

RADBOUD UNIVERSITY NIJMEGEN



FACULTY OF SCIENCE

Convolutional Neural Networks applied to Keyword Spotting using Transfer Learning

THESIS IN AUTOMATIC SPEECH RECOGNITION
(LET-REMA-LCEX10)

Author:

Christoph SCHMIDL
s4226887
c.schmidl@student.ru.nl

Supervisor:

dr. L.F.M. TEN BOSCH

August 17, 2019

Contents

1	Introduction	2
1.1	Literature review	2
1.1.1	Raw waveform-based audio classification using sample-level CNN architectures	2
1.1.2	Transferable deep features for keyword spotting	2
1.1.3	Imagenet: A large-scale hierarchical image database	2
1.1.4	Imagenet classification with deep convolutional neural networks .	2
1.1.5	Speech Recognition: Keyword Spotting Through Image Recognition.	2
1.1.6	Convolutional neural networks for small-footprint keyword spotting	4
1.1.7	Small-footprint keyword spotting using deep neural networks . .	5
1.1.8	Speech Commands: A Dataset for Limited-Vocabulary Speech Recognition	5
2	Method	6
3	Set-up	7
3.1	Dataset	7
3.2	Preprocessing	8
3.2.1	Load Data	8
3.2.2	Check Wav Length	9
3.2.3	Create Spectrograms	9
3.2.4	Create Log Spectrograms	11
3.3	Models	13
3.3.1	Baseline models	13
3.3.2	CNN models	13
4	Experiments	14
4.1	Baseline	14
4.2	CNNs	14
5	Analysis and Results	14
5.1	Baseline	14
5.2	CNNs	15
5.2.1	Xavier initialization	15
5.2.2	Imagenet weight initialization	16
6	Discussion	18
7	Conclusion	18
8	Future Work	18
9	References	18
10	Appendix	18

1 Introduction

The task of keyword spotting (KWS) is interesting to different domains where a hands-free interaction experience is required or desired like Google’s feature of interacting with mobile devices (include ”OK Google” reference).

Different approaches to keyword spotting like:

- Deep Neural Networks (DNNs)
- Convolutional Neural Networks (CNNs)
- (Keyword/Filler) Hidden Markov Models (HMMs)

1. Problem
2. Background (literature overview)
3. Research Question, Hypotheses, intro to experiment

1.1 Literature review

This section contains the most prominent approaches to the KWS task which have been successfully applied in the past and serve as baseline models or inspirations for the proposed model in this thesis.

1.1.1 Raw waveform-based audio classification using sample-level CNN architectures

- Raw waveform-based audio classification using sample-level CNN architectures [14]

1.1.2 Transferable deep features for keyword spotting

- Transferable deep features for keyword spotting [16]

1.1.3 Imagenet: A large-scale hierarchical image database

- Imagenet: A large-scale hierarchical image database [6]

1.1.4 Imagenet classification with deep convolutional neural networks

- Imagenet classification with deep convolutional neural networks [12]

1.1.5 Speech Recognition: Keyword Spotting Through Image Recognition.

The authors of the paper ”Speech Recognition: Keyword Spotting Through Image Recognition” [11] transformed the KWS task which incorporates audio data into the domain of image classification. They used the Speech Commands Dataset [21] which contains spoken words of the length of one second in order to train and evaluate their model. According to [21], the Speech Commands Dataset V2 [10] comprises one-second audio clips which were sampled at 16KHz and containing ten words, namely ”Yes”, ”No”, ”Up”, ”Down”, ”Left”, ”Right”, ”On”, ”Off”, ”Stop”, and ”Go”, and have one additional special label for ”Unknown Word”, and another for ”Silence” (no speech detected). A vector representation of these one-second audio clips would therefore be of the form \mathbb{R}^{16000} .

The authors used three different models, namely:

- A Low Latency Convolutional Neural Network which is designed to reduce its memory footprint by limiting the number of model parameters. This model is similar to the model called "cnn-one-fstride4" which is used in [17] but differs in terms of filter size, stride, channels, dense and params. The model has been trained using Stochastic Gradient Descent and Xavier Initialization has been used in order to initialize the model weights.
- The MNIST TensorFlow CNN / Basic CNN where some tweaks have been performed to the first layer in order to fix dimension mismatches. A baseline architecture is described in [17] (3 module setup?).
 - <https://github.com/tensorflow/docs/blob/master/site/en/tutorials/estimators/cnn.ipynb>
 - <https://github.com/tensorflow/tensorflow/tree/master/tensorflow/examples/tutorials/mnist>
- Adversarially trained CNN which is inspired by MCDNN [5] and AlexNet [12]. One shallow CNN which has been used for prototyping and hyperparameter tuning. Dropout was counter-productive and therefore Virtual Adversarial Training was used, inspired by [8].

Evaluated parameters in this paper:

- Adversarial Training Results and Comparison with Vanilla CNN
- increase in training and validation accuracy over the first ten epochs for the low-latency convolution and VAT models: a zoomed-in version of Figure 12
- decrease of costs over 500 epochs for the lowlatency convolution and VAT models: a zoomed-out version of Figure 13
- increase of training and validation accuracy over 500 epochs for the low-latency convolution and VAT models: a zoomed-out version of Figure 10
- reduction of cost over the first ten epochs for the low-latency convolution and VAT models: a zoomed-in version of Figure 11
- the evolution of training cross-entropy loss (blue and green) and validation accuracy (red and orange) compared between Xavier and truncated normal initialization; Xavier converges much faster and may attain better results
- the evolution of training cross-entropy loss (blue and green) and validation accuracy (red and orange) compared between Adam and SGD optimization; Adam converges faster than SGD but reaches the same results
- effect of the number of frequency-counting buckets on the accuracy of the low-latency convolution model. The model did not benefit from the increase in available data caused by increasing the number of buckets.
- effect of the spectrogram window size on the accuracy of the low-latency convolution model. There is a local optimum, as there was for stride in Figure 18.
- effect of added background noise on the final accuracy of the low-latency convolution model. The horizontal axis is signal-noise ratio in linear units.

