
Learning Disentangled Representations for Audio Classification

Nehemiah Skandera
UC San Diego
nskandera@ucsd.edu

Abstract

1 Much of this CSE 192 course has been centered around how we can manipulate
2 and visualize the unique features of audio to develop our understanding of how
3 audio signals behave. Thus, we see that digital audio manipulation can be a chal-
4 lenging task - especially as it relates to Machine Learning. Furthermore, we can
5 generalize that the task of classifying audio is especially challenging. According
6 to the information I have found in my limited research, typical audio classification
7 methods struggle with details such as feature extraction, and the disentangling of
8 more pronounced audio frequencies. This project has been a journey of taking
9 steps to explore the combination of Convolutional Neural Networks (CNN) with
10 Support Vector Machines (SVM), and taking ideas from other papers to improve
11 the model - such as generating disentangled representations of data. The objective
12 of this project is to develop a model that learns disentangled versions of content,
13 rhythm, and pitch from audio signals in order to integrate the corresponding rep-
14 resentation into a hybrid CNN-SVM framework that attempts to enhance/discover
15 the performance of classification mechanisms.

16 1 Introduction

17 The classification of audio signals is fundamental to many real-world applications. For instance, we
18 use audio classification in speech recognition, oceanic research, and music recognition. However, we
19 often see the limits of these systems of classification; often, complexity is overlooked. What makes
20 this project (which was mostly a process of discovery) unique is that it introduces a novel approach
21 that combines CNNs and SVMs and leverages disentangled signal representations in attempting to
22 improve the performance of audio classification.

23 2 Related Works

24 2.1 Learning Disentangled Representations for Timber and Pitch in Music Audio

25 Hung et al. [2018] Disentangling timber and pitch in audio signals forms the basis for understanding
26 the importance of disentangled features in audio classification.

27 2.2 Unsupervised Speech Decomposition via Triple Information Bottleneck

28 Qian et al. [2021] This paper investigates the use of information bottlenecks (hence the name) for
29 unsupervised speech decomposition. Thus, providing information on how we can disentangle audio
30 representations and enhance the extraction of audio features.

3 Data Description

The dataset used in this project is gathered from room impulse responses (RIRs) in addition to noise sources from various databases. All of the data was gathered with a 16k sampling rate and 16-bit precision. Here is a description of the folders and data sources based on the data's README:

3.1 Real RIRs and Isotropic Noises

This set of real RIRs includes three databases:

- **RWCP Sound Scene Database**
- **2014 REVERB Challenge Database**
- **Aachen Impulse Response Database (AIR)**

Overall, there are 325 real RIRs. The isotropic noises available in these databases are used in conjunction with their associated RIRs. The data can be downloaded from the following links:

- AIR: http://www.openslr.org/resources/20/air_database_release_1_4.zip
- RWCP: <http://www.openslr.org/resources/13/RWCP.tar.gz>
- 2014 REVERB Challenge: http://reverb2014.dereverberation.com/tools/reverb_tools_for_Generate_mcTrainData.tgz
http://reverb2014.dereverberation.com/tools/reverb_tools_for_Generate_SimData.tgz

3.2 Simulated RIRs

This folder contains simulated RIRs. Details regarding this dataset can be found in the following '*simulated_rirs/README*' file. The simulated RIR dataset can be downloaded here:

- http://www.openslr.org/resources/26/sim_rir.zip

4 Architecture

4.1 System Overview

The proposed system for my project integrates feature extraction using Fourier Transforms and creates a disentangled representation for learning through CNNs. Classification is then accomplished using a hybrid CNN-SVM model.

4.2 Components

- **Feature Extraction:** Fourier Transforms are used to extract frequency-domain features from audio signals.
- **Disentangled Representations:** The CNN-based encoder functions by extracting the content, rhythm, and pitch from audio signals.
- **Hybrid Model:** The CNN decoder functions by integrating the disentangled features. Meanwhile, the SVM classifier is used for further classification of the Fourier-transformed features.

