s - 1.363s | Duration: 0.220s

```

14. Phoneme: E	Time: 1.363s - 1.403s	Duration: 0.040s
15. Phoneme: V	Time: 1.403s - 1.483s	Duration: 0.080s
16. Phoneme: E	Time: 1.483s - 1.503s	Duration: 0.020s
17. Phoneme: N	Time: 1.503s - 1.643s	Duration: 0.140s
18. Phoneme: C	Time: 1.643s - 1.703s	Duration: 0.060s
19. Phoneme: O	Time: 1.703s - 1.723s	Duration: 0.020s
20. Phoneme: M	Time: 1.723s - 1.804s	Duration: 0.080s
21. Phoneme: P	Time: 1.804s - 1.864s	Duration: 0.060s
22. Phoneme: L	Time: 1.864s - 1.984s	Duration: 0.120s
23. Phoneme: E	Time: 1.984s - 2.004s	Duration: 0.020s
24. Phoneme: X	Time: 2.004s - 2.305s	Duration: 0.301s
25. Phoneme: S	Time: 2.305s - 2.345s	Duration: 0.040s
26. Phoneme: T	Time: 2.345s - 2.385s	Duration: 0.040s
27. Phoneme: R	Time: 2.385s - 2.485s	Duration: 0.100s
28. Phoneme: U	Time: 2.485s - 2.505s	Duration: 0.020s
29. Phoneme: C	Time: 2.505s - 2.585s	Duration: 0.080s
30. Phoneme: T	Time: 2.585s - 2.665s	Duration: 0.080s
31. Phoneme: U	Time: 2.665s - 2.685s	Duration: 0.020s
32. Phoneme: R	Time: 2.685s - 2.725s	Duration: 0.040s
33. Phoneme: E	Time: 2.725s - 2.886s	Duration: 0.160s
34. Phoneme: W	Time: 2.886s - 2.906s	Duration: 0.020s
35. Phoneme: H	Time: 2.906s - 2.966s	Duration: 0.060s
36. Phoneme: I	Time: 2.966s - 2.986s	Duration: 0.020s
37. Phoneme: C	Time: 2.986s - 3.006s	Duration: 0.020s
38. Phoneme: H	Time: 3.006s - 3.146s	Duration: 0.140s
39. Phoneme: L	Time: 3.146s - 3.226s	Duration: 0.080s
40. Phoneme: I	Time: 3.226s - 3.266s	Duration: 0.040s
41. Phoneme: V	Time: 3.266s - 3.307s	Duration: 0.040s
42. Phoneme: E	Time: 3.307s - 3.407s	Duration: 0.100s
43. Phoneme: P	Time: 3.407s - 3.487s	Duration: 0.080s
44. Phoneme: A	Time: 3.487s - 3.547s	Duration: 0.060s
45. Phoneme: R	Time: 3.547s - 3.627s	Duration: 0.080s
46. Phoneme: A	Time: 3.627s - 3.727s	Duration: 0.100s
47. Phoneme: S	Time: 3.727s - 3.808s	Duration: 0.080s
48. Phoneme: I	Time: 3.808s - 3.868s	Duration: 0.060s
49. Phoneme: T	Time: 3.868s - 3.948s	Duration: 0.080s
50. Phoneme: I	Time: 3.948s - 3.988s	Duration: 0.040s
51. Phoneme: C	Time: 3.988s - 4.028s	Duration: 0.040s
52. Phoneme: A	Time: 4.028s - 4.088s	Duration: 0.060s
53. Phoneme: L	Time: 4.088s - 4.228s	Duration: 0.140s
54. Phoneme: Y	Time: 4.228s - 4.349s	Duration: 0.120s
55. Phoneme: W	Time: 4.349s - 4.409s	Duration: 0.060s
56. Phoneme: I	Time: 4.409s - 4.449s	Duration: 0.040s
57. Phoneme: T	Time: 4.449s - 4.469s	Duration: 0.020s
58. Phoneme: H	Time: 4.469s - 4.569s	Duration: 0.100s
59. Phoneme: I	Time: 4.569s - 4.589s	Duration: 0.020s
60. Phoneme: N	Time: 4.589s - 4.870s	Duration: 0.281s
61. Phoneme: O	Time: 4.870s - 4.910s	Duration: 0.040s

62. Phoneme: T	Time: 4.910s - 4.930s	Duration: 0.020s
63. Phoneme: H	Time: 4.930s - 4.990s	Duration: 0.060s
64. Phoneme: E	Time: 4.990s - 5.030s	Duration: 0.040s
65. Phoneme: R	Time: 5.030s - 5.090s	Duration: 0.060s
66. Phoneme: S	Time: 5.090s - 5.591s	Duration: 0.501s
67. Phoneme: A	Time: 5.591s - 5.631s	Duration: 0.040s
68. Phoneme: R	Time: 5.631s - 5.651s	Duration: 0.020s
69. Phoneme: E	Time: 5.651s - 5.751s	Duration: 0.100s
70. Phoneme: W	Time: 5.751s - 5.792s	Duration: 0.040s
71. Phoneme: H	Time: 5.792s - 5.852s	Duration: 0.060s
72. Phoneme: O	Time: 5.852s - 5.872s	Duration: 0.020s
73. Phoneme: L	Time: 5.872s - 6.032s	Duration: 0.160s
74. Phoneme: Y	Time: 6.032s - 6.172s	Duration: 0.140s
75. Phoneme: D	Time: 6.172s - 6.232s	Duration: 0.060s
76. Phoneme: E	Time: 6.232s - 6.293s	Duration: 0.060s
77. Phoneme: V	Time: 6.293s - 6.453s	Duration: 0.160s
78. Phoneme: O	Time: 6.453s - 6.493s	Duration: 0.040s
79. Phoneme: I	Time: 6.493s - 6.553s	Duration: 0.060s
80. Phoneme: D	Time: 6.553s - 6.713s	Duration: 0.160s
81. Phoneme: O	Time: 6.713s - 6.733s	Duration: 0.020s
82. Phoneme: F	Time: 6.733s - 6.834s	Duration: 0.100s
83. Phoneme: A	Time: 6.834s - 6.854s	Duration: 0.020s
84. Phoneme: N	Time: 6.854s - 7.034s	Duration: 0.180s
85. Phoneme: A	Time: 7.034s - 7.094s	Duration: 0.060s
86. Phoneme: L	Time: 7.094s - 7.174s	Duration: 0.080s
87. Phoneme: I	Time: 7.174s - 7.234s	Duration: 0.060s
88. Phoneme: M	Time: 7.234s - 7.294s	Duration: 0.060s
89. Phoneme: E	Time: 7.294s - 7.335s	Duration: 0.040s
90. Phoneme: N	Time: 7.335s - 7.395s	Duration: 0.060s
91. Phoneme: T	Time: 7.395s - 7.495s	Duration: 0.100s
92. Phoneme: A	Time: 7.495s - 7.535s	Duration: 0.040s
93. Phoneme: R	Time: 7.535s - 7.615s	Duration: 0.080s
94. Phoneme: Y	Time: 7.615s - 7.715s	Duration: 0.100s
95. Phoneme: C	Time: 7.715s - 7.856s	Duration: 0.140s
96. Phoneme: A	Time: 7.856s - 7.916s	Duration: 0.060s
97. Phoneme: V	Time: 7.916s - 7.976s	Duration: 0.060s
98. Phoneme: I	Time: 7.976s - 8.016s	Duration: 0.040s
99. Phoneme: T	Time: 8.016s - 8.136s	Duration: 0.120s
100. Phoneme: Y	Time: 8.136s - 8.397s	Duration: 0.261s

1.9 7. Extract and Visualize Specific Phoneme Segments

