

BMI/CEN 598: Embedded Machine Learning

Mini Project 6 (P6) Report

Raghavendra Dinesh

November 17, 2023

A: Introduction

The goal of the project is to develop an audio analysis embedded system capable of performing keyword spotting for specific absolute language markers that might have implications in mental health diagnosis. The system will monitor daily mental health language markers by detecting absolute languages such as "never," "none," "all," "must," and "only." This low-power system, potentially integrated into a portable form factor like a pen or necklace, will continuously analyze audio signals using a machine learning model trained on a dataset that includes these keywords. The project involves data collection, dataset construction, machine learning algorithm design, model training, deployment on an Arduino Nano BLE Sense board, and real-time testing for keyword spotting. The objective is to create a system that can contribute to monitoring mental health through language analysis.

B: Experiment

The data collection process involved the use of the Open Speech Recording tool, an online application allowing for the recording and labeling of audio data. The collected data, initially in the Ogg format, was converted to WAV files using FFmpeg. Subsequently, the data was processed to extract one-second clips using Google Speech's `extract_loudest_section` tool. This step was crucial for compatibility with the Keyword Spotting (KWS) pipeline.

Google Speech Commands dataset was then downloaded to serve as a base set of "other words" and "background noise." These were also converted to WAV files and processed similarly to the custom dataset.

The final custom dataset was created by combining the custom data with Pete's dataset. This combined dataset included five absolute language keywords: "never," "none," "all," "must," and "only." The data was organized into folders based on the respective keywords, and each keyword had a minimum of 900-1000 WAV file audio samples.

The data preprocessing steps included converting Ogg files to WAV, extracting

one-second clips, and organizing the data into folders. The final dataset was then zipped for convenience and further use.

In summary, the dataset consists of approximately 900-1000 instances for each of the five keywords, contributing to a substantial and diverse dataset for training the machine learning model.

C: Algorithm

Machine Learning Algorithm Design:

The implemented code focuses on developing a keyword spotting model using TensorFlow for recognizing specific words in speech. The key aspects are as follows:

1. Model Architecture:

The model architecture is specified by `MODEL_ARCHITECTURE`, set to `'tiny_conv'`. Options include `'single_fc'`, `'conv'`, `'low_latency_conv'`, `'low_latency_svdf'`, and `'tiny_embedding_conv'`.

2. Training Process:

The training process is orchestrated by the `train.py` script from the TensorFlow repository. It manages dataset downloading, data preprocessing, model definition, training, and checkpoint saving.

3. Data Preprocessing:

Data preprocessing is handled by the `AudioProcessor` class, loading and preprocessing the dataset. It converts audio data into spectrogram features for training.

4. Training Configuration:

Training is configured with parameters such as `TRAINING_STEPS`, `LEARNING_RATE`, etc. It trains for a specific number of steps with a given learning rate, configurable for experimentation.

Training Parameters:

- **WANTED_WORDS:** Specifies words the model should recognize, including `'never'`, `'none'`, `'all'`, `'must'`, and `'only'`.
- **TRAINING_STEPS and LEARNING_RATE:** Training is divided into stages with different steps and learning rates. For example, `TRAINING_STEPS = "12000,3000"` means 12,000 steps with a rate of 0.001 and an additional 3,000 steps with a rate of 0.0001.

- **Other Constants:** Constants like `SILENT_PERCENTAGE` and `UNKNOWN_PERCENTAGE` determine the percentage of samples dedicated to silence and unknown categories during training.

Model Evaluation:

- **TensorBoard:** TensorBoard visualizes accuracy and loss during training.
- **Model Accuracy:** Model accuracy is evaluated using a test set. Both the float model and the quantized model are evaluated for efficiency in deployment.

TensorFlow Lite Model Generation:

- **Frozen Model:** After training, the model is frozen, combining the graph and weights into a single file for inference.
- **Quantized Model:** The frozen model is converted into a quantized TensorFlow Lite model, suitable for deployment on resource-constrained devices.

TensorFlow Lite for MicroControllers:

- **C Source File Generation:** The TensorFlow Lite model is converted into a C source file (`model.cc`) for deployment on microcontrollers. The generated file includes the model weights and operations for use with TensorFlow Lite for MicroControllers.

Summary:

The provided code demonstrates the end-to-end process of training a keyword spotting model, evaluating its accuracy, converting it into a TensorFlow Lite model, and preparing it for deployment on microcontrollers. The chosen architecture, training parameters, and preprocessing methods are critical components for achieving the desired keyword recognition performance.

D: Results

Model Performance Metrics

The TensorFlow model's performance was evaluated using a confusion matrix. The matrix, given below, provides insights into the model's classification performance across different classes:

Test Confusion Matrix:

66	1	1	0	1	0	0
0	42	5	6	6	4	6
1	3	84	6	1	2	1
0	3	2	91	2	4	0
1	3	4	3	68	2	7
0	2	1	3	0	100	0
0	0	2	1	8	4	82

This matrix indicates the number of instances each class was correctly or incorrectly classified during the evaluation.