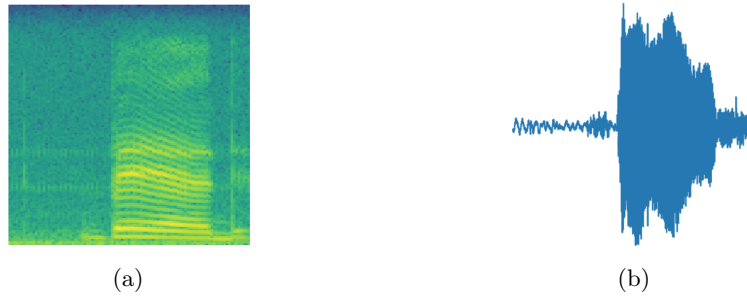


Figure 1: A comparison of the spectrogram (a) and the amplitude-vs.-time plot (b) for the same audio recording of a person saying the word “bed”.

- effect of the spectrogram window stride on the accuracy of the low-latency convolution model. For low values of stride, there is too much redundancy, while larger values result in lost information.

Conclusion

In this project we tackled the speech recognition problem by applying different CNN models on image data formed using log spectrograms of the audio clips. We also successfully implemented a regularization method “Virtual Adversarial Training” that achieved a maximum of 92% validation accuracy on 20% random sample of the input data. The significant work done in this project was the demonstration of how to convert a problem in audio recognition into the better-studied domain of image classification, where the powerful techniques of convolutional neural networks are fully developed. We also saw, particularly in the case of the low-latency convolution model, how crucial good hyperparameter tuning is to the accuracy of the model. A great number of hyperparameters must be tuned, including the many choices that go into network architecture, and any of the hyperparameters, poorly chosen, can make or break the overall performance of the model. Another contribution was the use of adversarial training to provide a regularization effect in audio recognition; this technique improved results relative even to well-established techniques such as dropout, and therefore has promising applications in the future.

Include summary about the approach of converting the long, one dimensional vector of audio data into a spectrograms and therefore making it a image classification problem.

1.1.6 Convolutional neural networks for small-footprint keyword spotting

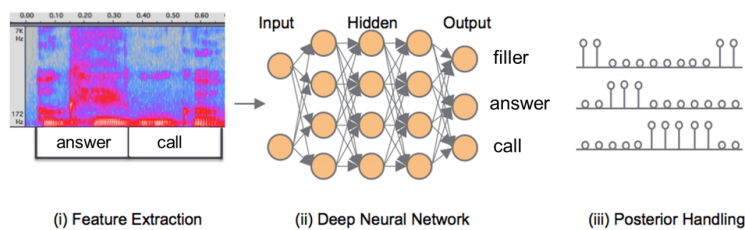


Figure 2: Framework of Deep KWS system, components from left to right: (i) Feature Extraction (ii) Deep Neural Network (iii) Posterior Handling

This framework originally comes from [4]. The only difference is the exchange of the DNN for a CNN.

- Convolutional neural networks for small-footprint keyword spotting [17]

Include summary about the different CNNs approaches which have been put into the 3 module framework of the below framework where the DNN has been exchanged for a CNN. How do the authors handle the long, one dimensional vector?

1.1.7 Small-footprint keyword spotting using deep neural networks

- Small-footprint keyword spotting using deep neural networks [4]

Include summary about the comparison between DNNs and HMMs and the general 3 module approach here: 1. Feature extraction. 2. Deep Neural Network 3. Posterior Handling. DNNs do not need a decoding algorithm like HMMs with Viterbi which makes it low latency.

1.1.8 Speech Commands: A Dataset for Limited-Vocabulary Speech Recognition

- Speech Commands: A Dataset for Limited-Vocabulary Speech Recognition [21]

Include summary and say why the Speech Commands Dataset is a good fit for this thesis. You probably do not need a Voice-activity detection (VAD) system here.

Top-One Error

The standard chosen for the TensorFlow speech commands example code is to look for the ten words "Yes", "No", "Up", "Down", "Left", "Right", "On", "Off", "Stop", and "Go", and have one additional special label for "Unknown Word", and another for "Silence" (no speech detected). The testing is then done by providing equal numbers of examples for each of the twelve categories, which means each class accounts for approximately 8.3% of the total. The "Unknown Word" category contains words randomly samples from classes that are part of the target set. The "Silence" category has one-second clips extracted randomly from the background noise audio files.

Applications

The TensorFlow tutorial gives a variety of baseline models, but one of the goals of the dataset is to enable the creation and comparison of a wide range of models on a lot of different platforms, and version one of has enabled some interesting applications. CMSISNN [21] covers a new optimized implementation of neural network operations for ARM microcontrollers, and uses Speech Commands to train and evaluate the results. Listening to the World [22] demonstrates how combining the dataset and UrbanSounds [23] can improve the noise tolerance of recognition models. Did you Hear That [24] uses the dataset to test adversarial attacks on voice interfaces. Deep Residual Learning for Small Footprint Keyword Spotting [19] shows how approaches learned from ResNet can produce more efficient and accurate models. Raw Waveformbased Audio Classification [14] investigates alternatives to traditional feature extraction for speech and music models. Keyword Spotting through Image Recognition [11] looks at the effect of virtual adversarial training on the keyword task.

Evaluation

One of this dataset's primary goals is to enable meaningful comparisons between different models' results, so it's important to suggest some precise testing protocols. As a starting point, it's useful to specify exactly which utterances can be used for training, and which must be reserved for testing, to avoid overfitting. The dataset

Data	V1 Training	V2 Training
V1 Test	85.4%	89.7%
V2 Test	82.7%	88.2%

Table 1: Top-One accuracy evaluations using different training data

download includes a text file called `validation_list.txt`, which contains a list of files that are expected to be used for validating results during training, and so can be used frequently to help adjust hyperparameters and make other model changes. The `testing_list.txt` file contains the names of audio clips that should only be used for measuring the results of trained models, not for training or validation. The set that a file belongs to is chosen using a hash function on its name. This is to ensure that files remain in the same set across releases, even as the total number in the same set across releases, even as the total number changes, so avoid set corss-contaimination when trying old models on the more recent test data. The Python implementation of the set assignment algorithm is given in the TensorFlow tutorial code [12] that is a companion to the dataset.