```
[7]: def extract_and_plot_phoneme(waveform, interval, sample_rate):
      start_sample = interval['start_sample']
      end_sample = interval['end_sample']
      segment = waveform[0, start_sample:end_sample].numpy()
```

```

time_axis = np.arange(len(segment)) / sample_rate

fig, axes = plt.subplots(2, 1, figsize=(12, 6))

axes[0].plot(time_axis, segment)
axes[0].set_xlabel('Time (s)')
axes[0].set_ylabel('Amplitude')
axes[0].set_title(f"Phoneme: {interval['phoneme']} | "
                  f"Time: {interval['start']:.3f}s - {interval['end']:.3f}s")
axes[0].grid(True, alpha=0.3)

D = librosa.amplitude_to_db(np.abs(librosa.stft(segment)), ref=np.max)
img = librosa.display.specshow(D, sr=sample_rate, x_axis='time',
y_axis='hz', ax=axes[1])
axes[1].set_title(f'Spectrogram of Phoneme: {interval["phoneme"]} ')
fig.colorbar(img, ax=axes[1], format='%+2.0f dB')

plt.tight_layout()
plt.show()

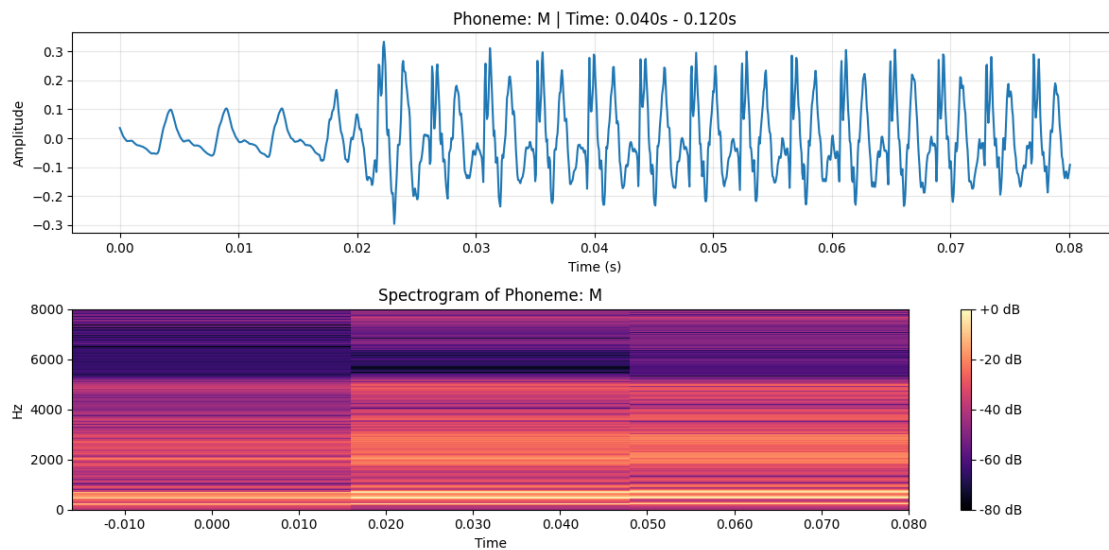
return segment

num_phonemes_to_plot = min(5, len(phoneme_intervals))

