

Interpolation Based Neural Audio Synthesis using Convolutional Autoencoders

Benedikt Langer, BSc



MASTERARBEIT

eingereicht am
Fachhochschul-Masterstudiengang

Mobile Computing

in Hagenberg

im Juni 2023

Advisor:

FH-Prof. DI Stephan Selinger
Alexander Palmanshofer, BSc MSc

© Copyright 2023 Benedikt Langer, BSc

This work is published under the conditions of the *Creative Commons License Attribution-NonCommercial-NoDerivatives 4.0 International* (CC BY-NC-ND 4.0)—see <https://creativecommons.org/licenses/by-nc-nd/4.0/>.

Declaration

I hereby declare and confirm that this thesis is entirely the result of my own original work. Where other sources of information have been used, they have been indicated as such and properly acknowledged. I further declare that this or similar work has not been submitted for credit elsewhere. This printed copy is identical to the submitted electronic version.

Hagenberg, June 27, 2023

Benedikt Langer, BSc

Contents

Declaration	iv
Preface	viii
Abstract	ix
Kurzfassung	x
1 Introduction	1
2 Related Works	2
2.1 Neural Audio Synthesis	2
2.2 Audio Style Transfer	7
2.3 Image Style Transfer	10
3 Approach	12
3.1 Motivation	12
3.2 Overview	13
3.3 Pre-processing	14
3.3.1 Spectrograms and STFT	15
3.4 ML-Model	17
3.4.1 Neural Networks - Introduction	18
3.4.2 Convolutional Neural Networks	18
3.4.3 Autoencoder	19
3.4.4 Optimizer	22
3.5 Synthesis of novel sounds	23
3.5.1 Interpolation in latent space	23
3.6 Post Processing	24
3.7 Dataset	25
3.7.1 NSynth Dataset	26
4 Experiment	27
4.1 Implementation Environment	27
4.2 Training	28
4.2.1 Training configuration	28
4.2.2 Training Execution	29

4.3	Initial Experiments	29
4.3.1	Whole Spectrograms as Input	30
4.3.2	Spectrograms of signalframes as Input	31
4.4	Experiments single frequency vectors	32
4.4.1	Experiments for Synthesis	34
4.5	Experiments with slices of spectrograms	36
4.5.1	Audio synthesis	38
4.5.2	Reconstruction and post-processing	38
4.6	Experiments with mel-scale	39
4.6.1	Pre-processing	39
5	Results	42
5.1	Results regarding models with single Frequency Vectors	42
5.1.1	Experiments of single reconstruction	43
5.1.2	Experiments with interpolation in embedding	45
5.2	Results of experiments with spectrogram frames	46
5.2.1	Experiments of single reconstruction	47
5.2.2	Experiments with interpolation in embedding	49
5.3	Results with mel-spectrograms	51
5.3.1	Experiments of single reconstruction	53
5.3.2	Experiments with interpolation in embedding	55
6	Discussion/Evaluation	59
6.1	Interpretation of results / observations	59
6.1.1	Quality of audio	59
6.1.2	Differences between instruments	59
6.1.3	Impact of model configurations	59
6.1.4	Impact of pre processing	59
6.1.5	Impact of post-processing	59
6.2	Comparison to other approaches	59
6.3	Limitations/difficulties	59
6.4	Outlook	59
7	Conclusion	60
8	Future Work	61
A	Additional Graphics	62
B	Supplementary Materials	63
B.1	PDF Files	63
B.2	Media Files	63
B.3	Online Sources (PDF Captures)	63
C	Questionnaire	64
D	LaTeX Source Code	65

References	66
Literature	66
Software	68
Online sources	68

Preface

Abstract

This should be a 1-page (maximum) summary of your work in English.

Kurzfassung

An dieser Stelle steht eine Zusammenfassung der Arbeit, Umfang max. 1 Seite. ...

Chapter 1

Introduction

Chapter 2

Related Works

There already exist a few good approaches around the topic of neural style transfer or generating audio using neural networks, that present rather good solutions. Some of these studies have proven, that with neural networks it is possible to generate synthesized audio up to a certain quality. Those approaches can be categorized into different areas, as their principle and methodology differ in certain ways. As this field is related to the technique of image style transfer, a lot of works apply those methods to audio, and respective audio spectrograms and thus, call it explicitly audio style transfer. Secondly because those solutions are defining a combination specifically of content and style. This topic will be discussed in more detail in section 2.2. All methods which do not include this principle of content and style can be categorized to the technique of neural audio synthesis or simply just audio synthesis (see 2.1). Those methods incorporate mostly autoencoder networks.

2.1 Neural Audio Synthesis

Neural audio synthesis is the field of creating or synthesizing novel sounds with the help of neural networks. The approach is similar and related to the field of audio style transfer. As described before, approaches in this domain differ in certain ways to neural style transfer. As a major difference, with respect to neural audio synthesis, no content or style sound is specified, which means, that for the creation of novel sounds, two sound sources are used equally. While audio style transfer also gets frequently applied on whole audio samples or musical pieces, in audio synthesis the focus is more directed on the application for single notes. Especially concerning autoencoder networks, neural audio synthesis also includes the tasks of learning important sound features for compression and recreation of the input data. How different approaches are designed and which machine learning techniques and which results could be obtained, will be described in the following sections.

Probably one of the most prominent solutions in the field of neural audio synthesis comes from *Engel et al.* [6]. With their work “Neural Audio Synthesis of Musical Notes with WaveNet Autoencoders” they have proposed a system that is capable of synthesizing audio as well as interpolating/morphing encoded audio data of two instruments to create new audio. In addition, there also exists a publicly available dataset called

“NSynth” [6] that contains a large amount of high quality audio samples. The dataset has been applied for the training of the model over the course of this specific project. In their work regarding the synthesis, *Engel et al.* developed and compared two different approaches with two different kind of networks. Basically they have a similar structure, as they are both designed as autoencoders. However they do accept different formats of audio data and thus have different components. While the one kind of network operates on time-continuous data the other one is trained on the spectral representation of audio samples. Throughout the work, the second technology using spectrograms is referenced as Baseline Model. This one consists of convolutional layers, with leaky ReLU activation, except the last layers in encoder and decoder. Furthermore batch normalization is applied as well. This spectral autoencoder is compared to the so called “WaveNet-style Autoencoders”, that are trained on continuous time signal. The autoencoder structure enables learning efficient encodings of the music data, which are representing essential features from the original audio.

In order to create new sounds they take the encoded data from the embedding space of two instruments and interpolate them linearly. In addition they applied the decoder part to reconstruct audio data. As a result they were able to create new sounds with the characteristics of two different audio signals. Comparing the performance of the two different networks used, they found the WaveNet-style autoencoder to be advantageous. This was proven both by the error scores for reconstructing the audios or auditory quality and by quantitative comparison with a pitch and quality classifier model. Nevertheless it was also concluded that the spectral baseline model has a strong performance.

The results regarding the WaveNet-style Autoencoder can be explored via Engel’s online AI-Experiment called “Sound Maker”¹.

In further publications and approaches, *Engel* continued the research on neural audio synthesis by applying other network structures for this purpose. With respect to generative adversarial networks (GANs) and recurrent neural networks (RNNs), two more works have been published in the field of neural audio synthesis [5, 10]. Similar to their work concerning WaveNet-style and convolutional autoencoder, they conducted experiments in (re)synthesizing audios, as an example interpolation of extracted features has been done. Eventually the discussed works show a suitability for audio synthesis and also highlight a major speedup in the computation of synthesized audio samples.

The work by *Natsiou et al.* does not explicitly mention the term neural audio synthesis in its title, but deals with it throughout the article [20]. In their work they research the reconstruction capacity of (stacked) convolutional autoencoders in terms of log-mel-spectrograms and carry out experiments on different configurations. In their experiments they evaluate the effectiveness of autoencoders in terms of neural audio synthesis whereas also feasible improvements through additional techniques are measured. Additionally, they mention that their work explores musical timbre compression. Here the synthesis gets specifically referenced to timbre synthesis. As audio spectrograms exist with different scales, this approach uses, in contrast to others, the log-mel scale. This scale proves to be beneficial, as it already captures the most significant properties with the characteristics of consuming less memory and computational power. For the training the authors used the NSynth Dataset proposed by *Engel et al.* [6], whereas just

¹“Sound Maker” <https://magenta.tensorflow.org/nsynth-instrument>

a sub-sample, containing samples of different instruments of one single pitch, was considered. The model(s) that were used throughout their experiments followed the general structure of a (stacked) convolutional autoencoder network, which is composed mainly of 2D convolutional layers. For experimental reasons, additional layers and techniques such as pooling layers, fully connected layers, dropout, kernel regularization got applied (added/removed).

In order to measure the results of their experiments they applied error metrics such as root mean squared error (RMSE) and structural similarity index (SSIM). As a result these metrics cannot accurately express anything about the quality. Because of this reason, they also introduced a precision and recall score and combined it in a F1 score. In order to generate sounds from the spectrograms, they were utilizing the preserved phase information, unless there was no modification of the embedding. In the latter case, the Griffin Lim phase estimation algorithm was applied, as no phase information is present. Regarding the results that could be obtained by reconstructing spectrograms (without modification in latent space), some interesting findings could be extracted. To their surprise, by reducing the size of the latent space, they found out that the smaller it is, more accurate spectrograms with a smoother distribution could be generated. Also, in some cases where kernel regularization got applied, the spectrograms were more accurate, while with dropout layers no improvement could be achieved. The use of (max) pooling also resulted in a more accurate time-frequency resolution with less noise than with just convolution layers. Finally, after the removal of the fully connected layers it showed that the quality was significantly better, as spatial information got preserved better.