Historical Evaluations

Version 1 of the dataset [9] was released August 3rd 2017, and contained 64,727 utterances from 1,881 speakers. Training the default convolution model from the TensorFlow tutorial (based on Convolutional Neural Networks for Small-footprint Keyword Spotting [17]) using the V1 training data gave a Top-One score of 85.4%, when evaluated against the test set from V1. Training the same model against version 2 of the data set [10], documented in this paper, produces a model that scores 88.2% Top-One on the training set extracted from the V2 data. A model trained on V2 data, but evaluated against the V1 test set gives 89.7% Top-One, which indicates that the V2 training data is responsible for a substantial improvement in accuracy over V1. The full set of results are shown in Table ??

- Convolutional recurrent neural networks for small-footprint keyword spotting [3]
- Honk: A PyTorch reimplementation of convolutional neural networks for keyword spotting [18]
- An experimental analysis of the power consumption of convolutional neural networks for keyword spotting [20]
- Transfer learning for speech recognition on a budget [13]
- Learning and transferring mid-level image representations using convolutional neural networks [15]
- Deep residual learning for small-footprint keyword spotting [19]

2 Method

The method is inspired by [11] where three different models have been evaluated on their capability to handle audio data transformed to images. One of the baseline models is the MNIST model which is also used in the Tensorflow Speech Recognition tutorial [2]

1. methodology, types of analyses, selection of the method

Xavier Glorot initialization [7]

Taken from [11]: For an $m \times x$ dimensional matrix M , $M_{i,j}$ is assigned values selected uniformly from the distribution $[-\epsilon, \epsilon]$, where

$$\epsilon = \frac{\sqrt{6}}{\sqrt{m+n}} \quad (1)$$

Xavier initialization is shown in equation 1

3 Set-up

3.1 Dataset

The TensorFlow Speech Recognition Challenge hosted by Kaggle [1] is using the Speech Commands Data Set v0.01 and contains 64,727 audio files and 31 class labels. The data set is publicly available [9] and is explained in more detail by Warden [21]. The characteristics of the original v0.01 data set are explained as follows by Warden:

Each utterance is stored as a one-second (or less) WAVE format file, with the sample data encoded as linear 16-bit single-channel PCM values, at a 16 KHz rate. There are 2,618 speakers recorded, each with a unique eight-digit hexadecimal identifier assigned as described above. The uncompressed files take up approximately 3.8 GB on disk, and can be stored as a 2.7 GB gzip-compressed tar archive.

The original class distribution of 31 words had to be reduced to 12 words to comply with the Kaggle competition guidelines, namely 10 concrete words **yes**, **no**, **up**, **down**, **left**, **right**, **on**, **off**, **stop**, **go** and two placeholder words **unknown**, **silence**. Words which do not belong to the 10 concrete words are merged into the **unknown** class label, while background noise and simple silence are merged into the **silence** class label. Provided files are further explained on the competition page [1] as follows:

- **train.7z**: Contains a few informational files and a folder of audio files. The audio folder contains subfolders with 1 second clips of voice commands, with the folder name being the label of the audio clip. There are more labels that should be predicted. The labels you will need to predict in Test are **yes**, **no**, **up**, **down**, **left**, **right**, **on**, **off**, **stop**, **go**. Everything else should be considered either **unknown** or **silence**. The folder **_background_noise_** contains longer clips of "silence" that you can break up and use as training input. The files contained in the training audio are not uniquely named across labels, but they are unique if you include the label folder. For example, **00f0204f_nohash_0.wav** is found in 14 folders, but that file is a different speech command in each folder. The files are named so the first element is the subject id of the person who gave the voice command, and the last element indicated repeated commands. Repeated commands are when the subject repeats the same word multiple times. Subject id is not provided for the test data, and you can assume that the majority of commands in the test data were from subjects not seen in train. You can expect some inconsistencies in the properties of the training data (e.g., length of the audio).
- **test.7z**: Contains an audio folder with 150,000+ files in the format **clip_000044442.wav**. The task is to predict the correct label. Not all of the files are evaluated for the leaderboard score.
- **sample_submission.csv**: A sample submission file in the correct format.

After merging words which do not comply with the competition guidelines, the class distribution changed from a balanced distribution to a rather unbalanced distribution towards the "unknown" label as depicted in table 2

Class	Frequency	Percentage
unknown	41039	0.634032
stop	2380	0.036770
yes	2377	0.036723
up	2375	0.036693
no	2375	0.036693
go	2372	0.036646
right	2367	0.036569
on	2367	0.036569
down	2359	0.036445
off	2357	0.036414
left	2353	0.036353
silence	6	0.000093

Table 2: Descending class distribution of merged data set

3.2 Preprocessing

The preprocessing pipeline is described in further detail in the following subsections. The pipeline is rather simple at the moment and provides room for improvement which is described in section 8. The overall preprocessing approach is depicted in figure 3. A Voice-activity detection (VAD) system has not been used based on the simple nature of the one-second audio clips.

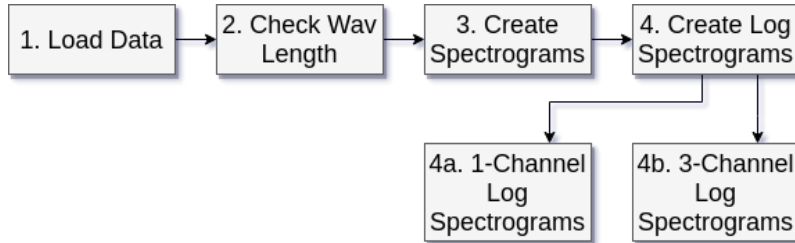


Figure 3: Preprocessing

3.2.1 Load Data

Based on the fact that there is no index file in a common format like csv which maps data entries to labels, the approach to loading the training data involves iterating labeled folders. The training folder contains one folder for each label which then contains the actual wav files for training purposes. By iterating through the different label folders and then copying the paths of the training files, a more convenient numpy array for training purposes is used. A subset of this array is depicted in table 3.

The different paths were iterated and the wav files were load into memory by using the `scipy.io.wavfile` package. The sample rate of the wav files remains at 16000 KHz.