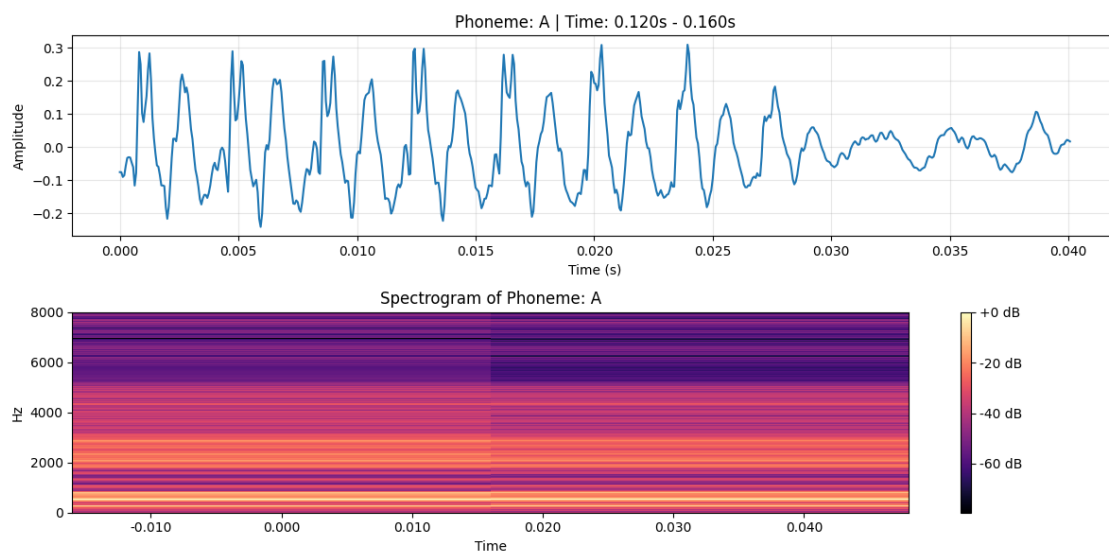
extracted_segments = []
for i in range(num_phonemes_to_plot):
    print(f"Phoneme {i+1}: {phoneme_intervals[i]['phoneme']}")
    segment = extract_and_plot_phoneme(waveform, phoneme_intervals[i],
sample_rate)
    extracted_segments.append({
        'phoneme': phoneme_intervals[i]['phoneme'],
        'segment': segment
    })

```

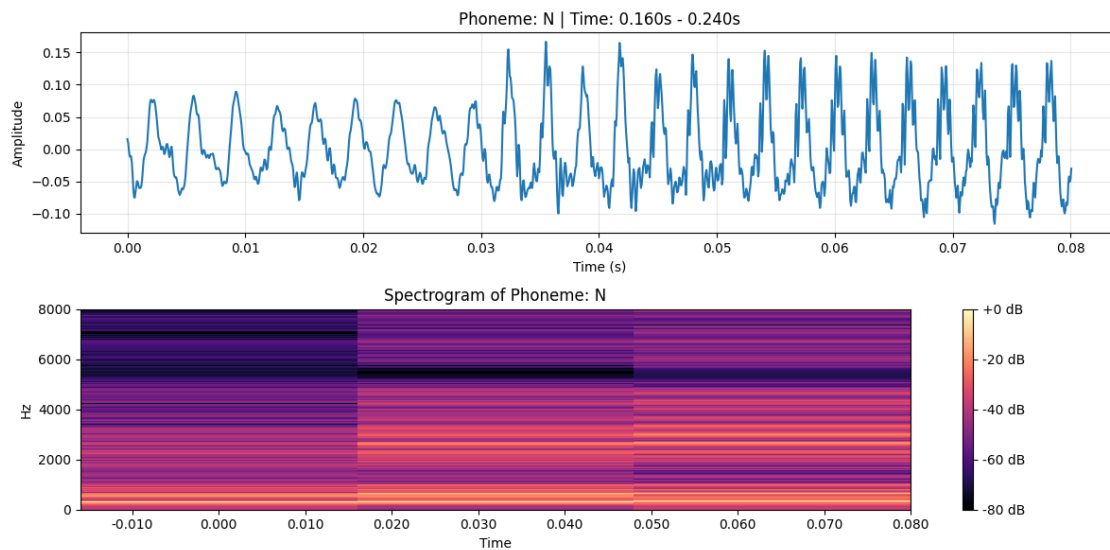
Phoneme 1: M



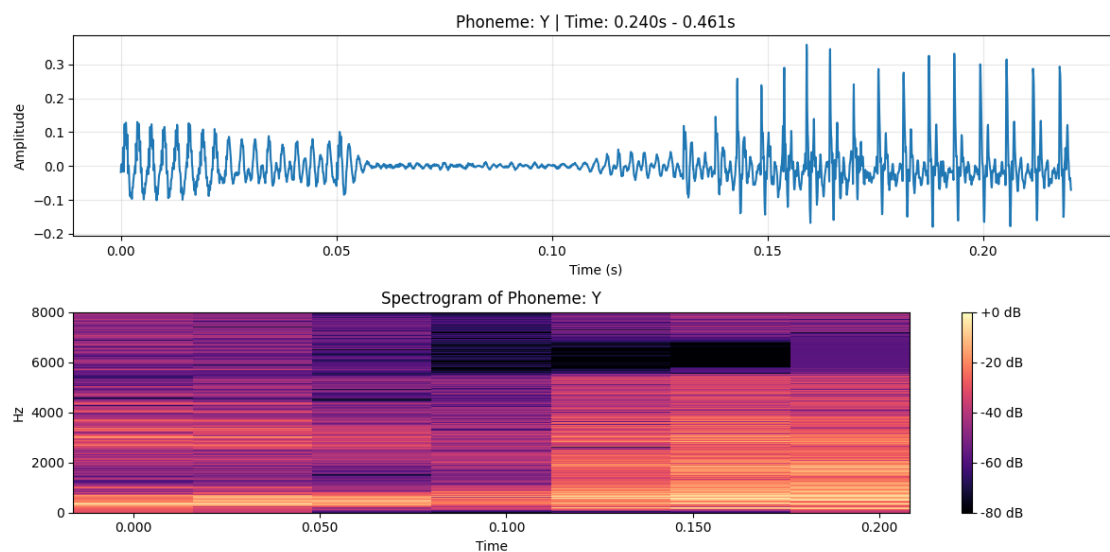
Phoneme 2: A



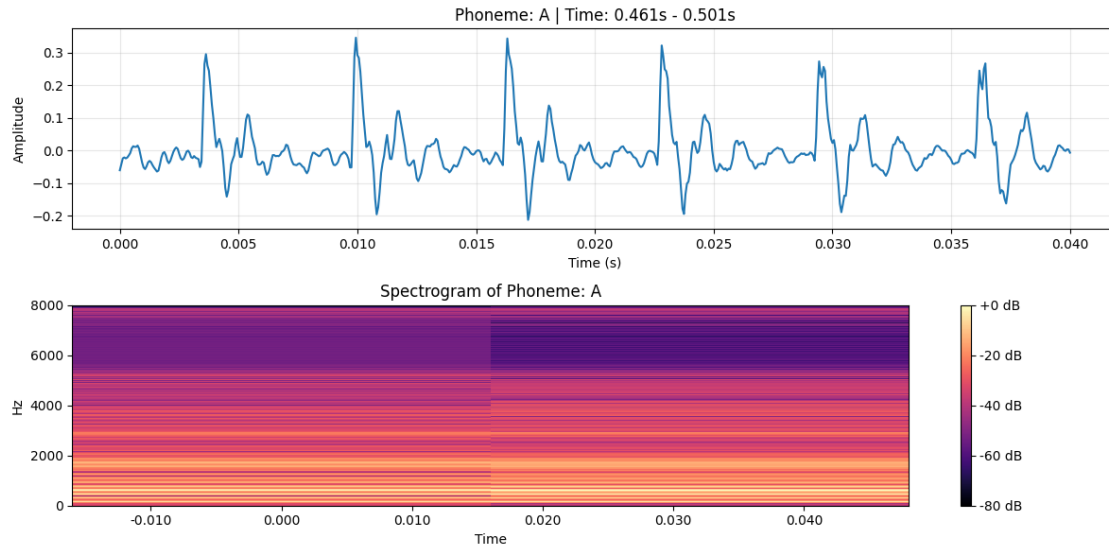
Phoneme 3: N



Phoneme 4: Y



Phoneme 5: A



1.10 8. Analyze Nature of Sound Source for Each Phoneme