Regarding audio synthesis, *Colonel et al.* published a few works where they investigated the suitability of autoencoder networks in connection with audio synthesis [2, 3, 17]. Starting in 2017 they proposed an autoencoder based audio synthesis through compression and reconstruction of audio spectrograms [17]. In contrast to the previously mentioned approaches, this one is based on fully connected layers without convolutions. Also a self-made dataset was created, with their own synthesizer. In contrast to e.g. the NSynth dataset, this one contains polyphonic notes and thus more complex harmonies. During the experiment they trained different parameterized networks, where they vary the depth and width of the network and its layers as well as the activation functions and different optimizers. The mean squared error (MSE) was used as the error metric. Comparing these scores regarding networks of one or two hidden layers on each side show, the network using the Adam optimizer worked out best in contrast to using Momentum as the optimizer. These networks, based on sigmoid activation functions worked best when less compression is applied. The concept of 4 hidden layers showed having a mix of ReLU and sigmoid activation functions worked out best. Additionally by applying regularization methods such as dropout and L2 penalty, the latter proved to be better, as the results were of better auditory quality. Further results showed that sigmoid activations led to fuller sound than with ReLU. Furthermore by using bias terms, it could be observed that noise was present in the results. As a consequence, despite of the better convergence, they chose to let them out. The authors concluded that using a network with 4 hidden layers and a composition of sigmoid and ReLU worked out best also in terms of auditory quality.

Another work by *Colonel et al.* was proposed in 2018, which actually states an

improvement of the method, described in their previous work from 2017 [2]. Those improvements contain the use of a phase reconstruction method not used before, which allows a direct activation of the latent space. For an improved model convergence, the autoencoder was designed asymmetrically, via input augmentation. This means that they padded the input magnitude data with different permutations, being first or second order difference or mel-frequency cepstral coefficients (MFCC). Whereas in the previous work only MSE was contemplated as the error metric, here, a comparison of several cost functions was used. Within the cost functions, the mean absolute error (MAE), as well as the spectral convergence cost function (SC) with L2 penalty, was considered. The penalization of the total spectral power proves to be advantageous as the power in the output is more accurate compared to others. In comparison to their work from 2017, they also left out additional bias terms, but decided to just use ReLU activations instead of a mixture with sigmoid.

Eventually improvements to their previous results could be achieved regarding the additional methods that were applied. Concerning the augmentation of the input data, a significant improvement regarding score could be reached, whereas augmenting with first order difference outperformed all other approaches. Concerning the generated sound, it could be observed that by padding with the MFCCs a different sound palette was present. In another comparison to their baseline, *Colonel et al.* introduced the possibility to omit the encoder part of the network. This directly enables the activation of the innermost 8 neuron layers while the decoder can generate novel sounds. As no phase information was present, estimations were done via a method called real-time phase gradient heap integration, which enables the generation of a playable sound. In addition to this work, the authors implemented a small program including a GUI, where it is possible to directly interact and activate the innermost neurons (eight control values in latent space) to generate new sounds.

In a more recent work, *Colonel et al.* implemented and compared autoencoder networks with different topologies regarding their performance for musical timbre generation [3]. This work already utilizes findings and methodologies from previous works. Based on a study from 2018, they implemented a mechanism to directly activate and control the latent space of a trained autoencoder with a graphical tool. They found out that this technique proved to be difficult in terms of controlling the latent space. To overcome this issue and improve the work, they added chroma-based input augmentation to improve the reconstruction performance in this approach. The chroma values are based on the 12 note (western) scale to represent the dominant note present in an audio sample. Besides this type of input augmentation, they also implemented a so-called skip connection, where the latent space is conditioned by the chroma value. In this work the chroma values are represented via a one-hot encoded representation for each training sample, where the maximum value is set to one while all others are set to zero. Consequently these one-hot encoded chroma representations tell the note played in a single-note audio. With this technique the authors could shape the timbre around a specific note class. For the networks topologies, they varied the size of the bottleneck-layer (8, 3 or 2 neurons), the activation functions, the input augmentation, the use of the chroma skip connection, as well as different datasets. It should be mentioned that the authors trained and experimented with the self-generated dataset from their previous works, containing five octaves of notes, a one octave subset of it and a separate violin

note dataset.

As a result, the network with an eight neuron bottleneck, with the chroma-based input augmentation worked out best. Thus, for the rest of the experiments *Colonel et al.* were using this technique. In case of the two neuron bottleneck network, the sigmoid activation functions without skip connection worked out best for the one octave dataset. The skip connection turned out to work best for the violin dataset (sigmoid and two neurons). Finally in the case of three neurons also the variant with the skip connection worked out best for both datasets. By analyzing the latent spaces some interesting observations could be made for the sake of audio synthesis. The authors applied a clustering method to see the distribution of the values in the latent space concerning their note and timbre. Using sigmoid activations turned out to limit the values to the range of (0,1) as well as distributing the values in a more uniform manner. Also, the skip connections led to a denser representation. By taking this as an advantage and moving forward with just sigmoid activations, sampling of the latent space (with a mesh grid e.g. 350x350 for two neurons bottleneck) was done to generate a new timbre. In combination with setting the additional chroma conditioning vector to a given note class, the decoder generates the timbre that matches the chroma vector and thus, the desired note is present in the output sample.

A comparative work on autoencoders, in terms of music sound modeling, has been published by *Roche et al.* in 2019 [25]. In this work they implemented four different types of autoencoder networks, that have been compared in terms of audio synthesis. Similar to the other techniques described earlier, this one also orientates itself on the principle of an autoencoder, to project the input data to a low-dimensional space, from which input can be (re)synthesized. In the described experiments, the proposed autoencoder networks consist of (shallow) autoencoders (AEs), deep autoencoders (DAEs), recurrent autoencoders (LSTM-AEs) and variational autoencoders (VAEs) which all got compared to principal component analysis (PCA) as baseline. As sound data for training and experimenting, they used a subsample of 10,000 different randomly selected notes from the publicly available NSynth dataset. The networks that where implemented got trained on the normalized log-magnitude spectra of those samples. Regarding the structure or the depth of the different networks, the researchers used two and three layers for the DAE on each side. In the case of the VAE just one version with two layers, and one version of the LSTM-AE with one layer on each side was applied. Regarding the size of the output from the encoder (latent space), they experimented with different values in a range from 4 to 100. The conducted experiments consist of a resynthesis-analysis where the reconstruction error (RMSE in dB) of the different methods got compared. Additionally to the RMSE so called PEMO-Q scores were introduced to calculate the objective measures of perceptual audio quality.

The results showed to their surprise, that PCA outperformed the shallow autoencoder network. Continuing with DAEs, the reconstruction performed almost 20% better than the shallow AE, having an encoding size of 12 and 16. Also the error decreased faster when the dimension of the latent space was reduced. Even better results with over 23% improvement compared to PCA could be achieved by using LSTM-AEs which brought them to the conclusion, that it is feasible to use more complex architectures. The fact that more compression and thus a small latent space can be generated, is even more important for sound synthesis. In comparison the reconstruction error from the VAE

was lying between the one of the DAE and shallow AE/PCA. As the size of the latent space influences the reconstruction error, it can be stated, that the bigger the size, the lower the error. Interestingly PCA outperforms all models having an encoding size of 100. In addition to the RMSE score, the perceptual audio quality got measured with the PEMO-Q score. The results are comparable to those with RMSE, with just the LSTM-AE having a slightly lower score as compared to RMSE. In this context it was also investigated how the latent space values can be used to be controlled by musicians, and thus, the correlation between those values has been calculated. Averaged over all samples per model, it was shown that the values from LSTM have the most correlation while VAE has the worst. Having less correlation makes VAE the better candidate in terms of using the latent values as control values for synthesis (less redundancy and clear perceptual meaning). In terms of audio synthesis, including the latent space variables, *Roche et al.* also demonstrated how it could be applied for sound interpolation like in the work of *Engel et al.*. For this task they selected the latent space vectors of two sounds with different characteristics and linearly interpolated each value. By decoding and in addition, applying the inverse STFT and Griffin Lim, new interesting sounds could be generated.

2.2 Audio Style Transfer

The works discussed in this section all orientate themselves on the techniques of image style transfer. As those techniques have a significant impact on the development of audio style transfer algorithms, two important works are discussed in section 2.3. Applying the method of image style transfer to audio also means, as audio is a time-continuous signal, that it has to be brought into a similar shape, which will be done mostly by generating spectrograms out of signals. As for image style transfer, a content and a style picture is needed, this principle also gets applied to audio style transfer. In image style transfer, the style (e.g. brush strokes, colors) and content of an other image (e.g. contours, scenery) are combined, to form a new stylized image [7]. This means that in the output image, the content image looks painted with a certain “style”. Mapping this principle to the audio domain, this means, that there has to be a specific content sound (sample) that gets stylized with a certain style of a sound (e.g. style of a specific instrument). In the image domain it is difficult to distinguish content from style, whereas it is a bigger question regarding audio, that appears in the different approaches. Most authors define the style as a musical instruments’ timbre or even a musical genre. Alongside, the content might be defined as global music structure with rhythmic constraints [9]. Those scientific questions also might be influenced if whole audio samples/musical pieces might be taken to get stylized or just some single notes from an instrument. Furthermore if speech is considered as audio data, style and content is differently defined as well. Here, style could be e.g. the emotion of the voice or the speakers identity and content of the spoken words in an sample. The following works show different solutions specific to the problem of Audio Style transfer in which they also get compared and assessed.