Path	Label
'../data/train/audio/down/fad7a69a_nohash_1.wav'	'down'
'../data/train/audio/go/fa7895de_nohash_0.wav'	'go'
'../data/train/audio/left/fa7895de_nohash_0.wav'	'left'
'../data/train/audio/no/a1c63f25_nohash_2.wav'	'no'
'../data/train/audio/off/4a1e736b_nohash_4.wav'	'off'
'../data/train/audio/on/4a1e736b_nohash_4.wav'	'on'
'../data/train/audio/right/b71ebf79_nohash_0.wav'	'right'
'../data/train/audio/_background_noise_/doing_the_dishes.wav'	'silence'
'../data/train/audio/stop/fa7895de_nohash_0.wav'	'stop'
'../data/train/audio/two/fa7895de_nohash_0.wav'	'unknown'
'../data/train/audio/up/4c841771_nohash_2.wav'	'up'
'../data/train/audio/yes/b71ebf79_nohash_0.wav'	'yes'

Table 3: Training data - Numpy array representation

3.2.2 Check Wav Length

In order to have a consistent dataset with clips of one second length, the length of each wav file has been checked. Two approaches to guarantee one second clips have been applied:

- *length < 1 second*: Pad the clip with constant zeros
- *length > 1 second*: Cut the clip from the beginning to the one second mark

3.2.3 Create Spectrograms

To transform audio data into images, the \mathbb{R}^{16000} audio vectors have been transformed into spectrograms. The code for the transformation is given in listing 1.

```

1 from scipy import signal
2 from scipy.io import wavfile
3 import numpy as np
4
5 def get_spectrogram(audio_path, num_channels=1):
6     (sample_rate, sig) = wavfile.read(audio_path)
7
8     if sig.size < sample_rate:
9         sig = np.pad(sig, (sample_rate - sig.size, 0), mode='constant')
10    else:
11        sig = sig[0:sample_rate]
12
13    # f = array of sample frequencies
14    # t = array of segment times
15    # Sxx = Spectrogram of x. By default, the last axis of Sxx corresponds
16    # to the segment times.
17    f, t, Sxx = signal.spectrogram(sig, nperseg=256, noverlap=128)
18    Sxx = (np.dstack([Sxx] * num_channels)).reshape(129, 124, -1)
19
20    return f, t, Sxx

```

Listing 1: Get spectrogram code

After a first inspection of the spectrograms given in figure 4, it is obvious that the spectrograms do not contain as much visible features as expected. Previous research [11] also suggested that using log spectrograms is more beneficial than using simple spectrograms. Spectrograms were reshaped into 129 x 124 dimensions which differs

from the log spectrograms but does not influence the experiment in any way because solely log spectrograms were used for training purposes.



Figure 4: Spectrogram samples of nine different classes

3.2.4 Create Log Spectrograms

In contrast to figure 4, the code depicted in listing 2 produces log spectrograms which contain more visual features as seen in figure 5 which should be beneficial for the model training. One potential problem however is shown in figure 5 at the "silence" class. The padding with zeroes presents itself with a dark bar at the beginning which might introduce more noise into the dataset. Nevertheless, for this experiment this potential problem is ignored and could be tackled in future work. Log spectrograms were reshaped into 99 x 161 dimensions.

```
1 from scipy import signal
2 from scipy.io import wavfile
3 import numpy as np
4
5 def get_log_spectrogram(audio_path, window_size=20, step_size=10, eps=1e
6   -10, num_channels=1):
7     (sample_rate, sig) = wavfile.read(audio_path)
8
9     if sig.size < 16000:
10         sig = np.pad(sig, (sample_rate - sig.size, 0), mode='constant')
11     else:
12         sig = sig[0:sample_rate]
13
14     nperseg = int(round(window_size * sample_rate / 1e3))
15     noverlap = int(round(step_size * sample_rate / 1e3))
16
17     # f = array of sample frequencies
18     # t = array of segment times
19     # Sxx = Spectrogram of x. By default, the last axis of Sxx corresponds
20     # to the segment times.
21     f, t, Sxx = signal.spectrogram(sig,
22                                     fs=sample_rate,
23                                     window='hann',
24                                     nperseg=nperseg,
25                                     noverlap=noverlap,
26                                     detrend=False)
27
28     log_spectrogram = np.log(Sxx.T.astype(np.float32) + eps)
29     log_spectrogram = (np.dstack([log_spectrogram] * num_channels)).
30     reshape(99, 161, -1)
31
32     return f, t, log_spectrogram
```

Listing 2: Get log spectrogram code

A distinction between 1- and 3-channel log spectrograms has been made because this project uses pre-trained CNN models which are trained on ImageNet data. ImageNet data is in its core based on 3-dimensional RGB data/images while (log) spectrograms are 1-dimensional images and are only depicted in green colors in figure 5 based on the settings in `matplotlib` package which shows grayscale images in a green spectrum. A quick fix to this problem is the duplication of the 1-dimensional spectrogram data and therefore mimicking 3-dimensional RGB data by having the grayscale data copied over three channels as depicted in figure 6.