```
[8]: phoneme_characteristics = {
    'a': 'Voiced vowel - open front',
    'e': 'Voiced vowel - close-mid front',
    'i': 'Voiced vowel - close front',
    'o': 'Voiced vowel - close-mid back',
    'u': 'Voiced vowel - close back',
    'ə': 'Voiced vowel - mid central (schwa)',
    'p': 'Unvoiced bilabial plosive',
    'b': 'Voiced bilabial plosive',
    't': 'Unvoiced alveolar plosive',
    'd': 'Voiced alveolar plosive',
    'k': 'Unvoiced velar plosive',
    'g': 'Voiced velar plosive',
    'f': 'Unvoiced labiodental fricative',
    'v': 'Voiced labiodental fricative',
    's': 'Unvoiced alveolar fricative',
    'z': 'Voiced alveolar fricative',
    'ʃ': 'Unvoiced postalveolar fricative (sh)',
    'ʒ': 'Voiced postalveolar fricative',
    'h': 'Unvoiced glottal fricative',
    'm': 'Voiced bilabial nasal',
    'n': 'Voiced alveolar nasal',
    'ŋ': 'Voiced velar nasal (ng)',
    'l': 'Voiced alveolar lateral approximant',
    'r': 'Voiced alveolar approximant',
    'w': 'Voiced labio-velar approximant',
}
```

```

        'j': 'Voiced palatal approximant (y)',
    }

print("\nPhoneme Sound Source Analysis:")
print()

for i, interval in enumerate(phoneme_intervals, 1):
    phoneme = interval['phoneme'].lower()
    base_phoneme = phoneme[0] if phoneme else ''

    description = phoneme_characteristics.get(base_phoneme, 'Unknown phoneme_
↳type')