One approach that applies this principle is the solution proposed by *Ramani et al.* in 2018 [24]. In their study they developed a neural network that is constructed as an convolutional autoencoder. Here they officially described their system as an audio style transfer algorithm. In this case the process of generating audio containing characteris-

tics of two audio signals is slightly different as in the work of *Engel et al.* They worked with two networks, namely a transformation network and a loss network. This architecture and methodology is especially inspired by the neural style transfer algorithm by *Johnson et al.* Both networks have the same structure and composition of layers. The loss network is trained to compress input spectrograms to lower dimensions, which means that the encoder part learns to preserve the high level features of the input. In addition, the decoder learns to reconstruct a spectrogram similar to the input of the network from the encoded data. For the training of the transformation network, the pre-trained weights of the loss network are used which speeds up the learning process. This means just optimization towards low level features/style has to be achieved. The trained transformation network is then able to transform an input spectrogram into a stylized spectrogram. The loss network is subsequently used to calculate the style loss but also content loss between the respective spectrograms and the output from the transformation network. This loss is minimized by backpropagation to the transformation network. Through this procedure it is possible to pass a single spectrogram through the transformation network. Subsequently the network outputs a new spectrogram containing the characteristics of itself (content) but also of a differently styled audio signal. Due to its architecture it also performed really quickly and could be used for real-time applications.

Verma et al. presented a new machine learning technique for the purpose of generating novel sounds, in their 2018 paper *Neural style transfer for audio spectrograms* [27]. In this approach they tried to apply the method for artistic image style transfer to audio and they specifically mention the approach proposed by *Gatys et al.* [7] (see section 2.3). Unlike *Gatys'* image style transfer approach, they adapted and trained an AlexNet architecture on the classification of audio samples. This kind of network is a so called convolutional neural network, where the audio gets converted into spectrograms, which can be seen as grey-scale images. An important aspect is, in this work they used the log-magnitude data of the STFT output. It also should be mentioned that they adapted the network towards a receptive size (kernel) of 3x3 instead of the larger ones in the original network, as it retains the resolution of the audio. Similarly to the image domain, the stylized output image gets initialized with random noise. Thus an input spectrogram consisting of a gaussian noise signal is utilized. This one gets iteratively optimized by minimizing the content, but also style loss via backpropagation. In the end this process creates a spectrogram combining the content of one audio with the style of a different audio sample. Additionally they found out that including additional loss terms for temporal and frequency energy envelopes, helped to improve the quality, as otherwise temporal dynamics would not get incorporated. For their experiments they imposed the style of a tuning fork onto a harp sound and also transferred the style of a violin sound onto a sample of a singing voice. In this way they developed a novel method for achieving cross-synthesis by using image style transfer methods.

More work was published by *Liu et al.* exploring the application of technologies given from the image domain for “mixing audio” [19]. This means, that this approach focuses on using audio as spectrograms. As the previous study solely focused on the technique by *Gatys et al.* this experiment explores two more approaches. While one is inspired by *Johnson et al.*, a convolutional autoencoder coupled with a VGG classification network, the other one uses an approach with (cycle)GAN (Generative Adversarial Network).

In this work they call Gatys' approach specifically “slow transfer”, as the iterative computation from gaussian noise was proven to be really slow. Different to the previous work by *Verma et al.*, the authors adapted a VGG network (1 input channel in first layer instead of 3) for the “slow transfer” method. This approach has also been used in *Gatys'* image style transfer. The transfer process is also similar to the previous work, as *Liu* uses a spectrogram initialized as gaussian noise in order to iteratively minimize the content loss in the higher layers, and the style loss in the lower layers. Setting this transfer process as a baseline model, they also adapted a faster style transfer method by coupling the VGG network with a convolutional autoencoder network. The purpose of this network is to take the content spectrogram as input and then outputting a spectrogram containing also the style features of a style spectrogram. Comparing the described work to other approaches, this proves to be similar to the one of *Ramani et al.* having a transformation network. The only difference is the second network, as here they are using a VGG classification network and no second autoencoder. Having the output of the autoencoder network (also called generative network), this one serves as the initial spectrogram on which the content and style loss is computed in the VGG network. The gradient descent then gets applied to the autoencoder network, resulting after few iterations, in a stylized spectrogram. The researchers have proven that their approach is faster than the one with gaussian noise. As mentioned before, for the third experiment they adapted a CycleGAN, which accepts audio spectrograms instead of images. In the image domain, this kind of network is able to apply style transfer to only a portion of the input images. Secondly the application of this method generates two images, as transfer is done in both directions, which means in case of audio that two new sounds are calculated.

As another result, this work showed that their approach generates the results in a shorter amount of time. For the purpose of evaluation, they listen to the outcome, but also apply objective mechanisms like visual assessment of spectrograms, consistency tests with classification and examination of signal clusters. Eventually the authors state that with the baseline approach features like the harmonic is not clear and high frequencies get discarded and that the faster transfer emphasizes on lower frequencies but is missing out on beginnings of the notes. With CycleGAN also the lower frequencies get emphasized while higher ones get discarded. The listenable results of each approach are provided online².

As the above mentioned approaches are working on single notes, the experiment of *Grinstein et al.* has been implemented for whole audio samples [9]. Within their work they were adapting several other approaches, with neural networks from the image domain. Besides neural networks, they also implemented a handcrafted sound texture model which was compared to the neural approaches. The latter one is composed of three sound processing steps, that in combination emulate the human auditory system. Taking a closer look at their work, especially on the neural networks, it can be said that they differ from existing ones in several ways. On the one hand, they do not use a random noise spectrogram, but use the content spectrogram which then gets stylized through their methods. On the other hand many audio style transfer approaches, are explicitly computing the result with a combined loss function, that incorporates both

²<https://www.xuehaoliu.com/audio-show>

style and content loss. *Grinstein et al.* do not make use of this concept, as they already initialize the future stylized spectrogram with the content spectrogram, like mentioned previously. On this target spectrogram, just the style loss gets optimized, as the content is already present. To mention here, they proved this method to have compelling results, as the global structure of the content sound is preserved.

In detail the authors investigated the use of three different network architectures for the purpose of audio style transfer. Concerning all three network types, they minimized the style loss on the content sound respective spectrogram. This style loss is equally computed as in *Gatys'* image style transfer approach. It is calculated by minimizing the error to a “style sound or spectrogram”, at specific layers in the network that extract the style. Via backpropagation the loss is minimized again at each layer, which results in a stylized content sound or spectrogram, after a few iterations. This workflow was applied to all three different network types and compared in the end. As their first network they used a VGG-19 network like *Gatys*, where the input spectrogram was replicated three times in order to match the input shape (RGB-like). By averaging all three channels in the end, they were able to obtain the final stylized spectrogram. The second network is called SoundNet, which is a convolutional network trained on unlabeled videos including sounds. This type of network operates on the raw waveform where no generation of spectrograms has to be done in advance. Finally a wide-shallow-random network was investigated with audio spectrograms consisting of just one layer CNN (like in the work of *Ulyanov and Lebedev* [29]). As the fourth and last method they made trials with a handcrafted sound texture model that emulates the human auditory system. Even if it is not a neural network, it consists of three layers doing cochlear filtering, envelope extraction with compressive non-linearity and modulation filtering.

Eventually the authors came to the following experimental conclusions: While using the VGG network no meaningful results could be obtained due to the noisiness. In comparison the SoundNet yielded more relevant results despite also containing some noise. Surprisingly the shallow random network performed best together with the sound texture model. For a better understanding, the results were provided online³.

The work of *Ulyanov and Lebedev* has to be mentioned here, as they are often referred to be one of the first that explored transfer algorithms for audio [29]. In fact, they took the architecture used in image style transfer by *Gatys et al.* and adapted it for usage on audio spectrograms. Rather than seeing the spectrogram as a picture, with the dimensions of frequency by time, they took the frequency values as channels for the CNN. The network itself was designed as a shallow network (1 layer) using 1D-Convolution with random weights. To obtain a final spectrogram containing content and style, an optimization is made on random noise, to minimize the loss values to a style and a content spectrogram. Instead of applying it on single notes, this approach also uses longer samples or music snippets.

2.3 Image Style Transfer

In the previous sections different works about neural audio synthesis as well as neural audio style transfer were explored and discussed. A lot of these works especially those

³<https://egrinstein.github.io/2017/10/25/ast.html>

proposing solutions for neural audio style transfer, took their inspirations from the image domain. For this reason this section provides a short excursion on relevant image style transfer works of *Gatys et al.* as well as *Johnson et al.* [7, 15]. *Gatys et al.* were the first to implement a system of neural style transfer applied on visual data. By using a convolutional network trained on object recognition and localization (VGG-19) with images, they were able to extract the texture as style but also the content of an image. They found that especially from higher layers in the network, high-level features of the images, can be reconstructed. This includes objects and their arrangements in the scene, without the exact pixel information, which will be used as content representation further on. Using a special feature space for texture synthesis, the style of a content image can be extracted by using the feature responses at certain layers in the network. By combining these two principles, respective style losses and content loss can be computed which will be used for the style transfer. The generation of the target image starts by initialization of a random noise image, on which those losses get minimized by using gradient descent. *Johnson et al.* developed an improved image style mechanism on the basis of the former methodology, that particularly shows improvements regarding computational speed. For the computation of content and style losses they use a VGG network with 16 layers pre-trained on image classification tasks. Here they added a special transformation network, that is designed as an autoencoder. This one takes a target image as input (content image) and synthesizes an image on which the style and content loss is calculated in the VGG network instead of random noise. Backpropagation was performed just in the transformation network, while the VGG network stayed fixed. As a consequence the transformation network produced a stylized image after training. By comparing *Johnson's* work to the method by *Gatys et al.*, it shows a significant improvement regarding computational speed but also yields promising results.