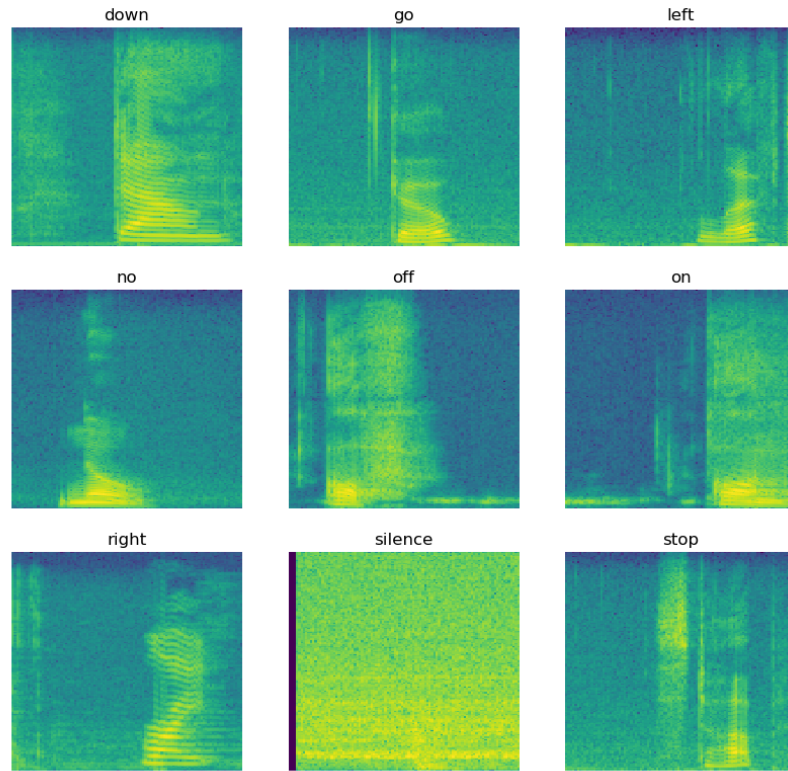


Figure 5: Log spectrogram samples of nine different classes

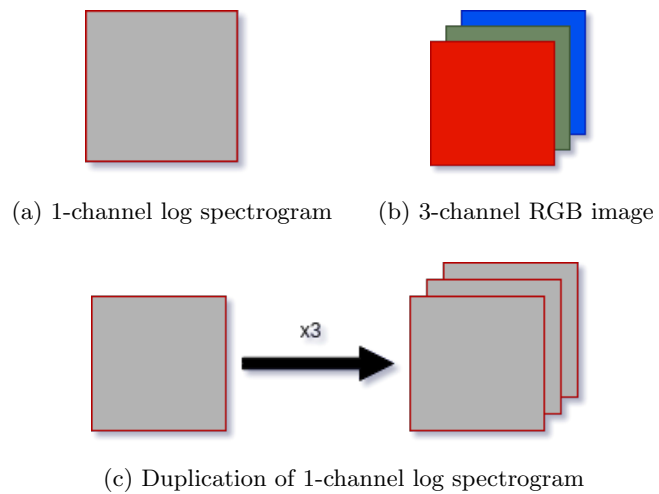


Figure 6: Conversion between 1-channel and 3-channel log spectrograms

3.3 Models

The choice of models can be divided into the following groups: **Baseline** and **CNN models with pre-trained ImageNet weights**. Table 4 shows an overview of the evaluated models in this project with the complexity of each model based on its amount of trainable parameters.

MNIST and the lightweight CNN model are used as a baseline because they showed good performance in previous work when applied to spectrograms. The other models, namely VGG16, Inception V3 and ResNet were used because they each won the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) in their respective years and differ in their architecture. The baseline models were initialized with Xavier initialization and their performance were not evaluated with ImageNet weight initialization.

Model	Total params	Trainable params	Non-trainable params
Leightweight CNN	723,968	723,454	514
VGG16	23,676,748	23,659,340	17,408
Inception V3	29,185,836	29,137,068	48,768
MNIST	55,038,988	54,929,868	109,120
ResNet50	75,182,988	75,029,516	153,472

Table 4: Model complexity ordered by amount of parameters

Models have been trained over 10 epochs each with either Xavier Glorot initialization or pre-trained ImageNet weights while all layers remained trainable which is against traditional transfer learning techniques. After 10 epochs the final accuracy has been evaluated in order to see if pre-trained ImageNet weights would give a boost to the convergence rate. Another evaluation has been made after the first epoch with regards to accuracy to see if pre-trained ImageNet weights would give a higher start accuracy based on the already learned image features from another domain than the provided spectrograms.

Drawing architectures with draw.io?

3.3.1 Baseline models

- MNIST: show source, architecture and what has been changed (BatchNormalization)
- Lightweight CNN: show source, architecture, and what has been changed

3.3.2 CNN models

- VGG16: show source, architecture, paper, won competition and what has been changed?
- Inception V3: show source, architecture, paper, won competition and what has been changed? What differs in this architecture from VGG16?
- ResNet50: show source, architecture, paper, won compeition and what has been changed? SkipConnection as main difference

4 Experiments

4.1 Baseline

4.2 CNNs

5 Analysis and Results

5.1 Baseline

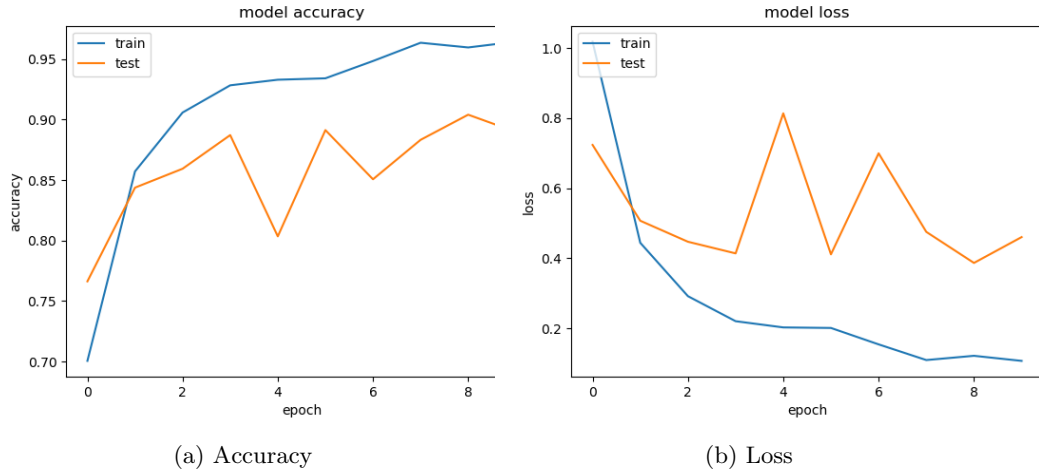


Figure 7: Accuracy and loss after 10 epochs for the MNIST model with Xavier Glorot initialization