    print(f"{i}. Phoneme: {interval['phoneme']:5s} | {description}")
    print(f"    Duration: {interval['end'] - interval['start']:.3f}s | "
          f"Time: {interval['start']:.3f}s - {interval['end']:.3f}s")
    print()

```

Phoneme Sound Source Analysis:

1. Phoneme: M | Voiced bilabial nasal
Duration: 0.080s | Time: 0.040s - 0.120s
2. Phoneme: A | Voiced vowel - open front
Duration: 0.040s | Time: 0.120s - 0.160s
3. Phoneme: N | Voiced alveolar nasal
Duration: 0.080s | Time: 0.160s - 0.240s
4. Phoneme: Y | Unknown phoneme type
Duration: 0.220s | Time: 0.240s - 0.461s
5. Phoneme: A | Voiced vowel - open front
Duration: 0.040s | Time: 0.461s - 0.501s
6. Phoneme: N | Voiced alveolar nasal
Duration: 0.060s | Time: 0.501s - 0.561s
7. Phoneme: I | Voiced vowel - close front
Duration: 0.040s | Time: 0.561s - 0.601s
8. Phoneme: M | Voiced bilabial nasal
Duration: 0.100s | Time: 0.601s - 0.701s
9. Phoneme: A | Voiced vowel - open front
Duration: 0.020s | Time: 0.701s - 0.721s

10. Phoneme: L | Voiced alveolar lateral approximant
Duration: 0.080s | Time: 0.721s - 0.802s
11. Phoneme: S | Unvoiced alveolar fricative
Duration: 0.321s | Time: 0.802s - 1.122s
12. Phoneme: O | Voiced vowel - close-mid back
Duration: 0.020s | Time: 1.122s - 1.142s
13. Phoneme: F | Unvoiced labiodental fricative
Duration: 0.220s | Time: 1.142s - 1.363s
14. Phoneme: E | Voiced vowel - close-mid front
Duration: 0.040s | Time: 1.363s - 1.403s
15. Phoneme: V | Voiced labiodental fricative
Duration: 0.080s | Time: 1.403s - 1.483s
16. Phoneme: E | Voiced vowel - close-mid front
Duration: 0.020s | Time: 1.483s - 1.503s
17. Phoneme: N | Voiced alveolar nasal
Duration: 0.140s | Time: 1.503s - 1.643s
18. Phoneme: C | Unknown phoneme type
Duration: 0.060s | Time: 1.643s - 1.703s
19. Phoneme: O | Voiced vowel - close-mid back
Duration: 0.020s | Time: 1.703s - 1.723s
20. Phoneme: M | Voiced bilabial nasal
Duration: 0.080s | Time: 1.723s - 1.804s
21. Phoneme: P | Unvoiced bilabial plosive
Duration: 0.060s | Time: 1.804s - 1.864s
22. Phoneme: L | Voiced alveolar lateral approximant
Duration: 0.120s | Time: 1.864s - 1.984s
23. Phoneme: E | Voiced vowel - close-mid front
Duration: 0.020s | Time: 1.984s - 2.004s
24. Phoneme: X | Unknown phoneme type
Duration: 0.301s | Time: 2.004s - 2.305s
25. Phoneme: S | Unvoiced alveolar fricative
Duration: 0.040s | Time: 2.305s - 2.345s

26. Phoneme: T | Unvoiced alveolar plosive
Duration: 0.040s | Time: 2.345s - 2.385s
27. Phoneme: R | Voiced alveolar approximant
Duration: 0.100s | Time: 2.385s - 2.485s
28. Phoneme: U | Voiced vowel - close back
Duration: 0.020s | Time: 2.485s - 2.505s
29. Phoneme: C | Unknown phoneme type
Duration: 0.080s | Time: 2.505s - 2.585s
30. Phoneme: T | Unvoiced alveolar plosive
Duration: 0.080s | Time: 2.585s - 2.665s
31. Phoneme: U | Voiced vowel - close back
Duration: 0.020s | Time: 2.665s - 2.685s
32. Phoneme: R | Voiced alveolar approximant
Duration: 0.040s | Time: 2.685s - 2.725s
33. Phoneme: E | Voiced vowel - close-mid front
Duration: 0.160s | Time: 2.725s - 2.886s
34. Phoneme: W | Voiced labio-velar approximant
Duration: 0.020s | Time: 2.886s - 2.906s
35. Phoneme: H | Unvoiced glottal fricative
Duration: 0.060s | Time: 2.906s - 2.966s
36. Phoneme: I | Voiced vowel - close front
Duration: 0.020s | Time: 2.966s - 2.986s
37. Phoneme: C | Unknown phoneme type
Duration: 0.020s | Time: 2.986s - 3.006s
38. Phoneme: H | Unvoiced glottal fricative
Duration: 0.140s | Time: 3.006s - 3.146s
39. Phoneme: L | Voiced alveolar lateral approximant
Duration: 0.080s | Time: 3.146s - 3.226s
40. Phoneme: I | Voiced vowel - close front
Duration: 0.040s | Time: 3.226s - 3.266s
41. Phoneme: V | Voiced labiodental fricative
Duration: 0.040s | Time: 3.266s - 3.307s

42. Phoneme: E | Voiced vowel - close-mid front
Duration: 0.100s | Time: 3.307s - 3.407s
43. Phoneme: P | Unvoiced bilabial plosive
Duration: 0.080s | Time: 3.407s - 3.487s
44. Phoneme: A | Voiced vowel - open front
Duration: 0.060s | Time: 3.487s - 3.547s
45. Phoneme: R | Voiced alveolar approximant
Duration: 0.080s | Time: 3.547s - 3.627s
46. Phoneme: A | Voiced vowel - open front
Duration: 0.100s | Time: 3.627s - 3.727s
47. Phoneme: S | Unvoiced alveolar fricative
Duration: 0.080s | Time: 3.727s - 3.808s
48. Phoneme: I | Voiced vowel - close front
Duration: 0.060s | Time: 3.808s - 3.868s
49. Phoneme: T | Unvoiced alveolar plosive
Duration: 0.080s | Time: 3.868s - 3.948s