Previously described image domain methods significantly inspired the development of audio style transfer algorithms, which are presented in section 2.2. Finally, it should be noted that works by [24] and [19] adapted and applied the method of *Johnson et al.* while all others in section 2.2 mainly used the methodology proposed by *Gatys et al.*

Chapter 3

Approach

In the last chapter some approaches have been outlined and discussed, that have successfully implemented methods regarding the creation of audio signals using a neural network approach. As previously written, these have been mainly categorized in neural audio synthesis and neural audio style transfer. This work is mainly influenced from the area of neural audio synthesis, and can be categorized as such, as the methodology and workflow is strongly related to those works. Nevertheless regarding certain components, it is also influenced by the style transfer methods, despite not defining a specific content or style audio respective loss functions.

This chapter will therefore dive into the methodology and exact workflow of this works' solution, to the problem that also will help to derive the answers to the defined research questions. First a motivation should provide the reader with the intended idea and an overview of the applied methods, to get a general understanding of the idea (see section 3.1). Later on the separate steps and components that are needed to reach the desired functionalities, are described in detail, starting with the pre processing. Further on the ML-model (neural network) will be described, as well as the step that is done to synthesize new sounds. Further on, the required steps for (re)synthesizing a listenable audio as well as a description of the used dataset for training and also all experiments conducted later on are given (see chapter 4).

3.1 Motivation

Like mentioned in the beginning of this thesis, this work aims to explore the possibilities of machine learning techniques such as neural networks for applications in the audio domain for sound generation. This idea is mainly inspired by the idea of taking two distinct audio sources and mixing their characteristics in order to generate a new sound. As seen in the previous chapter, this idea is strongly related to the image domain, where the "synthesis" of a new picture based on two source images, is commonly known as image style transfer (section 2.3). This technique, having a content image to be stylised with a certain style from another image, would mean for the application in the audio domain, to have a sound style to be transferred onto a content or target sound. Such methods are specifically known as audio style transfer and can either be applied to single notes or also whole audio samples or songs. Having the principle of content and style

this would mean, that of one sound the global structure and rhythmical components get preserved while imposing style (e.g. the timbre) on it to generate audios. The details to these approaches have already been outlined in the previous chapter, when describing some existing work around this topic.

Neural audio synthesis is another method for neural sound generation, which does not apply the principles of style and content audio. In the previous chapter, some insights could be gained, how neural audio synthesis can look like, as well as how it can be achieved using different methods and neural networks. Most of those methods were showing promising results, either concerning the auditory quality but also the possibilities that arise in experimenting and designing sounds. Those methods were applying so called autoencoder networks most of the time, that can be used for dimensionality reduction of input data, as they have a so called bottleneck in the middle [12]. Because of this structure, the compressed data in this bottleneck, can be seen as a representation for essential features that either can be combined/interpolated or directly synthesized. To generate synthesized audio, the solutions described in section 2.1 took advantage of the “decompression part” of the network in order to generate audio data. The exact workflow and methodologies for sound creation have already been mentioned in the chapter related works (see chapter 2, sectionn 2.1).

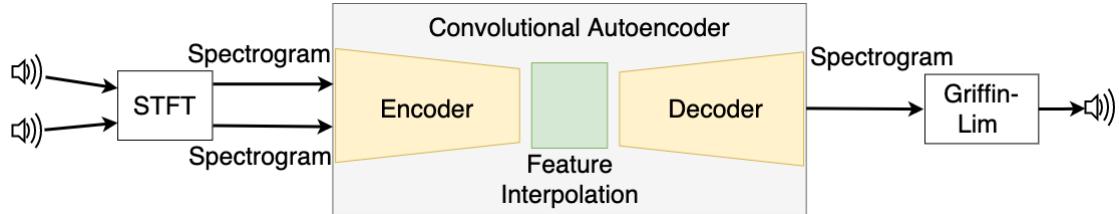
Out of those methodologies, when having the idea of using two instruments’ characteristics, to generate audio, the approach of *Engel et al.* [6] using convolutional and WaveNet-style autoencoders yielded the most promising and interesting results. This can be said especially in terms of output quality but also concerning its implementation and reproducability. With a provided interactive web application, the results of these solutions can be explored, whereas different sounds can be mixed based on a certain ratio. The results in the web application are based on the WaveNet-style autoencoder but according to the scientific article, the convolutional model (baseline) also provides strong results. Implementing an approach with a WaveNet-style network would also go beyond the scope, not at least as the computational costs would be too high. As also some audio style transfer methods, especially the approach by *Ramani et al.* [24], are using convolutional autoencoders, this kind of network was chosen to be preferable, to be applied in this work. How (convolutional) neural networks work, especially concerning autoencoders, will be described later on with special focus, on how those functionalities help to carry out neural audio synthesis, but also to gain general knowledge and a better understanding.

3.2 Overview

Based on the motivation and existing approaches, this work aims to propose a system that uses a convolutional autoencoder network, for the task of neural audio synthesis. This systems’ goal is to take two distinct audio samples as input, whereas the significant features of those are extracted and interpolated, to (re)generate a novel sound in the end. In figure 3.1 the general workflow of the toolchain is depicted in order to get an understanding of how this system is built up.

Starting on the very left, two audios are taken and have to be brought into a suitable representation for this type of network. As audio is in its raw form a time-continuous signal and the input for convolutional networks are of a different shape (e.g. images) some

Figure 3.1: Overview of the proposed solution.



pre-processing has to be done. In this case the short-term Fourier transform (STFT) is applied in order to generate a spectrogram, that shows the frequency spectra over time. The frequency spectra contain on the one hand the magnitude (power) of the frequencies but also the phase information. For this purpose, only the magnitude data is used, as recent publications stated that it contains the most descriptive data of an audio (spectrogram). The autoencoder model then takes the magnitude data as input, from which a compressed representation with the essential features gets generated by the lefthand (encoder) side. Having those features of two different audio samples, those get linearly interpolated, to generate one feature vector representing the “mixed” features of two instruments. This new vector gets passed through the righthand (decoder) side of the network, which regenerates again spectral magnitude data of the same dimension as the input. In order to obtain a “playable” audio sound, it gets transformed back into time domain with the Griffin-Lim algorithm [8] or the inverse short-term Fourier transform (ISTFT). The latter will be applied if there was no interpolation in the embedded space, as the phase information can be reused. Corresponding terminologies as well as a detailed insight into each step and its functionalities are given down below in the following points.

3.3 Pre-processing

Pre-processing is the task of preparing raw data for a specific purpose. Moreover it is an important component of machine learning techniques, with respect to training neural networks. Deciding which pre-processing technique(s) to use on the one hand depends heavily on the type of ML problem that has to be solved or even the training method that is chosen, but of course also on the type of data itself. As written before, for this work a neural network consisting of convolutional layers has been chosen to be applied to the problem of neural audio synthesis. As convolutional neural networks are known for image processing tasks, they also can be applied for audio data, which has already been outlined in recent works in this field (see chapter 2). In contrast to image data which most of the time has a 3D shape (width x length x RGB-colors), raw audio has a different structure in its data representation, as it is a time-continuous signal (1D-shape). In order to bring the audio data in a similar shape, it has to undergo some pre-processing steps. Some representations of audio that have an image-like shape and that got proven beneficial regarding neural audio synthesis, include e.g. log-magnitude spectrograms or Mel-spectrograms but also chromagrams and Constant-Q Transform as stated by [1]. Taking recent works into account, this work chooses to use the first

representation as this one also got used more frequently and had promising results. As for comparison in the experimental part of this work, also the use of Mel-spectrograms will be assessed and discussed concerning the synthesis task and the performance of the neural network.

As the practical part of the project to this thesis is implemented in Python, a special library was used for the pre-processing part. For this part the library *librosa* [28] was used, as it provides practical functions for audio processing, that were considered useful for this work. The functionalities of calculating spectrograms (STFT) but also transforming spectral data back into time domain to generate playable audio data (ISTFT, Griffin-Lim) are of special interest here.

3.3.1 Spectrograms and STFT

Spectrograms represent a 2D-representation of an time-continuous signal, which essentially shows the presence and change of frequencies over time. Like previously said, there exist different forms of spectrograms e.g. log-magnitude and log-mel spectrograms. Especially speaking of the log-magnitude spectrogram whose calculation is based on the short-time Fourier transform (STFT) and thus on the Fourier transform. The Fourier transform takes a frame of N values of an (audio) signal and transforms it from the time domain into the frequency domain. Generally said, that the bigger the frame, the better is the frequency resolution. What this means in terms of calculating the spectrogram, will get outlined shortly. What's also important to mention at this point is, that the result of the Fourier transform consists of an array of N complex numbers, which are mirrored around the middle. Every complex number in this array stands for a so called frequency bin in the signal. The real part of these numbers would represent the power/magnitude of this “bin” and the imaginary part gives information about the phase. Coming back to the frequency resolution, this for example means that when taking a one-second signal with a sampling rate SR and performing the Fourier transform with length $N = SR$ on it, this would yield an array of N values. The first value in the result depicts the signal's offset whereas all values from 1 to $N/2$ are the frequency bins with a resolution of 1 Hz per bin. This means that each of these bins shows the magnitude and also phase of each frequency from 1 to $N/2$ Hz. The ongoing values in the result show the same values except they are mirrored, as they depict the negative frequencies. Because of this behaviour, the second part can be omitted for further use. Now these values just show the frequency spectra of one time frame and do not incorporate more information about the change. To overcome this shortcoming the Fourier transform can be applied to a series of frames of the signal in order to obtain multiple frequency spectra over time that are depicted as a spectrogram.

