

Lab Evaluation 1 for
UCS749 : Speech Processing and Synthesis

SUBMITTED BY

Aiman Gupta

102103488

PROJECT TITLE

Recognise My Voice Commands

COURSE COORDINATOR

Raghav B. Venkataramaiyer



THAPAR INSTITUTE
OF ENGINEERING & TECHNOLOGY
(Deemed to be University)

THAPAR INSTITUTE OF ENGINEERING AND TECHNOLOGY,
PATIALA

SESSION: July-December, 2024

DATE OF SUBMISSION: 21/03/2024

Summary of the paper

The paper introduces the "Speech Commands Dataset," a collection of spoken words for training keyword spotting systems. Since it is intended for low-resource devices, effective models must be able to understand commands such as "Yes," "No," and "Stop." The dataset, which has baseline models with an accuracy of 88.2%, aids in improving model comparability and reproducibility. More than 100,000 words from 2,618 speakers make up this collection.

Dataset

The Speech Commands Dataset described in the research paper contains over 105,829 audio recordings of 35 different words. The recordings are in WAV format, sampled at 16 kHz, and each audio file contains a single spoken word. The dataset was collected from 2,618 speakers, ensuring a wide variety of pronunciations and accents. It primarily focuses on small-vocabulary keyword spotting tasks, with words like "Yes," "No," "Up," "Down," "Left," "Right," "On," "Off," "Stop," "Go," and additional words such as digits (zero to nine).

The entire dataset is approximately 3.8 GB when uncompressed, or 2.7 GB when stored as a gzip-compressed tar archive. Additionally, it includes background noise files to simulate real-world conditions and improve model robustness.

Key details:

- Words: 35 target words, including digits, commands, and auxiliary words.
- Recordings: 105,829 utterances.
- Speakers: 2,618 speakers, each assigned a unique identifier.
- File Format: 1-second WAV files, 16-bit PCM, 16 kHz sample rate.
- Size: ~3.8 GB uncompressed, ~2.7 GB compressed.

This dataset is designed to train and evaluate models for on-device keyword recognition in low-resource environments.

Snapshots :

Loading Dataset

```
import rarfile
import os

def extract_rar(rar_path, extract_to):
    # Ensure the extract directory exists
    if not os.path.exists(extract_to):
        os.makedirs(extract_to)

    # Open and extract the .rar file
    with rarfile.RarFile(rar_path) as rf:
        rf.extractall(path=extract_to)

    print(f"Extracted to {extract_to}")

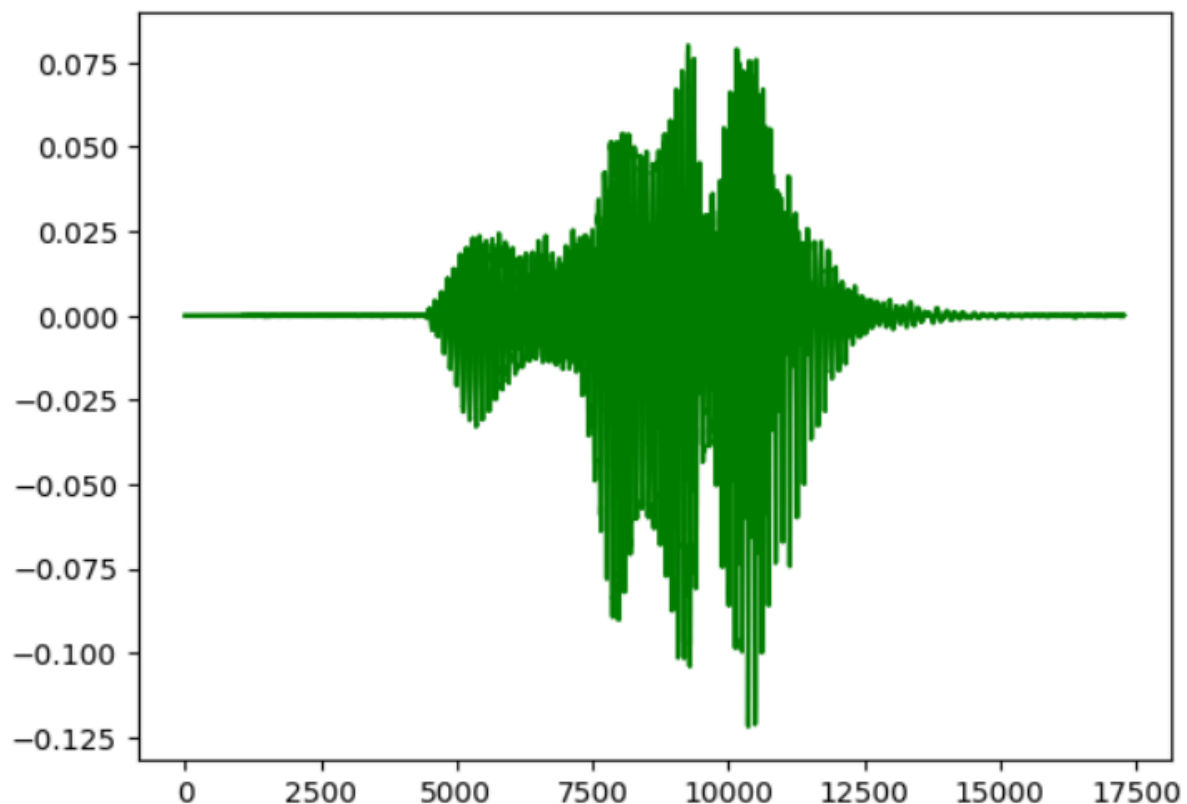
# Path to your .rar file and where to extract
rar_file_path = "/content/recordings_v2.rar"
extract_to_path = "/content"

# Extract the .rar file
extract_rar(rar_file_path, extract_to_path)
```