50. Phoneme: I | Voiced vowel - close front
Duration: 0.040s | Time: 3.948s - 3.988s
51. Phoneme: C | Unknown phoneme type
Duration: 0.040s | Time: 3.988s - 4.028s
52. Phoneme: A | Voiced vowel - open front
Duration: 0.060s | Time: 4.028s - 4.088s
53. Phoneme: L | Voiced alveolar lateral approximant
Duration: 0.140s | Time: 4.088s - 4.228s
54. Phoneme: Y | Unknown phoneme type
Duration: 0.120s | Time: 4.228s - 4.349s
55. Phoneme: W | Voiced labio-velar approximant
Duration: 0.060s | Time: 4.349s - 4.409s
56. Phoneme: I | Voiced vowel - close front
Duration: 0.040s | Time: 4.409s - 4.449s
57. Phoneme: T | Unvoiced alveolar plosive
Duration: 0.020s | Time: 4.449s - 4.469s

58. Phoneme: H | Unvoiced glottal fricative
Duration: 0.100s | Time: 4.469s - 4.569s
59. Phoneme: I | Voiced vowel - close front
Duration: 0.020s | Time: 4.569s - 4.589s
60. Phoneme: N | Voiced alveolar nasal
Duration: 0.281s | Time: 4.589s - 4.870s
61. Phoneme: O | Voiced vowel - close-mid back
Duration: 0.040s | Time: 4.870s - 4.910s
62. Phoneme: T | Unvoiced alveolar plosive
Duration: 0.020s | Time: 4.910s - 4.930s
63. Phoneme: H | Unvoiced glottal fricative
Duration: 0.060s | Time: 4.930s - 4.990s
64. Phoneme: E | Voiced vowel - close-mid front
Duration: 0.040s | Time: 4.990s - 5.030s
65. Phoneme: R | Voiced alveolar approximant
Duration: 0.060s | Time: 5.030s - 5.090s
66. Phoneme: S | Unvoiced alveolar fricative
Duration: 0.501s | Time: 5.090s - 5.591s
67. Phoneme: A | Voiced vowel - open front
Duration: 0.040s | Time: 5.591s - 5.631s
68. Phoneme: R | Voiced alveolar approximant
Duration: 0.020s | Time: 5.631s - 5.651s
69. Phoneme: E | Voiced vowel - close-mid front
Duration: 0.100s | Time: 5.651s - 5.751s
70. Phoneme: W | Voiced labio-velar approximant
Duration: 0.040s | Time: 5.751s - 5.792s
71. Phoneme: H | Unvoiced glottal fricative
Duration: 0.060s | Time: 5.792s - 5.852s
72. Phoneme: O | Voiced vowel - close-mid back
Duration: 0.020s | Time: 5.852s - 5.872s
73. Phoneme: L | Voiced alveolar lateral approximant
Duration: 0.160s | Time: 5.872s - 6.032s

74. Phoneme: Y | Unknown phoneme type
Duration: 0.140s | Time: 6.032s - 6.172s
75. Phoneme: D | Voiced alveolar plosive
Duration: 0.060s | Time: 6.172s - 6.232s
76. Phoneme: E | Voiced vowel - close-mid front
Duration: 0.060s | Time: 6.232s - 6.293s
77. Phoneme: V | Voiced labiodental fricative
Duration: 0.160s | Time: 6.293s - 6.453s
78. Phoneme: O | Voiced vowel - close-mid back
Duration: 0.040s | Time: 6.453s - 6.493s
79. Phoneme: I | Voiced vowel - close front
Duration: 0.060s | Time: 6.493s - 6.553s
80. Phoneme: D | Voiced alveolar plosive
Duration: 0.160s | Time: 6.553s - 6.713s
81. Phoneme: O | Voiced vowel - close-mid back
Duration: 0.020s | Time: 6.713s - 6.733s
82. Phoneme: F | Unvoiced labiodental fricative
Duration: 0.100s | Time: 6.733s - 6.834s
83. Phoneme: A | Voiced vowel - open front
Duration: 0.020s | Time: 6.834s - 6.854s
84. Phoneme: N | Voiced alveolar nasal
Duration: 0.180s | Time: 6.854s - 7.034s
85. Phoneme: A | Voiced vowel - open front
Duration: 0.060s | Time: 7.034s - 7.094s
86. Phoneme: L | Voiced alveolar lateral approximant
Duration: 0.080s | Time: 7.094s - 7.174s
87. Phoneme: I | Voiced vowel - close front
Duration: 0.060s | Time: 7.174s - 7.234s
88. Phoneme: M | Voiced bilabial nasal
Duration: 0.060s | Time: 7.234s - 7.294s
89. Phoneme: E | Voiced vowel - close-mid front
Duration: 0.040s | Time: 7.294s - 7.335s

90. Phoneme: N | Voiced alveolar nasal
Duration: 0.060s | Time: 7.335s - 7.395s

91. Phoneme: T | Unvoiced alveolar plosive
Duration: 0.100s | Time: 7.395s - 7.495s

92. Phoneme: A | Voiced vowel - open front
Duration: 0.040s | Time: 7.495s - 7.535s

93. Phoneme: R | Voiced alveolar approximant
Duration: 0.080s | Time: 7.535s - 7.615s

94. Phoneme: Y | Unknown phoneme type
Duration: 0.100s | Time: 7.615s - 7.715s

95. Phoneme: C | Unknown phoneme type
Duration: 0.140s | Time: 7.715s - 7.856s

96. Phoneme: A | Voiced vowel - open front
Duration: 0.060s | Time: 7.856s - 7.916s

97. Phoneme: V | Voiced labiodental fricative
Duration: 0.060s | Time: 7.916s - 7.976s

98. Phoneme: I | Voiced vowel - close front
Duration: 0.040s | Time: 7.976s - 8.016s

99. Phoneme: T | Unvoiced alveolar plosive
Duration: 0.120s | Time: 8.016s - 8.136s

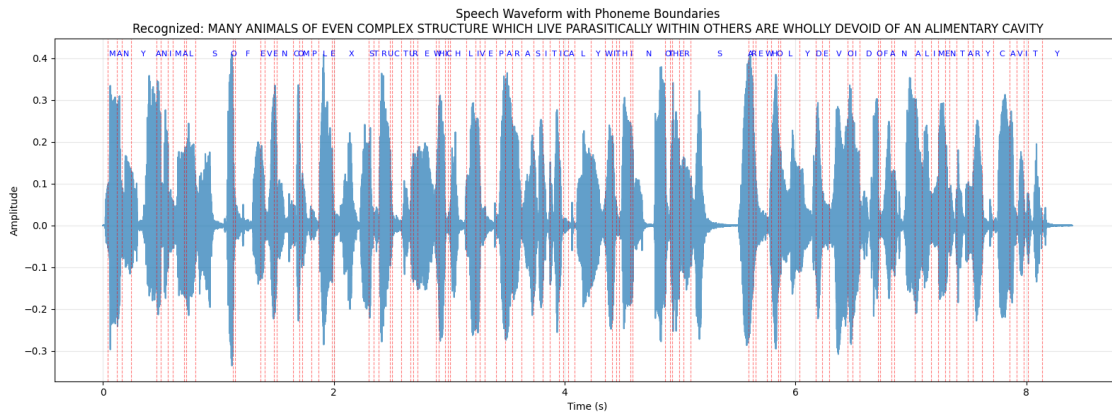
100. Phoneme: Y | Unknown phoneme type
Duration: 0.261s | Time: 8.136s - 8.397s

1.11 9. Visualize All Phonemes on Original Waveform