The calculation of multiple frequency spectra over time is done via the so called short-time Fourier transform. This form of calculation is widely used for pre-processing of audio data for ML-tasks (see chapter 2). When applying this transform, a few parameters have to be considered, as those influence the result but also the quality for the later workflow. As *librosa* is used for the sound processing steps, the mentioned parameters are specifically concerning this library. One of the the most important parameters is `n_fft` as it specifies the actual length of the signal frame, on which the FFT (fast Fourier transform) gets applied. This parameter therefore influences the frequency- but

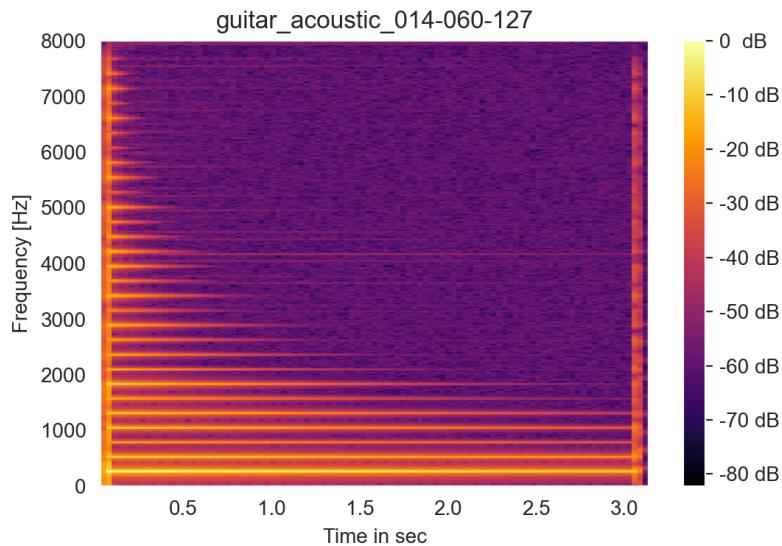
also time-resolution in the final spectrogram. Regarding the official Librosa documentation, this should be a value of a power of two, as it speeds up the computation of the FFT. Another important parameter would be the `hop_length`, which defines how many audio values are between the beginning of the first and the following frame. This means that when defining this parameter to `n_fft/2` this would yield a 50% overlap of the following frame. Modifying this parameter would mean to either increase or decrease the overlap and also the amount of time columns in the result, as more overlapping frames occur. The overlap of the frames is also coherent with the chosen window function for the STFT. As every time when the FFT gets applied to a frame, this one gets multiplied with a so called “window”. Multiplying a signal frame with a window has to be done, as the FFT assumes, that the transformed signal is periodic (repeating itself infinitely) [11]. This becomes problematic when the input signal does contain frequencies that may not directly fall into a frequency bin, due to the FFT’s frequency resolution. Due to the assumed cyclic continuation the Fourier transform will ‘think’ that there is a discontinuity and will spread therefore the power over the spectrum. There are various window functions such as “Hann”, “Hamming”, “Blackman”, etc. which start at (almost) zero, rise to a maximum in the middle but fall again to (almost) zero at the end (symmetric). Multiplying the signal frame with such a window function helps to overcome this issue as it removes the discontinuity. Which window to chose depends on the use-case of the application, whereas throughout this work a “Hann” window was chosen. It got chosen, because it is good suitable for most sound processing tasks. Moreover it reaches zero at both ends and thus eliminates the discontinuity problem. Coming back to the relation with the window overlap, if no overlap would be used a lot of information of the signal would get lost. This is because when multiplying the signal frames with the window functions, this would result in very small or zero values at the beginning and end of the frame [11]. When having overlapping frames this issue would be corrected. Important here is again the amount of overlap as this is dependent on the window and its wideness. Using the “Hann” window a common value for the overlap would be 50% which was also considered throughout this work. This value is also beneficial later on when performing the inverse transformation, back into time domain, but more on that in section 3.6. Beside of these parameters some more exist, for example for specifying the padding of the signal whereas in this work a constant padding on both sides of the signal has been used, which is also the default setting.

Now having knowledge about the STFT and its parameters, it can be applied to a signal to generate a spectrogram. For example by using an audio sample with a sample rate of 16 kHz and applying the STFT with a `n_fft` value of 1024 and a `hop-length` of 512 this would result in a spectrogram with a frequency resolution of 15,625 Hz and time resolution of 64 ms. As explained before, the values of the result consist of complex numbers which contain the magnitude but also the phase at each frequency bin. By setting this result absolute, or calling the function `librosa.magphase(spectrogram)` the real magnitude data can be obtained, whereas the latter also retrieves phase information in a separate vector. The magnitude here displays the energy values of the spectrogram, whereas for further processing and also to be better displayable those get converted into a dB-scale¹. This function also takes a reference value that is set to 0 dB, which in

¹normally the magnitude would need to get squared to obtain the power, but in this case magnitude without squaring was taken

this case will be the maximum value of the magnitude spectrum. As for post-processing when converting the dB-scaled data back into energy, also a reference value is needed, this one gets preserved, in order to get the same scaling as in the original input. Finally when having the log-mag spectrograms in dB, those were considered for the training of the neural network afterwards. An example of a log-mag spectrogram can be seen in figure 3.2. This spectrogram shows the frequency spectra of a guitar sample over time. It can be seen, that at the beginning and at the end there are broadband spectra, which represent the guitar stroke (transient) at the beginning and the noise of damping the strings at the end. As also the phase information was obtained when calculating the magnitude data, this one was also preserved next to the energy reference value for the recreation of signals, as it is needed there (this will mainly affect the recreation of single samples, without interpolation in embedded space, but more on that later on).

Figure 3.2: STFT log-mag spectrogram of a guitar note.



This section describes the general workflow of the pre-processing from taking a signal and converting it into a spectral representation. This workflow is a basis on which different experiments with different parameterization (size of `n_fft`, etc.) of the calculation of the spectrograms but also with additional steps (log-mel scale, additional framing, etc.) are being made. Those steps will get mentioned later on in chapter 4 when describing the experimental part of the thesis.

3.4 ML-Model

The main or core component of every machine learning project is of course the model itself, as it achieves the main task of prediction or inference to a given problem. Those models exist as different technologies that perform regression tasks or even classification tasks. Dependent on the use case, but also the kind of data that is present, different

models are better suited or not. To enumerate technologies, there exist the KNN algorithm, Decision-Trees, Random Forests, Support-Vector-Machines (SVM) but also neural networks which can be applied in a variety of use cases. Especially the latter, the neural networks, are able to achieve a variety of different tasks, as they are highly adaptive regarding their topology, used layers, but also their size and shape. This variety of different tasks spreads across different domains, including images and audio.

3.4.1 Neural Networks - Introduction

Generally said a neural network can be seen as a graph of connected nodes with numeric values that can achieve transformations between patterns using message-passing algorithms [16]. Those nodes are commonly structured in layers, where there especially exist certain nodes or even layers that are seen as input nodes/layers and some as output nodes/layers. Between the input and the output there can also exist so called hidden layers, expanding the depth of the network. The links between the nodes, that are also called neurons, are connected via links, that are parameterized with weights, that get optimized using learning algorithms. Each neuron receives its weighted input (activities) of its connected predecessors, which get converted into a single output that gets broadcast to all its connected successors [13]. The latter involves a so called input-output function which is also commonly known as activation function (e.g. ReLU, Sigmoid, Softmax, etc.). Important to know is that the weights on the connections define how much this value influences the input of the connected node. When a neural network is trained to a specific problem (e.g. classifying certain images), using predetermined training data, the output of the neural network is compared with the desired one, resulting in a certain error (metric). To minimize this error, the weights in the network are adapted by back-propagating the error through all the layers, to the beginning. On this way it changes the influence of certain connections and therefore the overall outcome. This procedure is repeated on all training data over several iterations, until the error becomes low to produce the desired output. The initialization of the network and its weights is often random, which also means that every training run starts and progresses differently.

Depending on the problem at hand, neural networks can be trained using labeled data (supervised) but also just by minimizing a cost function (unsupervised) [21]. More details on the learning will get mentioned later on when explaining the model itself.

3.4.2 Convolutional Neural Networks

In this approach, due to the promising usage in existing solutions, a convolutional neural network has been chosen as the model. Convolutional neural networks are a type of network, that get primarily used for tasks in the image domain for example to recognize patterns in pictures or classification, but also like seen in chapter 2 for image style transfer. They have the advantage over traditional neural networks, that they can deal with the dimensionality of pictures (width by height by colors/depth). The layers containing convolutional nodes have kernels as learnable parameters [21]. Those kernels, if taking a 2D-convolutional layer, are normally small in width and height (e.g. 3x3) but span the whole depth (channels) of the input (in the case of RGB pictures depth of three). Those kernels calculate the scalar product for each value contained in the kernel and the input map. This yields, having e.g. a 3x3 kernel operating on a 3x3 field, in a single

value. This value is the weighted sum of the kernel's values and from those of the input vector. This operation will be applied to each 3x3 field along the spatial dimension of the input, resulting in a smaller activation map. One can also apply padding around the input to preserve the dimensionality. Furthermore a stride can be defined, which defines how much these convolved fields overlap, as using a bigger stride would result in a smaller overlap. Having also a smaller overlap would result in a much smaller activation map. For example taking a 7x7 input, by applying a 3x3 kernel with no strides and also no padding, this would yield an output field of 5x5. With a padding of 1x1 the output would be of the same size. Finally if the stride would be 2, the output field would be of 3x3. Through training of the network and back propagating the error, those values in the kernel get adapted in order to learn certain important features or patterns on which a classification or pattern recognition can be made (easier). As it got described for 2D convolutions, depending on the input dimensionality it can also be applied as 1D convolutions or even 3D convolution.

