

# **LAB EVAL: Conversational AI - Speech Processing and Synthesis**

**Simran Minhas - 102117020**

## **Summary of Research Paper**

The Speech Commands dataset is a crucial resource designed for building and assessing keyword spotting systems. It offers a consistent framework for developing models that identify specific spoken words from a predefined set of keywords, ensuring a uniform basis for evaluating and comparing different models.

## **Key Findings**

1. Standardized Comparison:
  - This dataset provides a standardized approach for comparing keyword spotting models, which is essential for tracking progress and evaluating model performance effectively.
2. Keyword Spotting Focus:
  - The dataset is specifically tailored for keyword spotting, focusing on distinguishing between audio clips with target keywords and those without. This specialization aids in creating models that perform well in practical scenarios, minimizing false positives.
3. Enhanced Accuracy in Version 2:
  - The second version of the dataset has led to better model accuracy. Models trained with Version 2 data achieved a Top-One accuracy of 88.2% on the training set and 89.7% on the Version 1 test set, reflecting improved model performance.
4. Comparison with Other Datasets:
  - While datasets such as Mozilla's Common Voice are valuable for general speech tasks, they are not optimized for keyword spotting. The Speech Commands dataset addresses this need specifically, making it ideal for keyword spotting tasks.
5. Training and Evaluation:
  - The dataset supports the training and evaluation of various models. The advancements seen with Version 2 provide a more effective benchmark for evaluating keyword spotting performance.

## **Data Collection and Processing**

- Collection Method:
  - The dataset was gathered using an open-source web application that utilized the Web Audio API for recording. This method was chosen for its simplicity, aiming to reduce participant dropout and ensure comprehensive data collection.
- Data Processing:
  - Each recording is a one-second WAV file. Post-processing involves filtering out quiet or silent segments and extracting the loudest part of each recording to enhance data quality and relevance.

## **Model Training and Evaluation**

- Model Training:
  - The default convolutional model from TensorFlow's tutorial was trained using both Version 1 and Version 2 data. This provided a baseline for performance evaluation.
- Evaluation Metrics:
  - Models were tested using the datasets from both versions to evaluate their ability to distinguish between clips containing target keywords and those that do not.

## **Conclusions**

1. Dataset Utility:
  - The Speech Commands dataset is highly valuable for training and evaluating keyword spotting models. Its structured approach facilitates meaningful comparisons and progress in model development.
2. Improvements with Version 2:
  - The second version of the dataset shows notable improvements over the original, enhancing model accuracy and reliability.
3. Benchmarking:
  - This dataset offers a solid benchmark for comparing different models, advancing the field of keyword spotting by ensuring consistent data testing.

In summary, the Speech Commands dataset is an essential resource for developing and refining keyword spotting systems. It offers significant improvements in model performance and provides a standardized platform for evaluation and comparison.