4.3 Implementation Details

The system uses Python libraries such as 'joblib' for faster computing times, TensorFlow, scikit-learn, and scipy for audio processing. Key components of the system include data preprocessing, feature extraction, model training, and model analysis.

69 5 Experiments

70 5.1 Experimental Setup

71 The dataset used for the experiments is described in Section 3. The dataset uses both real and
72 simulated RIRs, in addition to isotropic noises. The dataset was divided into training and testing
73 subsets, and the system’s performance was evaluated based on classification accuracy.

74 5.2 Methodology

75 The methodology can be broken into steps:

- 76 • **Step 1:** Download and extract the dataset
- 77 • **Step 2:** Load dataset RIR data
- 78 • **Step 3:** Extract Fourier transforms as features
- 79 • **Step 4:** Apply PCA for dimensionality reduction
- 80 • **Step 5:** Data Augmentation
- 81 • **Step 6:** Create Encoders
- 82 • **Step 7:** Create CNN Decoder
- 83 • **Step 8:** SVM Kernel System
- 84 • **Step 9:** Bringing everything together
- 85 • **Step 10:** System Evaluation

86 5.3 Evaluation Metrics

87 Accuracy is used as a metric to evaluate the performance of the CNN, SVM, and combined models.

88 6 Results

89 6.1 Results Presentation

- 90 • **SVM Accuracy:** The SVM classifier achieved an accuracy of 92% on the test set.
- 91 • **CNN Accuracy:** The CNN model achieved an accuracy of 61% on the test set.
- 92 • **Combined Model Accuracy:** The combined model achieved an accuracy of 85%.

93 6.2 Confusion Matrix

True \ Predicted	RVB2014	RWCP	air
RVB2014	22	5	0
RWCP	0	40	0
air	0	8	9

Table 1: Confusion matrix showing the number of true vs. predicted classifications for each class.

94 **6.3 Classification Report**

Class	Precision	Recall	F1-score	Support
RVB2014	1.00	0.81	0.90	27
RWCP	0.75	1.00	0.86	40
air	1.00	0.53	0.69	17
Accuracy	0.85			
Macro avg	0.92	0.78	0.82	84
Weighted avg	0.88	0.85	0.84	84

Table 2: Classification report showing precision, recall, F1-score, and support for each class.