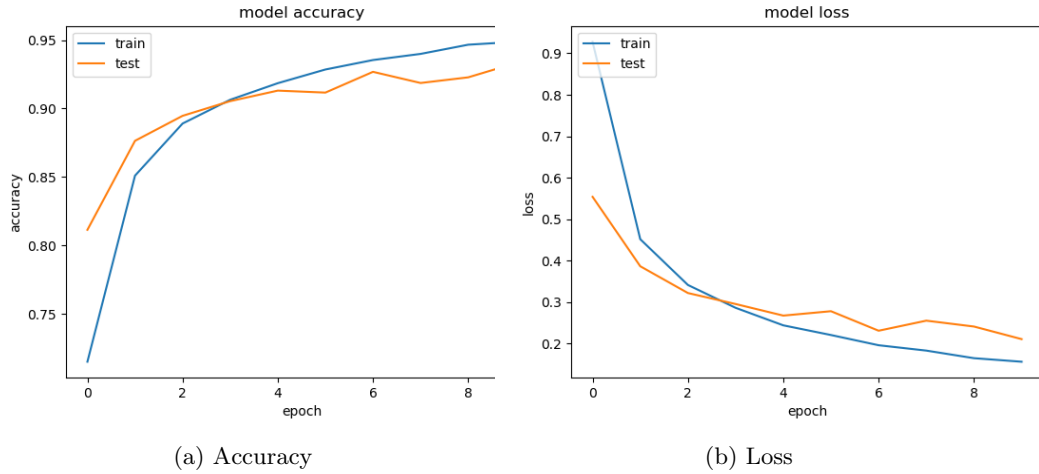


Figure 8: Accuracy and loss after 10 epochs for the Lightweight CNN model with Xavier Glorot initialization

Model	Train Acc	Train Loss	Val Acc	Val Loss	Time (sec)
MNIST	0.9595	0.1213	0.9039	0.3866	1064.72
Leightweight CNN	0.9487	0.1564	0.9332	0.2109	532.73

Table 5: Final baseline results with Xavier Glorot initialization

Model	Train Acc	Train Loss	Val Acc	Val Loss
MNIST	0.7005	1.0187	0.7662	0.7237
Leight CNN	0.7150	0.9276	0.8114	0.5540

Table 6: Baseline results after one epoch with Xavier Glorot initialization

5.2 CNNs

5.2.1 Xavier initialization

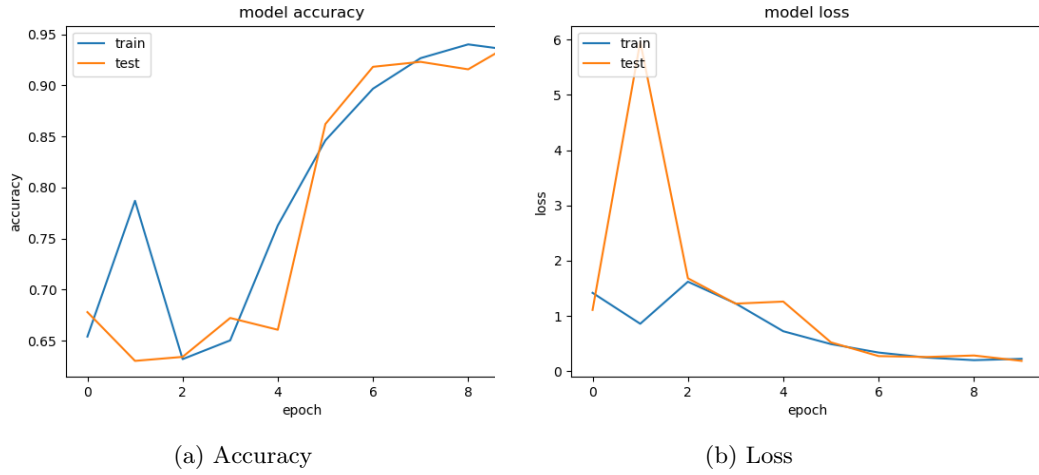


Figure 9: Accuracy and loss after 10 epochs for the Inception V3 model with Xavier Glorot initialization

Model	Train Acc	Train Loss	Val Acc	Val Loss	Time (sec)
Inception V3	0.9338	0.2238	0.9427	0.1868	2332.23
VGG16	0.9723	0.0900	0.9583	0.2011	2530.47
ResNet50	0.9548	0.1421	0.9445	0.1873	3807.93

Table 7: Final CNN results with with Xavier Glorot initialization

Model	Train Acc	Train Loss	Val Acc	Val Loss
Inception V3	0.6539	1.4176	0.6778	1.1090
VGG16	0.6519	1.3202	0.6906	1.3379
ResNet50	0.6916	1.2478	0.6351	5.7467

Table 8: CNN results after one epoch with Xavier Glorot initialization

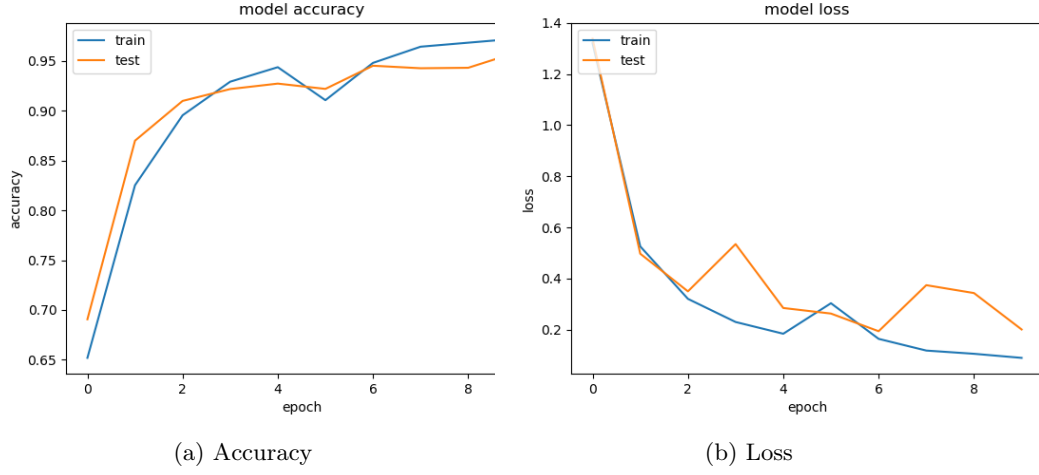


Figure 10: Accuracy and loss after 10 epochs for the VGG16 model with Xavier Glorot initialization