With the knowledge that such convolutions are successful on image data, those can be also applied on audio provided in the shape of a spectrogram. As described in the previous section (see 3.3) spectrograms can be described in a “picture-like” shape having the dimensions of frequency by time with a depth of one. Speaking of that, the spectrogram can be seen as a grey-scale image. The energy in different frequency bands over time with its variations, can be seen as “recognizable” patterns. Taking a 2D convolution, the kernel (e.g. 3x3) takes a 3x3 frame of frequency by time which results in one value. Summed up, this results in an activation map smaller or equal (if zero padding is used) than the input. After training on different samples and iterating several times, this activation map would contain the most significant features or characteristics of the spectrogram for example when training on classification. Depending on the chosen hyperparameter for the amount of output channels (depth), the resulting activation map can be of depth one or even deeper.

As mentioned before, each neuron in a neural network has an activation function, which is also the case for convolutional neural networks. As there exist different kinds of activation functions, for this approach mostly ReLU (Rectified Linenar Unit) activation functions got used but also LeakyReLU. Additional to this activation function, Batch-Normalization gets applied after each convolutional and ReLU non-linearity component. The choice is mainly based on already existing approaches like from *Ramani et al.*[24] and *Engel et al.*[6], as it was proven to yield promising results in combination. According to *Ioffe and Szegedy* [14] applying Batch Normalization also improves the training speed (number of iterations/epochs) and enables the use of much higher learning rates. Furthermore it acts as a regularizer, so that overfitting-reducing technologies, such as Dropout, can be omitted.

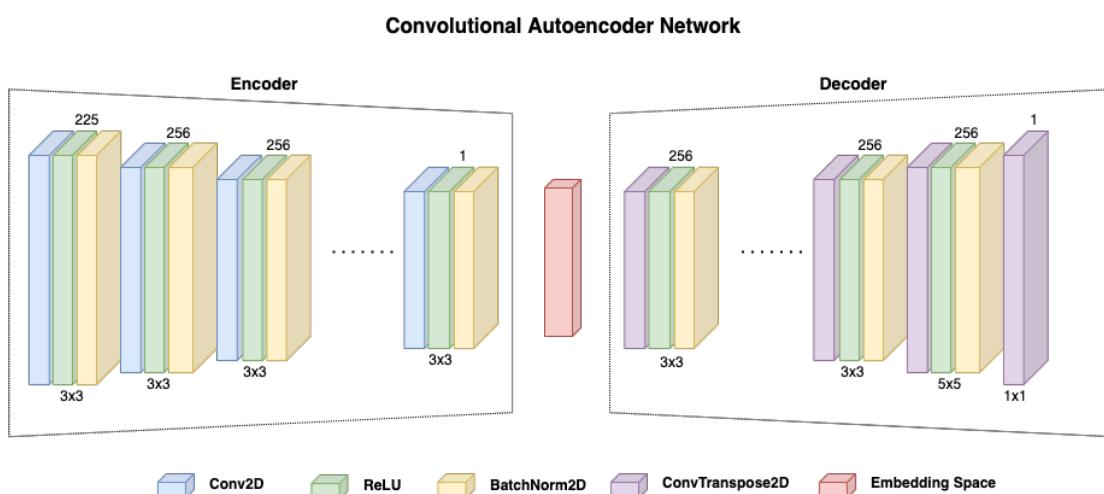
3.4.3 Autoencoder

Neural networks exist in different shapes and compositions, depending on the desired work it should fulfill. Whereas some networks e.g. for classification of a given input, might reduce the width of the layers towards the end, some are designed to have a so called bottleneck. This means that this kind of network contains a smaller central layer then the input, but at the end again a bigger layer (eventually same size of input layer)

[12]. Those networks are called autoencoders, in which the first part of the network (until the smaller central layer) gets called “encoder”. The second part beginning at the small central layer, that is getting bigger towards the end, gets called “decoder”. These parts are called this way, as for the encoder part, it “encodes” the high-dimensional input data to a low-dimensional representation (output of small central layer). The counterpart is therefore called “decoder” as it decodes the low-dimensional data, to bring it again to a higher-dimensional representation (mostly same as input). The lower dimensional output of the small central layer, can be named differently, as for example “code”, what *Hinton et al.* does when describing the concept of autoencoders in his publication [12]. Some other common names would be e.g. latent space, embedding space, encoding or embedded data. This principle of dimensionality reduction therefore means, that the encoder part extracts the most important or characteristic data, from which the decoder is able to reconstruct the input data.

In the related works chapter 2 a few works have made use of this principle to extract “characteristic” features for audio data and synthesize audio from it, by altering them. This knowledge encouraged this work to also make use of this principle to synthesize audio by using autoencoders. Especially when having the extracted audio features as encodings, those can be modified (easier), to in order synthesize novel sounds from them using the decoder part of the network. Combined with the knowledge to apply convolutional networks to audio spectrograms and the advantages of autoencoders to extract features from the input, the model in this work is designed as a convolutional autoencoder. Also the works of *Ramani et al.* or *Engel et al.* made use of these kind of network to generate novel sounds. Figure 3.3 shows the basic autoencoder network structure which is used throughout this thesis. To be mentioned, the whole implementation of the neural network for this work has been achieved with the deep learning library PyTorch [22].

Figure 3.3: Basic autoencoder structure used throughout this thesis.



In figure 3.3, the depicted autoencoder, gives an insight of which components and layers it is composed of. The amount of layers as well as the parameters shown in this

sketch, are just an example, as throughout the experimental part, those get varied. Therefore, the dotted lines in this figure act as a placeholder for possibly more layers on each side of the network. The numbers above the individual layers represent the amount of output channels (depth) and those on the bottom represent the kernel size. It has been mentioned that in this approach the convolutional layer gets equipped with ReLU activation functions and batch normalization. This also gets visualized in this figure, where each layer is shown as a combination of sub-layers that incorporate those three components. Speaking of the composition of the layers, it can be seen in the decoder part, that instead of a convolutional layer, a convolutional transpose layer is used. That is because of the nature of convolutions, as their output is always smaller or equally sized, which has been mentioned under point 3.4.2. As the decoder part generates output that is bigger than its input, the layers have to perform upsampling. In convolutional nets this is typically done via the convolutional transpose layer.

This transposed convolution is not the reverse operation of a convolution, but it is more of a operation to recover the shape from the convolutions input. [4] Taking the convolution example from above, by taking as input a 5x5 field and applying the transposed convolution with a kernel of 3x3, this would result in a field of 7x7. The equivalent of this operation would be a convolution with the same kernel, on an input of 5x5, with 2x2 padding (padded input = 9x9), as this would result in a 7x7 output too. When in the convolution, striding was applied, this also works for the transposed convolution. Again when having the 3x3 input, by applying a 3x3 kernel with a stride of 2, this again results in a 7x7 output field. To be mentioned, the striding parameter for the transposed convolutions defines how much zeros are added between the values of the input. This means, when taking the previous example, that the 3x3 input gets one column respective one row of zeros inserted after each value to result in a 5x5 field.

With this knowledge, convolutional transpose layers are best suited to be in the decoder part of convolutional autoencoder networks. Coming back to the architecture of the convolutional autoencoder in figure 3.3, those convolutional transpose sublayers are also coupled with ReLU activations and batch normalization, with an exception to the last layer.

Having this kind of autoencoder for this approach and the idea to extract features of spectrograms using convolutions, this models task is to encode and decode spectral audio data. The network is therefore configured to produce an output of the same dimensionality as the input of the encoder. Like mentioned above, the output should be a reconstruction of the input, that gets inferred by decoding the extracted features in the embedded space. To achieve this goal of reconstructing the input, the network has to be trained through minimizing a specific error. In some cases also maximization is desired, but is dependent on the error score and the desired outcome. In the case of an autoencoder, this works by comparing the output of the decoder with the input of the encoder, by calculating the difference. Generally said, depending on the outcome and the goal of a training, there exist different metrics like mean squared error (MSE), root mean squared error (RMSE), mean absolute error (MAE) but also more specific formulas. Throughout the literature, those error functions are also often called cost (functions) or loss (functions). For the experimental part of this work, the choice has been made to use MSE as the error metric, as it has been used successfully by some existing works.

To minimize this error an optimization of the network has to be made, which will be done through backpropagating the error score. Through this step, the parameters (weights, bias, convolutional kernels, etc.) get adapted according to the error, which in the best case improves the score and thus the output of the network. In this work the output of the network are reconstructions of spectral data. As having the encoder-decoder structure, this output data is reconstructed through the decoder, which takes the encoder output as input. As mentioned before, this generated data of the encoder is a compressed vector, consisting of the most significant information extracted from the input. Depending on the configuration of the network, especially the encoder part, this vector can be of different sizes. This is because when using more convolutional layers with strides, the output of each layer gets more downsampled, which results in a smaller encoding. As a consequence, the decoder part has to learn to regenerate spectral data from this encoded vector. This vector therefore is a compressed representation of the models input, which contains the most significant characteristics of this specific spectrogram and further on, of the sound. All in all regarding the learning, it can be said, that the encoder learns to extract essential features, from which the decoder learns to regenerate the input as well as possible.

Altering those encodings has therefore the consequence, that the output of the decoder is different, and thus will be a novel, synthesized sound. More on how it should get altered is discussed later on. Choosing the right compression is also an important part that influences the outcome and thus the quality of the desired model output. Too much compression may lead to the fact that the decoder has too little information to infer the desired spectral output, which in order results in poorer quality of the resulting audio. On the other hand having less compression results in embeddings being not significantly smaller than the input, containing less important data too (e.g. noise). It also may become more difficult to alter those encoded vectors. Throughout this work, different amounts of compression have been applied in the experimental part, which get discussed later, including the impacts and observations made on that. Knowing those properties and behaviour of the autoencoder model strengthens the idea of applying autoencoders for the use of neural audio synthesis.