95 **6.4 Plots**

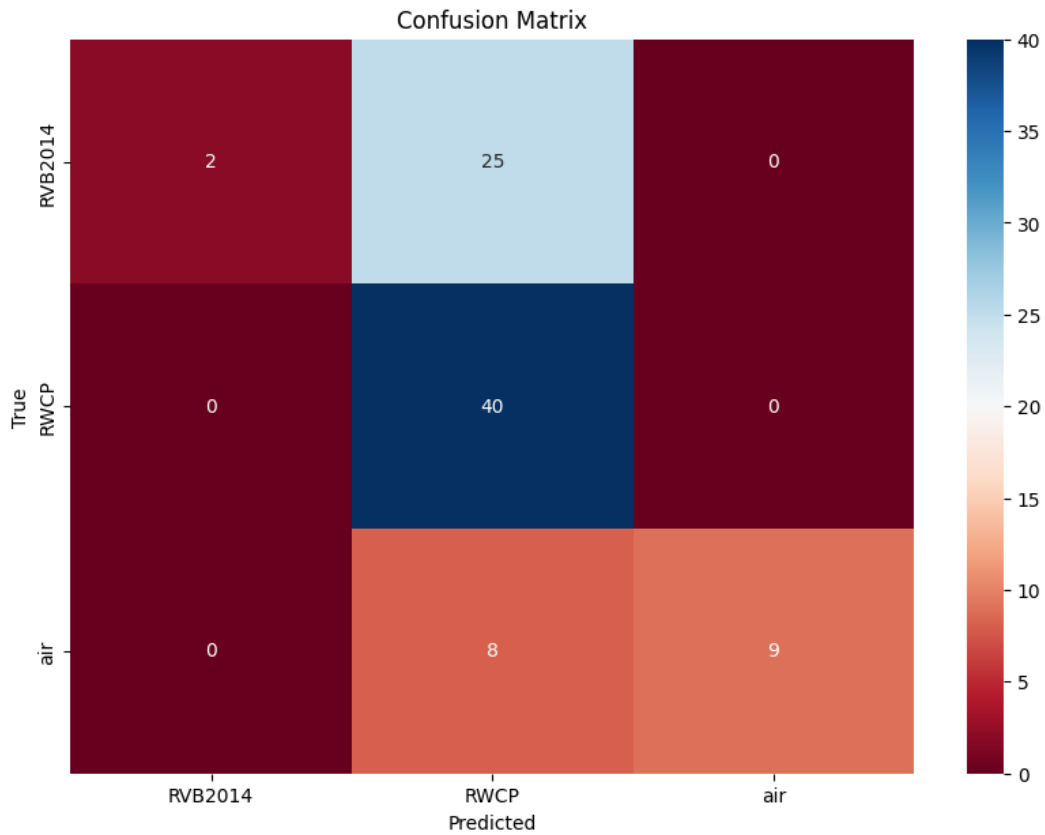


Figure 1: Confusion matrix heatmap showing true vs. predicted classifications.

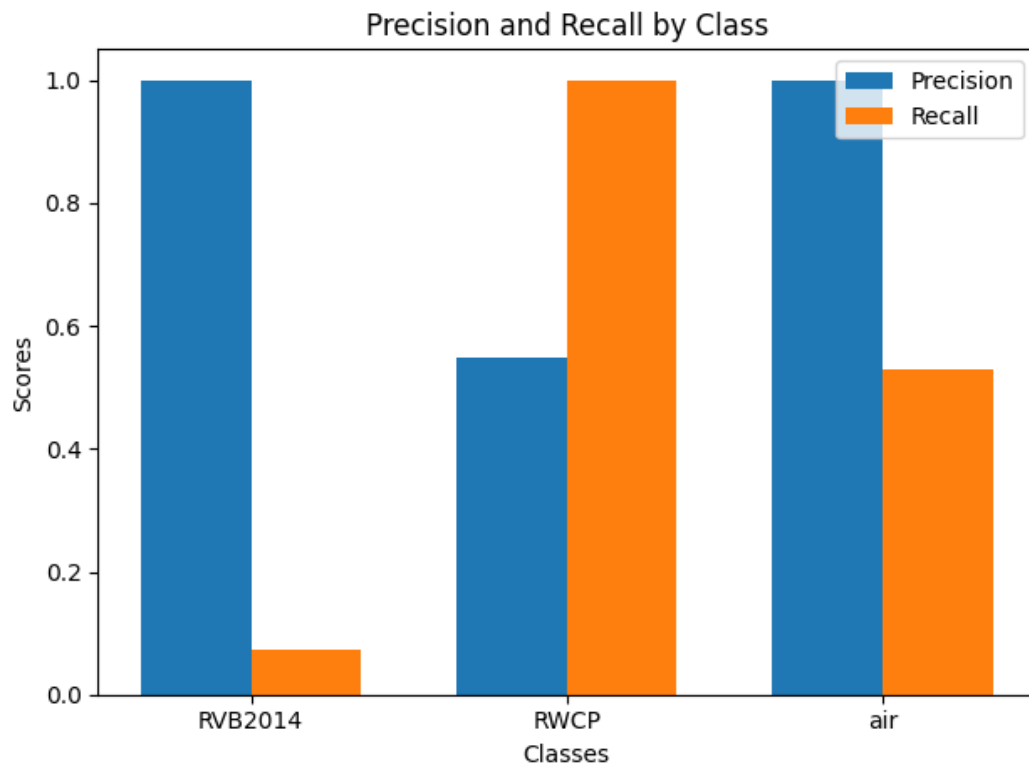


Figure 2: Precision and recall bar chart by class.