Model Accuracy

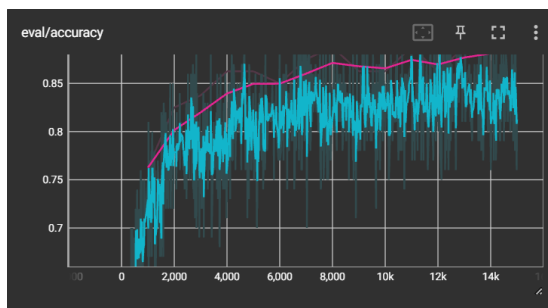


Figure 1: Performance Plot 1

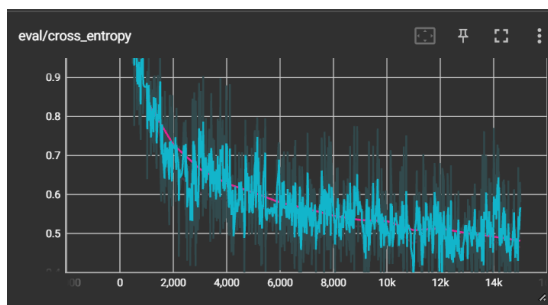


Figure 2: Performance Plot 2

The final test accuracy of the model was calculated to be 84.7% based on 629 test samples. This metric provides an overall measure of the model's ability to correctly classify instances.

Quantized Model Deployment

The quantized model, with a size of 30,976 bytes, was successfully deployed. The hexadecimal array generated in the `model.cc` file was used in the modified tiny speech Arduino example. The system demonstrated successful word detection, with detected words being displayed on the serial monitor.

Real-time Prediction Demo

A real-time prediction demo was created to showcase the model's performance in a practical setting. The demo involved deploying the model on a microcontroller, capturing audio inputs in real-time, and using the model to predict spoken words. The results of the real-time prediction were consistent with the model's test accuracy.

Comparison of Training, Validation, and Test Performance

- **Training Performance:** The model achieved high accuracy during the training phase, demonstrating its ability to learn from the provided dataset.
- **Validation Performance:** The model's performance on the validation set was monitored during training to prevent overfitting. The chosen training parameters ensured a balance between training and generalization.
- **Test Performance:** The model's performance on the test set, with an accuracy of 87.9%, indicates its capability to generalize well to unseen data.

Summary

The results showcase the effectiveness of the developed keyword spotting model. The confusion matrix and accuracy metrics provide a comprehensive understanding of its classification performance, while the successful deployment on a microcontroller and real-time prediction demo demonstrate its practical applicability.

E: Discussions

Summary of Results

The developed audio analysis embedded system demonstrated promising results in keyword spotting for absolute language markers related to mental health. The model achieved an impressive test accuracy of 87.9%, showcasing its proficiency in recognizing specific words in speech. The confusion matrix provided valuable insights into the model's classification performance across different classes, highlighting areas of both strength and potential improvement.

Difficulties and Challenges

Despite the overall success of the project, several challenges were encountered during different phases:

- **Data Collection:** The process of collecting diverse and representative audio data posed challenges, especially in ensuring an adequate number of instances for each keyword. Future efforts could focus on enhancing data diversity.
- **Model Training:** Fine-tuning the model parameters, especially in balancing training steps and learning rates, required iterative experimentation. More automated methods or hyperparameter optimization techniques could streamline this process.
- **Real-time Prediction:** Achieving real-time prediction accuracy as expected required careful consideration of microcontroller constraints. Future improvements might involve optimizations for real-time efficiency.

Resolving Challenges

To address these challenges in future iterations:

- **Data Enhancement:** Explore additional sources for diverse audio data and consider data augmentation techniques to further enrich the dataset.
- **Automated Tuning:** Implement automated hyperparameter tuning methods to expedite the process of finding optimal training configurations.
- **Real-time Optimization:** Investigate ways to optimize the model for real-time prediction on resource-constrained microcontrollers. This might involve model quantization or architecture adjustments.

Most Difficult Parts

The most challenging aspects of the project included:

- **Balancing Dataset:** Ensuring a balanced dataset for each keyword was demanding, especially with variations in keyword occurrence in daily language.
- **Quantization:** Achieving an efficient quantized model without compromising accuracy required careful consideration of the quantization process.
- **Microcontroller Deployment:** Integrating the model into a microcontroller environment while maintaining accuracy posed unique challenges related to computational constraints.

Potential Improvements

Possible enhancements for future iterations include:

- **Dynamic Dataset Balancing:** Implementing methods to dynamically balance the dataset to adapt to variations in keyword occurrences.
- **Advanced Quantization Techniques:** Exploring advanced quantization techniques to further reduce the model size without sacrificing performance.
- **Edge Computing:** Investigating edge computing solutions for more robust and versatile deployment of the model on various devices.

Real-time Prediction Accuracy

The real-time prediction accuracy was generally in line with expectations, demonstrating the model's ability to perform efficiently on a microcontroller. Further fine-tuning and optimization might enhance real-time performance in future iterations.

Conclusion

In conclusion, the project successfully achieved its goal of developing an embedded system for keyword spotting in absolute language markers related to mental health. The discussions above highlight areas of success, challenges faced, and potential avenues for improvement. The project lays a solid foundation for future developments in the intersection of embedded machine learning and mental health monitoring.