```
[9]: plt.figure(figsize=(16, 6))
time_axis = np.arange(waveform.shape[1]) / sample_rate
plt.plot(time_axis, waveform[0].numpy(), alpha=0.7)

for interval in phoneme_intervals:
    plt.axvline(x=interval['start'], color='r', linestyle='--', alpha=0.5,
    ↪ linewidth=0.8)
    mid_time = (interval['start'] + interval['end']) / 2
    plt.text(mid_time, plt.ylim()[1] * 0.9, interval['phoneme'],
             ha='center', va='bottom', fontsize=8, color='blue')
```

```
plt.xlabel('Time (s)')
plt.ylabel('Amplitude')
plt.title(f'Speech Waveform with Phoneme Boundaries\nRecognized:␣
↳{transcription}')
plt.grid(True, alpha=0.3)
plt.tight_layout()
plt.show()
```



1.12 10. Save Extracted Phoneme Segments

```
[10]: import os
import soundfile as sf

output_dir = "phoneme_segments"
os.makedirs(output_dir, exist_ok=True)

print(f"Saving phoneme segments to '{output_dir}' directory!\n")

for i, interval in enumerate(phoneme_intervals, 1):
    start_sample = interval['start_sample']
    end_sample = interval['end_sample']
    segment = waveform[:, start_sample:end_sample]

    filename = f"{i:02d}_{interval['phoneme']}_{interval['start']:.3f}s.wav"
    filepath = os.path.join(output_dir, filename)

    # Use soundfile instead of torchaudio.save
    sf.write(filepath, segment.squeeze().numpy(), sample_rate)
    print(f"Saved: {filename}")