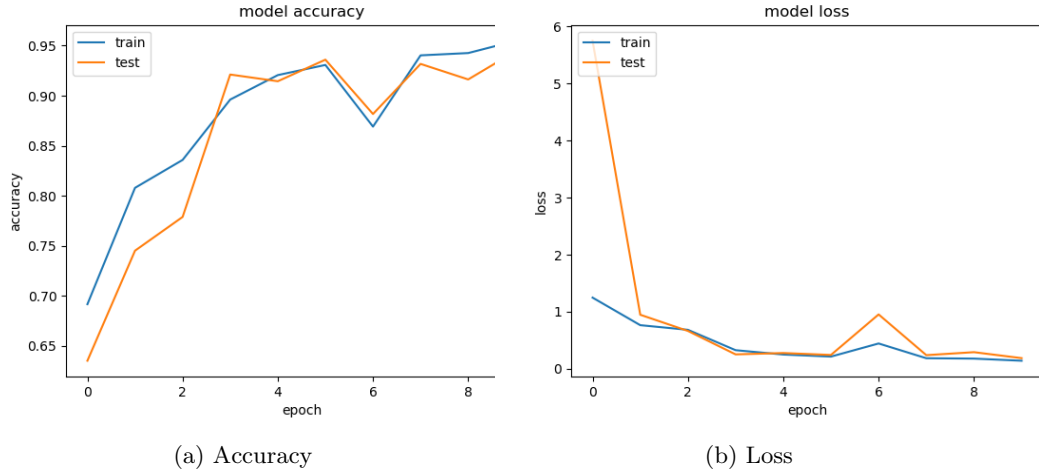


Figure 11: Accuracy and loss after 10 epochs for the ResNet50 model with Xavier Glorot initialization

5.2.2 Imagenet weight initialization

Test

Model	Train Acc	Train Loss	Val Acc	Val Loss	Time (sec)
Inception V3	0.6334	1.6597	0.6401	1.6433	2184.54
VGG16	0.9745	0.0912	0.9611	0.1730	2541.83
ResNet50	0.9134	0.2991	0.9252	0.3055	3653.67

Table 9: Final CNN results with Imagenet initialization

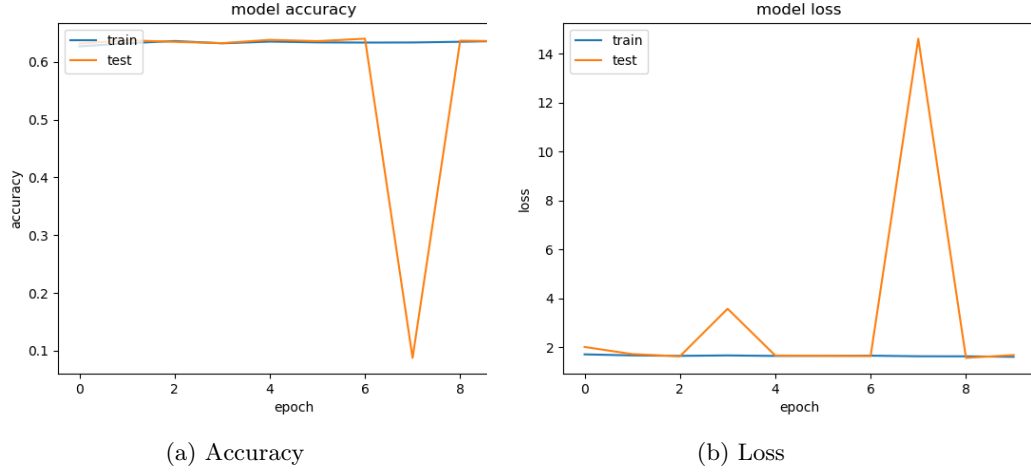


Figure 12: Accuracy and loss after 10 epochs for the Inception V3 model with Imagenet initialization

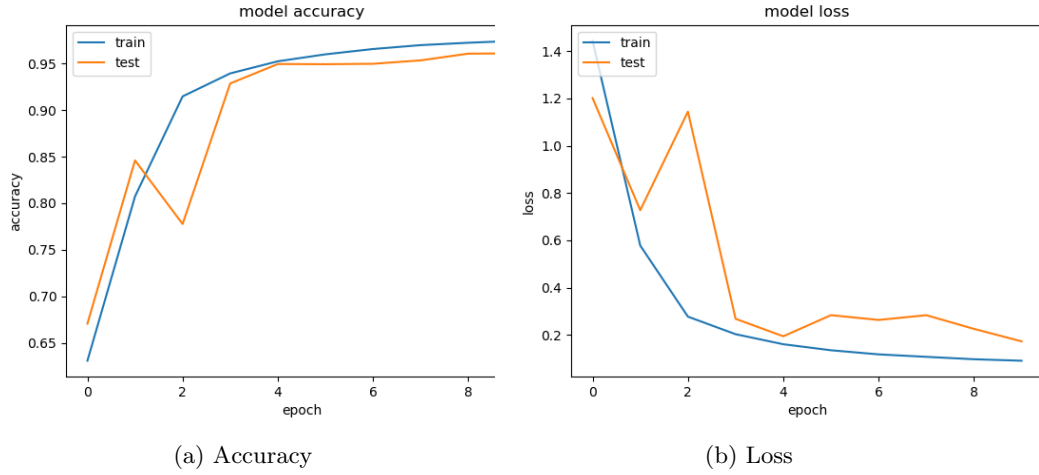


Figure 13: Accuracy and loss after 10 epochs for the VGG16 model with Imagenet initialization

Model	Train Acc	Train Loss	Val Acc	Val Loss
Inception V3	0.6268	1.7132	0.6315	2.0153
VGG16	0.6307	1.4430	0.6706	1.2011
ResNet50	0.6389	1.4214	0.5931	5.1381