3.4.4 Optimizer

Coming back to the training process, where the network gets optimized, in order to minimize a certain error function. For this optimization, different strategies exist, where hyperparameters such as learning-rate or weight-decay play an essential role. Those optimizations are on a large scale, stochastic gradient-based techniques, to which algorithms such as stochastic gradient descent (SGD) or Adam can be counted. Those algorithms influence and improve the convergence which means to find a minimal error. Throughout this implementation the Adam optimizer [18] has been chosen, as it is used widely in recent publications where promising results could be achieved. Also regarding the training process in this work, Adam optimizer has proven to be advantageous, in contrast to SGD. The parameters that have been found to have the most impact on the optimization process during training, are the previous mentioned learning-rate and weight-decay. To be mentioned, the learning-rate specifies how “fast” the model actually learns while weight-decay works as a penalty for the weights optimization to prevent

overfitting. If the learning rate is chosen rather large, then the network learns faster, but because of its big steps or jumps, it could miss the optimal solution in the solution space. Also it could happen that when a local optimum is found, that it “jumps” out again. In this case a smaller learning rate would be desirable, as it makes smaller steps. Choosing it too low would end up in a slow training where also large areas of the solution space are not visited. The latter leads to a training process stuck in a local minimum. Therefore it is important to choose the right size of learning rate, but this issue depends also on the problem size and type of network. In this work different learning rates have been applied throughout the experimental part where also different findings could be made, but this will be shown and discussed later on in this thesis.

With this knowledge, it can be stated, that a high learning rate could be advantageous at the beginning of the training process, in order to rather find a global optimum and explore the solution space. In order to prevent jumping out of a minimum, a smaller learning rate would be desirable later on in the training. For this case, there exist some mechanisms to decrease the learning rate later on in the training, especially when detecting oscillations of the error due to a too large learning rate.

In this work’s implementation, for the start of the training a specific starting learning rate has been set, while throughout the training it gets adapted, when no more optimization and eventual oscillation gets detected.

3.5 Synthesis of novel sounds

Having now covered the important properties of the pre-processing but also of the applied machine learning model, this section explains the methodology to synthesize novel sounds. In the chapter 2 where related works got discussed, an insight could be gained, on how different works tried to synthesize audio with their neural networks. For example *Colonel et al.*[2, 3, 17] suggested in their works, to synthesize novel sounds through directly activating the innermost layer. This means after training, for the sound creation process, the encoder part gets omitted. For example in the case where they had a network with 8 neurons at the encoders bottleneck, 8 different values could get defined, within a certain range. The decoder part then created, based on its training on recreating spectrograms, spectral data that got converted back to time domain to form a synthesized playable sound.

3.5.1 Interpolation in latent space

Having this as one possibility, to use in particular autoencoders as a tool for audio synthesis, there also exist some more interesting approaches. One of those got applied by *Engel et al.* [6] but also *Roche et al.* [25] where they also made use of the latent space encodings. Contrary to omitting the encoder part, those works aimed to utilize the whole network, as no direct activation of the innermost layer is considered. Instead the authors proposed to take the encoded values of two different audio samples (possibly two distinct instrument) and combine them via linear interpolation. The interpolation process yields a vector of interpolated values. This new vector can be seen as a modified encoding, which gets processed by the decoder part, resulting in spectrogram-like vector. The result of this process is then a synthesized spectrogram, that aims to contain features

of the two input sounds.

As the latter methodology corresponds the most with the initial idea, to synthesize audio based on the characteristics of two instruments (e.g. guitar and synthesizer), this method was chosen to be implemented to carry out experiments on the creation of novel sounds.

To explain this method in more detail, the encoded vectors of two audio samples serve as the basis. Then each of this vector is taken, to interpolate a value that lies on a linear line between the value at a given index in one vector with the corresponding value of the same index in the second vector. To mention at this point, those vectors are of the same length. The result then is a new array containing the interpolated values of those two vectors. This procedure gets repeated for each encoded output of one sample with the corresponding output of the second sample. Knowing the fact that those encoded values represent the most significant features and probably the characteristics of the note, the result can be seen as a combination of those characteristics. Decoding those will then end up in having a spectrogram containing the characteristics of both instruments. How this procedure and its results look like, with the actual experiments, is shown later on in this work.

3.6 Post Processing

It has been discussed, that regarding neural audio synthesis, the audio data can appear in different shapes. As there exist approaches, that focus on time domain signal like the WaveNet-style autoencoder from *Engel et al.*, there also exist those who operate on spectrograms using convolutional networks. As mentioned before, this work emphasizes the use of a convolutional autoencoder, that takes spectrograms or spectral data as input. In the case of this autoencoder, the output is of the same shape as the input and therefore also a spectrogram. To generate again a playable or listenable audio, this one has to get converted back into time domain. There exist many different methods, to achieve this, while this also depends heavily on the data that is present. As stated in section 3.3, when spectrograms are calculated via the STFT, the output vectors are complex valued. To be mentioned, that without modification or further utilization, via the inverse STFT (ISTFT), this result vector can get converted back into time domain without loss. For achieving this task, it is necessary that the magnitude but also the phase information has to be present, combined in a complex number. As this autoencoder just operates on the magnitude data of spectrograms, there is no phase information present in the output of the network. At this point, multiple ideas can be applied, depending of what is done throughout the process. This means that when there is no modification in latent space, i.e. value interpolation, the original phase information can be reused while applying the ISTFT. As mentioned previously the phase information gets preserved for exactly this case.

In the second case, where modification steps are performed, like here the interpolation of two sounds' embeddings, there is no phase information present that can be used. To overcome this issue, there exist techniques that can approximate the phase information. One of those techniques, which is also probably one of the most prominent ones, is the Griffin-Lim [8] algorithm that tries to estimate the time domain signal based on just the magnitude data. With this algorithm the phase gets randomly initialized, and with

alternating forward- and inverse STFT, estimated. For the calculations of the audio signal, again the python library *librosa* [28] has been used, as it provides the inverse STFT but also implements the Griffin-Lim algorithm. According to the documentation of *librosa*, a so called “fast” Griffin-Lim algorithm is applied, which got developed by *Perraudin et al.* [23]. The difference here is, that this one utilizes a additional momentum parameter which helps to accelerate the convergence of the estimation.

Some further note here, as when the pre-processing steps (see section 3.3) have been examined, that some certain parameters have to be taken care of. These are especially the `n_fft`, `hop_length` but also the `window`. It has been mentioned, that e.g. the `hop_length` was chosen to be half of the `n_fft`, to ensure a 50% overlap of the STFT frames. Furthermore the “Hann”-window was chosen, to be multiplied with the signal frames, to avoid discontinuities. Again, those also appear in the inverse STFT as well as Griffin-Lim. In order to obtain the best result, it is important to apply the same values with those parameters. As the inverse Fourier transform gets applied to each vector in the spectrogram, this results in single time domain frames that have the length of `n_fft` (original length). With the knowledge of the window and the overlap, those single frames, get “overlap-added”, which results again in the full length audio sample.

Before applying the the inverse calculations on the autoencoders output, it has to be considered that the autoencoder works on the db-scaled magnitude. As a consequence the output therefore is also db-scaled. To apply the inverse calculations, the output becomes energy again. For this case a reference value has been preserved from the pre-processing stage, in order to obtain the (almost) same scaling again in the output signal. For the experimental part some more steps also got applied, like scaling the energy according to the average energy present in the original spectrograms. This step should also correct and improve the resulting sounds. More on that in chapter 4, but also later on when examining and discussing the results.

3.7 Dataset

As machine learning models including neural networks, have to get trained, in order to deliver accurate results, data is needed on which it should get trained on. To get compelling results, it is not only important to have an appropriate model configuration or pre-processing chain. The choice of an appropriate dataset therefore is also of high significance. As seen in related works, datasets can either be self-generated or taken from a publicly available data source. As in the case of this work, the model operates on audio data, a dataset of musical notes is desirable. Generating sufficient data, by oneself, is a task that would take a significant amount of time. Not only as this dataset preferably should contain a large amount of audio samples, those also should be highly diverse such as different instrument sources or different pitches. Not only is it an advantage for the training of the neural network to have lot of samples and diversity, as it helps to improve the learning process and generalization within the neural networks. Moreover it is also advantageous for this kind of work, as when having many different instrument sources and available notes, more interesting combinations with regard to audio synthesis can be made.

3.7.1 NSynth Dataset

Exactly for this kind of approach a large dataset consisting of instrument samples called “NSynth”, has been made publicly available by *Engel etx al.* [6]. This dataset consists of a total of 306.043 musical notes that have a unique pitch, timbre but also envelope and has been created for the idea of neural audio synthesis. This amount of musical notes incorporate monophonic audio snippets, sampled at a rate of 16 kHz, of 1.006 different instruments. Every note is of a specific pitch, ranging over every note of a standard midi piano (21-108). This results in having 88 different pitched notes, in the best case, as not every instrument is capable of producing all different pitches. The average amount of pitches is therefore, according to the scientific publication, 65.4 per instrument. With more detail, each audio sample belongs to a certain instrument family which could for example be a keyboard, guitar, organ, bass, brass and so on. Further on they can be distinguished by their source, being either produced acoustically, electronically or synthetically. All these specifications make this dataset highly attractive for this kind of work, besides being publicly available² for free and already successfully used for neural audio synthesis. Regarding the use for machine learning tasks, on their website, they provide the dataset already split up into a training, validation and test set which do not overlap at all. To specify, the training set consists of 289.205, the validation set of 12.678 and the test set of 4.096 examples. For this purpose, it has been chosen, to take these splits as they are, for the training, validation and finally testing stage.