Extracted to /content

Waveform

Shape of waveform: torch.Size([2, 17280])
Sample rate of waveform: 16000
Label of waveform: zero



Labels

```
▶ labels = sorted(list(set(datapoint[2] for datapoint in train_set)))  
labels
```

```
↔ ['backward',  
   'bed',  
   'bird',  
   'cat',  
   'dog',  
   'down',  
   'eight',  
   'five',  
   'follow',  
   'forward',  
   'four',  
   'go',  
   'happy',  
   'house',  
   'learn',  
   'left',  
   'marvin',  
   'nine',  
   'no',  
   'off',  
   'on',  
   'one',  
   'right',  
   'seven',  
   'sheila',  
   'six',  
   'stop',  
   'three',  
   'tree',  
   'two',  
   'up',  
   'yacht']
```

Testing and training epochs

```
Train Epoch: 16 [0/700 (0%)] Loss: 0.001486  
75%| 15.079545454545755/20 [01:05<00:20, 4.09s/it] Train Epoch: 16 [50/700 (7%)] Loss: 0.225093  
76%| 15.147727272727575/20 [01:05<00:23, 4.94s/it] Train Epoch: 16 [100/700 (14%)] Loss: 0.054471  
76%| 15.215909090909395/20 [01:05<00:24, 5.18s/it] Train Epoch: 16 [150/700 (21%)] Loss: 0.057191  
76%| 15.261363636363942/20 [01:06<00:25, 5.34s/it] Train Epoch: 16 [200/700 (29%)] Loss: 0.146788  
77%| 15.329545454545762/20 [01:06<00:26, 5.61s/it] Train Epoch: 16 [250/700 (36%)] Loss: 0.002141  
77%| 15.375000000000309/20 [01:06<00:26, 5.69s/it] Train Epoch: 16 [300/700 (43%)] Loss: 0.164154  
77%| 15.420454545454856/20 [01:06<00:26, 5.75s/it] Train Epoch: 16 [350/700 (50%)] Loss: 0.038748  
77%| 15.488636363636676/20 [01:07<00:29, 6.56s/it] Train Epoch: 16 [400/700 (57%)] Loss: 0.013447  
78%| 15.556818181818496/20 [01:07<00:24, 5.59s/it] Train Epoch: 16 [450/700 (64%)] Loss: 0.091980  
78%| 15.61363636363668/20 [01:08<00:21, 4.99s/it] Train Epoch: 16 [500/700 (71%)] Loss: 0.499241  
78%| 15.6818181818185/20 [01:08<00:19, 4.51s/it] Train Epoch: 16 [550/700 (79%)] Loss: 0.039321  
79%| 15.738636363636683/20 [01:08<00:18, 4.32s/it] Train Epoch: 16 [600/700 (86%)] Loss: 0.090431  
79%| 15.772727272727593/20 [01:08<00:18, 4.31s/it] Train Epoch: 16 [650/700 (93%)] Loss: 0.025524  
80%| 16.045454545454874/20 [01:09<00:14, 3.61s/it]  
Test Epoch: 16 Accuracy: 167/175 (95%)
```

Prediction

```
Recording started for 1 seconds.  
Recording ended.  
Predicted: learn.
```