Table 10: CNN results after one epoch with Imagenet initialization

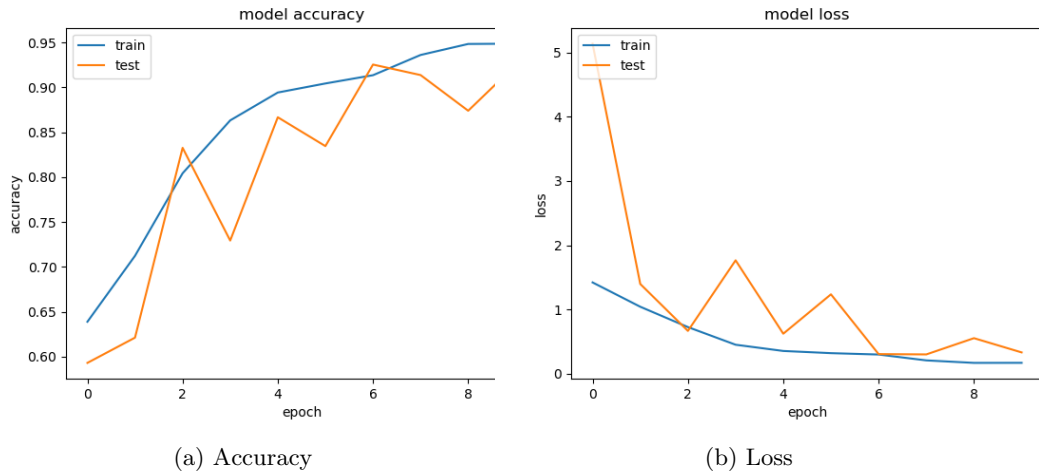


Figure 14: Accuracy and loss after 10 epochs for the ResNet50 model with Imagenet initialization

6 Discussion

7 Conclusion

8 Future Work

- Use other image representations like spectrograms based on MFCC
- Use recursive layer freezing to investigate which amount of layers would be optimal for transfer learning
- Use additional preprocessing steps like discarding bad audio
- Look further into Batch Normalization
- Look into correlation between accuracy and dimensions of spectrograms

9 References

10 Appendix

- the experiment(s) may be carried out in collaboration with others. In that case: specify in the “author’s statement” everybody’s contribution
- the thesis itself is written individually and assessed individually
- the ASR performance itself is not relevant for the assessment of the thesis
- the RQ, the literature embedding of the RQ, the description of the method, the justification and set-up of the experiment are relevant for the assessment
- the general university guidelines apply (e.g., with respect to plagiarism)
- there is no minimum number of pages for the thesis

	experimental	theoretical
aspect	(max. points)	(max. points)
Research Question (RQ)	20	20
Literature embedding of the RQ	20	40
Method	20	
Justification experiment(s)	10	
Set-up experiment(s)	30	
Discussion and Conclusion	30	70
Use of figures and tables	10	10
Overall completeness	20	20
Overall clarity, transparency	20	20
Overall coherence (from intro to conclusion)	20	20
<i>Total</i>	<i>200</i>	<i>200</i>

Figure 15: Weighted grading

References

- [1] Tensorflow speech recognition challenge. <https://www.kaggle.com/c/tensorflow-speech-recognition-challenge>. Accessed: 2019-08-05.
- [2] Tensorflow speech recognition tutorial. https://www.tensorflow.org/tutorials/sequences/audio_recognition. Accessed: 2019-08-05.
- [3] Sercan O Arik, Markus Kliegl, Rewon Child, Joel Hestness, Andrew Gibiansky, Chris Fougner, Ryan Prenger, and Adam Coates. Convolutional recurrent neural networks for small-footprint keyword spotting. *arXiv preprint arXiv:1703.05390*, 2017.
- [4] Guoguo Chen, Carolina Parada, and Georg Heigold. Small-footprint keyword spotting using deep neural networks. In *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4087–4091. IEEE, 2014.
- [5] Dan Cireşan, Ueli Meier, and Jürgen Schmidhuber. Multi-column deep neural networks for image classification. *arXiv preprint arXiv:1202.2745*, 2012.
- [6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [7] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256, 2010.
- [8] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- [9] Speech commands dataset version 1. http://download.tensorflow.org/data/speech_commands_v0.01.tar.gz. Accessed: 2019-08-05.

- [10] Speech commands dataset version 2. http://download.tensorflow.org/data/speech_commands_v0.02.tar.gz. Accessed: 2019-08-05.
- [11] Sanjay Krishna Gouda, Salil Kanetkar, David Harrison, and Manfred K Warmuth. Speech recognition: Keyword spotting through image recognition. *arXiv preprint arXiv:1803.03759*, 2018.
- [12] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [13] Julius Kunze, Louis Kirsch, Ilya Kurenkov, Andreas Krug, Jens Johannsmeier, and Sebastian Stober. Transfer learning for speech recognition on a budget. *arXiv preprint arXiv:1706.00290*, 2017.
- [14] Jongpil Lee, Taejun Kim, Jiyoung Park, and Juhan Nam. Raw waveform-based audio classification using sample-level cnn architectures. *arXiv preprint arXiv:1712.00866*, 2017.
- [15] Maxime Oquab, Leon Bottou, Ivan Laptev, and Josef Sivic. Learning and transferring mid-level image representations using convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1717–1724, 2014.
- [16] George Retsinas, Giorgos Sfikas, and Basilis Gatos. Transferable deep features for keyword spotting. In *Multidisciplinary Digital Publishing Institute Proceedings*, volume 2, page 89, 2018.
- [17] Tara Sainath and Carolina Parada. Convolutional neural networks for small-footprint keyword spotting. 2015.
- [18] Raphael Tang and Jimmy Lin. Honk: A pytorch reimplementation of convolutional neural networks for keyword spotting. *arXiv preprint arXiv:1710.06554*, 2017.
- [19] Raphael Tang and Jimmy Lin. Deep residual learning for small-footprint keyword spotting. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5484–5488. IEEE, 2018.
- [20] Raphael Tang, Weijie Wang, Zhucheng Tu, and Jimmy Lin. An experimental analysis of the power consumption of convolutional neural networks for keyword spotting. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5479–5483. IEEE, 2018.
- [21] Pete Warden. Speech commands: A dataset for limited-vocabulary speech recognition. *arXiv preprint arXiv:1804.03209*, 2018.