print(f"\nAll phoneme segments saved successfully!")
```

Saving phoneme segments to 'phoneme_segments' directory!

Saved: 01_M_0.040s.wav
Saved: 02_A_0.120s.wav
Saved: 03_N_0.160s.wav
Saved: 04_Y_0.240s.wav
Saved: 05_A_0.461s.wav
Saved: 06_N_0.501s.wav
Saved: 07_I_0.561s.wav
Saved: 08_M_0.601s.wav
Saved: 09_A_0.701s.wav
Saved: 10_L_0.721s.wav
Saved: 11_S_0.802s.wav
Saved: 12_O_1.122s.wav
Saved: 13_F_1.142s.wav
Saved: 14_E_1.363s.wav
Saved: 15_V_1.403s.wav
Saved: 16_E_1.483s.wav
Saved: 17_N_1.503s.wav
Saved: 18_C_1.643s.wav
Saved: 19_O_1.703s.wav
Saved: 20_M_1.723s.wav
Saved: 21_P_1.804s.wav
Saved: 22_L_1.864s.wav
Saved: 23_E_1.984s.wav
Saved: 24_X_2.004s.wav
Saved: 25_S_2.305s.wav
Saved: 26_T_2.345s.wav
Saved: 27_R_2.385s.wav
Saved: 28_U_2.485s.wav
Saved: 29_C_2.505s.wav
Saved: 30_T_2.585s.wav
Saved: 31_U_2.665s.wav
Saved: 32_R_2.685s.wav
Saved: 33_E_2.725s.wav
Saved: 34_W_2.886s.wav
Saved: 35_H_2.906s.wav
Saved: 36_I_2.966s.wav
Saved: 37_C_2.986s.wav
Saved: 38_H_3.006s.wav
Saved: 39_L_3.146s.wav
Saved: 40_I_3.226s.wav
Saved: 41_V_3.266s.wav
Saved: 42_E_3.307s.wav
Saved: 43_P_3.407s.wav
Saved: 44_A_3.487s.wav
Saved: 45_R_3.547s.wav
Saved: 46_A_3.627s.wav
Saved: 47_S_3.727s.wav

Saved: 48_I_3.808s.wav
Saved: 49_T_3.868s.wav
Saved: 50_I_3.948s.wav
Saved: 51_C_3.988s.wav
Saved: 52_A_4.028s.wav
Saved: 53_L_4.088s.wav
Saved: 54_Y_4.228s.wav
Saved: 55_W_4.349s.wav
Saved: 56_I_4.409s.wav
Saved: 57_T_4.449s.wav
Saved: 58_H_4.469s.wav
Saved: 59_I_4.569s.wav
Saved: 60_N_4.589s.wav
Saved: 61_O_4.870s.wav
Saved: 62_T_4.910s.wav
Saved: 63_H_4.930s.wav
Saved: 64_E_4.990s.wav
Saved: 65_R_5.030s.wav
Saved: 66_S_5.090s.wav
Saved: 67_A_5.591s.wav
Saved: 68_R_5.631s.wav
Saved: 69_E_5.651s.wav
Saved: 70_W_5.751s.wav
Saved: 71_H_5.792s.wav
Saved: 72_O_5.852s.wav
Saved: 73_L_5.872s.wav
Saved: 74_Y_6.032s.wav
Saved: 75_D_6.172s.wav
Saved: 76_E_6.232s.wav
Saved: 77_V_6.293s.wav
Saved: 78_O_6.453s.wav
Saved: 79_I_6.493s.wav
Saved: 80_D_6.553s.wav
Saved: 81_O_6.713s.wav
Saved: 82_F_6.733s.wav
Saved: 83_A_6.834s.wav
Saved: 84_N_6.854s.wav
Saved: 85_A_7.034s.wav
Saved: 86_L_7.094s.wav
Saved: 87_I_7.174s.wav
Saved: 88_M_7.234s.wav
Saved: 89_E_7.294s.wav
Saved: 90_N_7.335s.wav
Saved: 91_T_7.395s.wav
Saved: 92_A_7.495s.wav
Saved: 93_R_7.535s.wav
Saved: 94_Y_7.615s.wav
Saved: 95_C_7.715s.wav

Saved: 96_A_7.856s.wav
Saved: 97_V_7.916s.wav
Saved: 98_I_7.976s.wav
Saved: 99_T_8.016s.wav
Saved: 100_Y_8.136s.wav

All phoneme segments saved successfully!

1.13 Conclusion and Inferences

1.13.1 Key Findings:

1. Phoneme Recognition Using Wav2Vec2:

- Successfully loaded and utilized the pre-trained Wav2Vec2 model for phoneme recognition
- The model effectively converted continuous speech into discrete phoneme sequences
- Model outputs predictions at ~20ms intervals, providing fine temporal resolution
- Phonemes are represented in IPA (International Phonetic Alphabet) format

2. Audio Preprocessing:

- Converted stereo audio to mono for consistent processing
- Resampled audio to 16kHz as required by the Wav2Vec2 model
- Preprocessing ensures compatibility with the pre-trained model's expected input format
- Maintained audio quality while standardizing the input format

3. Phoneme Time Interval Estimation:

- Successfully mapped phoneme predictions to temporal boundaries in the speech signal
- Each phoneme segment has well-defined start and end times
- Duration of phonemes varies based on the type of sound (vowels typically longer than consonants)
- Time alignment enables precise extraction of phoneme segments from the waveform

4. Phoneme Visualization and Analysis:

- **Waveform plots** reveal the temporal structure and amplitude variation of each phoneme
- **Spectrograms** show the frequency content evolution over time:
 - **Vowels:** Display clear harmonic structure with visible formants
 - **Plosives:** Show burst patterns with brief silence followed by energy release
 - **Fricatives:** Exhibit high-frequency noise-like patterns
 - **Nasals:** Show low-frequency emphasis with formant structure
- Overlaying phoneme boundaries on the original waveform aids in understanding speech segmentation

5. Nature of Sound Source Classification:

- **Voiced sounds** (vowels, nasals, voiced consonants): Produced with vocal cord vibration, showing periodic waveforms
- **Unvoiced sounds** (fricatives like /s/, /f/, plosives like /p/, /t/, /k/): Produced without vocal cord vibration, showing aperiodic noise patterns

- **Plosives:** Characterized by stop closure followed by burst release
- **Fricatives:** Show continuous airflow causing turbulent noise
- **Nasals:** Feature nasal resonance with energy concentration in lower frequencies

6. Practical Applications:

- **Phonetic analysis:** Detailed study of pronunciation and articulatory features
- **Speech synthesis:** Understanding phoneme characteristics for text-to-speech systems
- **Language learning:** Visualizing phoneme differences for pronunciation training
- **Speech recognition:** Feature extraction for improved ASR systems
- **Speech therapy:** Analyzing and correcting pronunciation issues
- **Linguistic research:** Studying phonetic variations across speakers and dialects

1.13.2 Observations:

- Deep learning models like Wav2Vec2 provide robust phoneme recognition without manual feature engineering
- Different phoneme types exhibit distinct spectrotemporal characteristics
- Time-aligned phoneme extraction enables detailed acoustic analysis of speech units
- Spectrogram visualization is crucial for understanding the acoustic properties of phonemes
- Combination of waveform and spectrogram provides comprehensive phoneme characterization

1.13.3 Technical Insights:

- The Wav2Vec2 model leverages self-supervised learning on large speech corpora
- CTC (Connectionist Temporal Classification) layer enables sequence-to-sequence mapping
- Model outputs probability distributions over phoneme classes at each time step
- Phoneme boundaries are inferred from transitions in predicted phoneme sequences
- Saved phoneme segments can be used for further analysis or building phoneme databases

1.13.4 Conclusion:

This assignment successfully demonstrated phoneme extraction and visualization from continuous speech using a state-of-the-art deep learning model. The experiment revealed how modern speech processing techniques can automatically segment and classify phonemes with high accuracy. The Wav2Vec2 model effectively recognized phonemes and provided temporal alignment, enabling detailed acoustic analysis. Visualization through waveforms and spectrograms helped understand the distinct characteristics of different phoneme types - voiced vs. unvoiced, vowels vs. consonants, and various manners of articulation. These techniques form the foundation for advanced speech processing applications including automatic speech recognition, speech synthesis, and phonetic analysis systems.