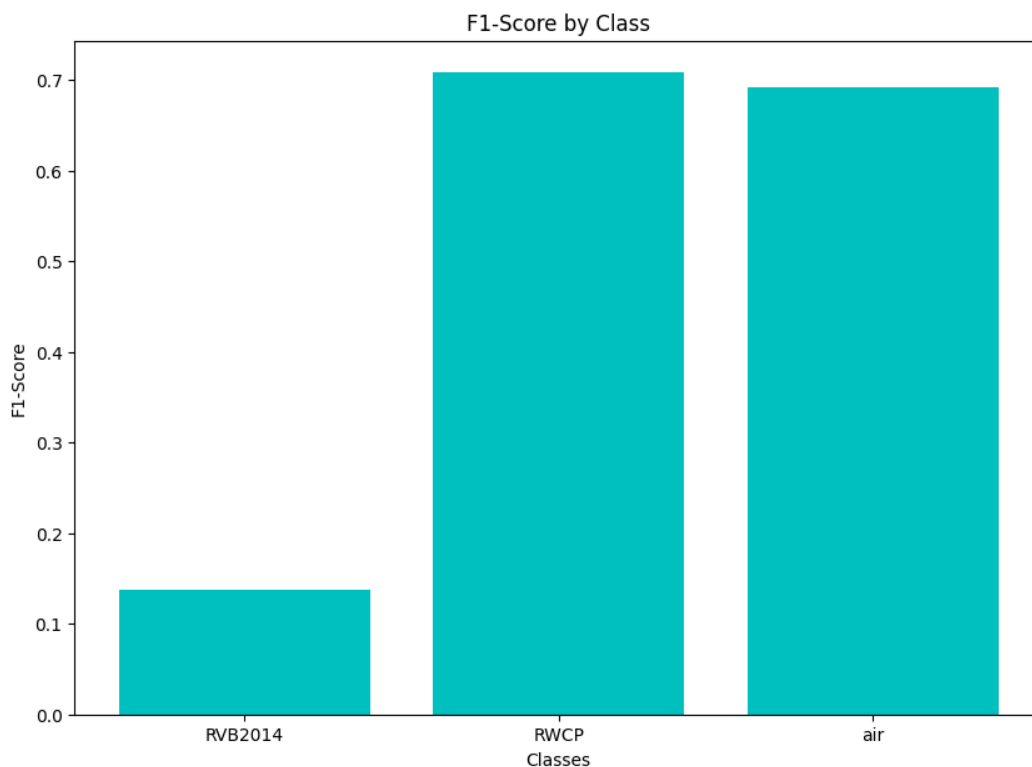


Figure 3: F1-score bar chart by class.

7 Conclusion

For the most part, the project was a process. I started with the idea of a simple CNN-SVM hybrid audio classification system. Over time, the project evolved. First, I was able to learn some new techniques for audio processing from a couple of papers, then I had to go through the process of making my models work cohesively together. This was the longest process. Finally, I had to find a way to train the model - given the computing resources I had. Initially, I was hoping to gain access to the San Diego Supercomputer (SDSC). However, this proved to not be an option. Thus, I learned several ways to optimize the resources that I had to perform computations. Thus, we arrive at the end product. I believe my project demonstrates the potential effectiveness of combining CNNs with SVMs for audio classification of disentangled signals. Overall, I was pleasantly surprised by the overall accuracy result of 85%. Future work could be done in exploring additional feature extraction methods and more complex hybrid models. Maybe this could eventually lead to instrument recognition?

8 Supplementary Materials

8.1 Code

The complete code for this project is available at https://github.com/NehemiahSkan/CSE_192_Project.

References

- Yun-Ning Hung, Yi-An Chen, and Yi-Hsuan Yang. Learning disentangled representations for timber and pitch in music audio. *Research Center for IT Innovation, Academia Sinica, Taiwan; KKBOX Inc., Taiwan*, 2018.

117 Kaizhi Qian, Yang Zhang, Shiyu Chang, David Cox, and Mark Hasegawa-Johnson. Unsuper-
118 vised speech decomposition via triple information bottleneck. [https://arxiv.org/abs/2004.](https://arxiv.org/abs/2004.11284)
119 11284, 2021.