# PROBLEM REPORT

UCS749: CONVERSATIONAL

AI: SPEECH PROCESSING AND

SYNTHESIS

Prepared by: Kartik Koul 4CO26, Roll Number 102103819

Submitted to: Dr B.V. Raghav Assistant Professor, CSED

#### ABOUT THE PAPER

The paper describes the development process of a unique dataset designed to further the training and evaluation of keyword spotting models. It recognises and states the requirements for such locally stored models that detect the triggers for the virtual assistants available on modern devices. The paper succinctly details the processes of data collection, cleaning, and evaluation.

### **ABOUT THE AUTHOR**

Pete Warden is an influential AI researcher and engineer known for his contributions to TensorFlow and TensorFlow Lite. He focuses on enabling machine learning on low-power devices, such as smartphones and IoT hardware, and is recognized for creating the SpeechCommands dataset for voice recognition tasks.

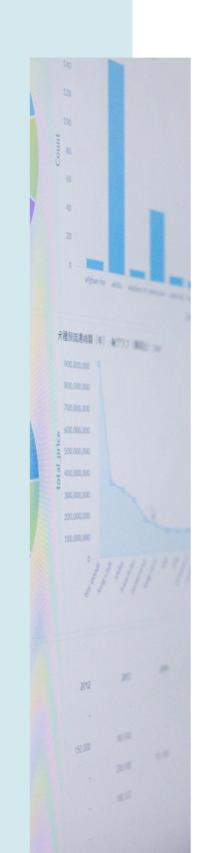
#### ABOUT THE DATASET

#### Following is a short note on the dataset itself:

The SpeechCommands dataset is a collection of audio recordings designed for training and evaluating machine learning models in simple voice recognition tasks. It consists of thousands of 1-second clips of 35 spoken words, such as "yes," "no," "stop," and digits. Created by Pete Warden, it is commonly used for keyword spotting and real-time voice control in low-power devices like smartphones and embedded systems. The dataset aids in developing models for recognizing commands in limited-resource environments.

#### PROBLEM STATEMENT

This evaluation involves a comprehensive study of the SpeechCommands dataset, starting with summarizing Pete Warden's paper that introduces the dataset. The dataset must then be downloaded and statistically analyzed to provide valuable insights. A classifier is to be trained to distinguish between the voice commands in the dataset, with performance evaluation against standard benchmarks. Additionally, 30 samples of each command must be recorded in one's voice to create a new dataset. The classifier will then be fine-tuned to adapt to the user's voice, with its performance being reported.



## A BRIEF ANALYSIS OF THE DATASET

- 1. Class Distribution: The dataset contains an uneven number of samples across different classes, which may require balancing techniques. For example, common words like "yes" and "no" tend to have more examples than less frequent words.
- 2. **Diverse Speakers**: The dataset includes recordings from a diverse group of speakers, which ensures robustness in recognizing voice commands across various accents, pitches, and speaking speeds. This is an important factor for generalizing the classifier to different users.
- 3. **Audio Characteristics**: Each audio clip is sampled at 16 kHz and is trimmed or padded to 1 second in length.
- 4. **Background Noise**: The dataset includes specific "background noise" files, which help models become resilient to real-world noise interference.
- 5. **Speaker IDs**: Each recording includes metadata about the speaker (user id), which allows for speaker-dependent or independent modeling.
- 6. Statistical Features:

## APPROACH TAKEN FOR MODEL DEVELOPMENT:

This code uses a Convolutional Neural Network (CNN) to implement a keyword recognition system. It starts by loading audio files from subfolders, where each subfolder name represents a class label. The audio is preprocessed to ensure uniform length and normalized amplitude. The data is then split into training and validation sets.

The CNN architecture consists of three convolutional layers with increasing filter sizes (32, 64, 128), each followed by max pooling. This structure helps the model learn hierarchical features from the raw audio waveform. The flattened output feeds into a dense layer with dropout for regularization before the final classification layer.

The model is trained using the Adam optimizer and sparse categorical cross-entropy loss. It runs for 10 epochs, which should be enough to see if model is learning effectively without overfitting. After training, the model is saved for future use. This approach is straightforward and effective for audio classification tasks, though one might want experiment with to hyperparameters or model architecture optimize performance for their specific dataset.

#### SUBMISSIONS:

- 1. A colaboratory file detailing the preprocessing of the Speech Commands dataset and the model definition and execution.
- 2. A colaboratory file with a rudimentary statistical analysis of the dataset.
- 3. The project report.
- 4. The developed model.

#### **FURTHER WORK:**

The model submitted has been developed on the SpeechCommands dataset only as of now. I am recording audio for my own dataset, but have been unable to complete it thus far. Once complete, further fine tuning will be carried as per the tasks detailed.

The checksum also has to be incorporated in the demo notebook with the self-developed dataset.