²<https://magenta.tensorflow.org/datasets/nsynth>

Chapter 4

Experiment

This section describes the experiments conducted in order to be able, to answer the defined research questions. In the previous chapter 3 the general methodology used technologies got described, to gain an understanding of the implementation and its components. Those can be parameterized in order to influence the outcome of this work. Furthermore some additional steps can be introduced that have a significant impact on the result. Based on this knowledge and implementation, some experiments were conducted. Those experiments, should deliver some interesting insights, on how different configurations influence the workflow but also the final result being a synthesized audio. Those experiments span almost every stage from the pre-processing until post-processing. This section therefore describes, how the proposed methods were utilized, in order to answer the questions.

4.1 Implementation Environment

In the previous chapter it has already been mentioned, that this approach was developed in python using specific libraries. To shortly mention, for pre- and post-processing the python audio-library *librosa* [28] has been chosen, as it provides all necessary functionalities that are needed for the approach and experiment. For all steps regarding the neural network model such as configuration, training, inference etc., *PyTorch* [22] has been utilized. The project though has been implemented and applied on two different machines, depending on the task that had to be done. Generally speaking, the toolchain compromised of all stages, has been developed on a local machine running python, except the training itself.

Speaking of that, the training has been mainly performed on a remote “jupyter-notebook” that has access to high-performance GPU resources. Not at least, as in the case of training a convolutional neural network, this is a rather time and computational power-consuming task. Of course, this not only depends on the kind and complexity of the network, but also on the amount of data that is used for the training. Furthermore using GPU-acceleration means to significantly have more computation power and speed. Using the local machine, just the CPU could be utilized for training, which would mean that training is done sequentially and thus significantly slower while also the local machine has to be awake constantly. For the training on the remote instance,

the pre-processing also was done there as the data is directly loaded there. As just the training was performed remotely, all other steps, including the evaluation towards audio (re)synthesis having the trained model, were done locally. Not at least, as no time-consuming tasks had to be made, but also as it was more convenient, as the remote service is not always accessible.

4.2 Training

The training, as mentioned previously, was performed on a remote “jupyter notebook”-service with access to a GPU. As outlined in the previous chapter, for the whole experiments, the NSynth data set proposed by *Engel et al.*[6], was used. This dataset is already split into a training, validation and test part. For the training on the remote notebook, the training and validation data set was used therefore, which in order also was pre-processed there. As a side note, in the beginning of the project, the training was just held locally, with a small subset of the already small test set (mostly of one instrument). This was just done to make a low-level proof that the autoencoder model can produce meaningful results.

4.2.1 Training configuration

To take a closer look onto the training process, this one consists of several important stages and components. First of all the PyTorch-model, defined as a class, was initialized. As a metric is needed, to measure the error of the output, the mean squared error (MSE) was utilized. This error metric calculates the difference between all values of the desired and actual output, squares them, and takes the average over all. Furthermore to optimize the network, as explained before, the Adam optimizer got applied, in which the (starting) learning rate, but also the weight decay was defined. The right learning rate depends here heavily on the amount of training data but also complexity of model. In the case of this work, this means it is in the range between $1e - 5$ to $1e - 7$. In chapter 3 it was also mentioned, that a technique to minimize the learning rate, during the training process was utilized. Generally speaking, this function called `torch.optim.lr_scheduler.ReduceLROnPlateau(...)` reduces the learning rate, by a given factor, if within a certain patience period (epochs) no optimization was detected. This method, improves the training process as further convergence can be achieved. Of course for the training process, the pre-processed training and validation data set had to be loaded. To mention, the pre-processed data was calculated and stored on disk, in advance, to not always have to run through it. For the training it therefore was loaded and brought into the desired shape for the corresponding model. This shape was varied throughout the experiments, to observe its impact on the training process but also quality of the output. As this also depends on the chosen network, this gets described later on in this chapter. For the training process, this was done with the training dataset, but also on the validation dataset. The latter is an important step during the training, to validate the models performance on never seen data. The data in the right shape, had to be converted to a tensor and in further notice, to a (custom) dataset object. Finally a so called “DataLoader” had to be initialized, either for the training but also for the validation. In this DataLoader the batch size can be specified, which is an important

parameter regarding the training. Throughout a few batch sizes, have been tried out, whereas in the end, the preferred batch size was 32. This means that the input data, was portioned in equally sized chunks of 32 tensors, that were fed into the network at once. Further on it can be set, that after each training epoch, the dataset was shuffled. Setting this to true, the samples in the batches are in a different order and constellation. Otherwise, the batches consist always of the same data in the same order. Throughout the experiments, this setting was proven to be advantageous regarding the convergence of the model, as the error could be more minimized (see chapter 5).

4.2.2 Training Execution

Having the configured dataloaders, which are also iterables, it is possible to iterate over the batches of the dataset. Before the training begins, a number of epochs has to be set, which defines how often the training should be performed on the whole dataset. In advance it cannot be said, how many epochs are needed, therefore a preferably big number was chosen e.g. 1000. To clarify, the training never ran until the final epoch, as it got stopped at a point, where the result is sufficient (more on that shortly). In each iteration, the corresponding batch of data got input into the network with a forward pass which calculated an output. This output was compared to the desired value, which in this case was the same as the input, and further on the error was calculated. Further on the gradients got computed and optimization via parameter update was done. The loss value was added up, and after all batches are done, the average loss got computed, to show the progress.

After all train batches have ran through and optimization was done, the model was validated using the held out validation set. This one was equally batched, and ran through the same process, except there was no optimization. Here the error was just calculated to see, how the model performs on data that was not used for training and therefore never was seen before, by the network. This technique helps to prevent to overfit the training data, which would be signalized through an increasing validation error despite training error becomes smaller. Having the validation error after each epoch, this value was used for the learning rate scheduler to decrease the learning rate if the error does not decrease in a certain period.

To ensure to have a sufficient trained network, the error scores got observed periodically. If the validation score did not improve more, i.e. the network convergence stagnated, and was sufficient, the training was stopped. Important to know that after each training iteration, which is also called epoch, the state of the model got saved, for further use. As also the scores for each model state were known, it was therefore possible to determine the best model, having the lowest MSE-score on the validation data. This one then was further tested and analyzed towards the applicability for audio (re)synthesis. Those further steps, which do not involve training, as well as the pre-processing regarding the evaluation data, were performed locally.

4.3 Initial Experiments

In the beginning of this project, initial experiments, were done in order to make a very basic proof of concept implementation. Like mentioned before, those implementations,

where entirely held on the local machine, not at least as there was no access to GPU-accelerated training. This was also possible, due to just taking a small subset of the test dataset. Of this test dataset, samples of one instrument (*keyboard_synthetic*) were taken out, and considered for test-wise training.

4.3.1 Whole Spectrograms as Input

Pre-processing

For the first experiments, all the samples of the subset, have been converted to log-magnitude spectrograms as a whole. As already known, those samples have all the same length of 4 seconds. Some samples are padded with zeros, as those do not contain audio data over the full length. As parameters for the STFT a `n_fft` of 512 and `hop_length` of 256 were chosen, this resulted in spectrograms having 257 frequency bins with a resolution of 31,25 Hz. Regarding the time-resolution it could be said, that each frequency vector represented 16 milliseconds of the original signal with a 50% overlap. This procedure resulted in spectrograms having a dimension of 257x250.

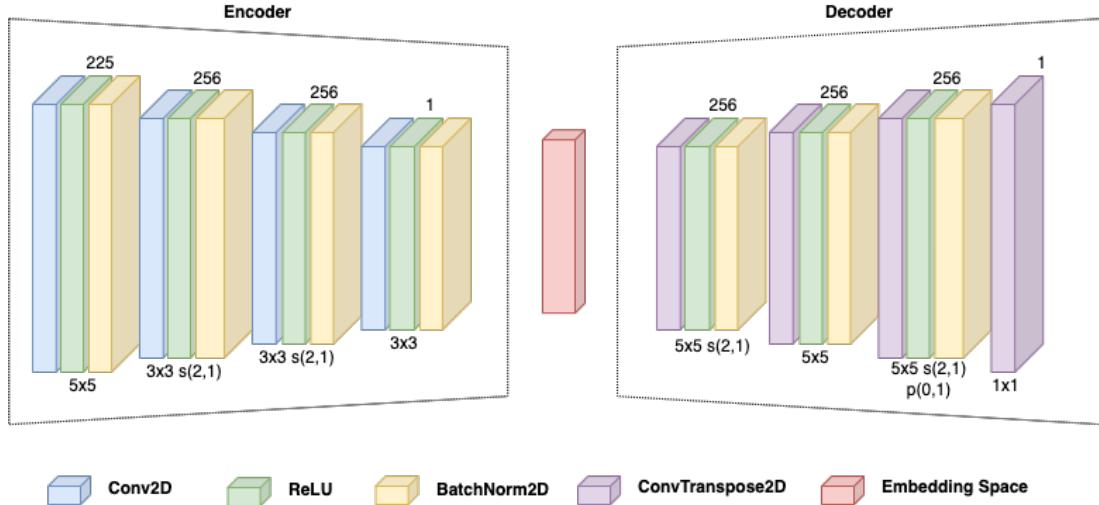
Model and training

As the spectrograms can be seen as grey-scale images, and thus 2D convolutions were applied, a third dimension was added resulting in 257x250x1 spectrograms. Finally to form a dataset ready for training, all the spectrograms of the mentioned subset were concatenated to a 4D array, which was converted into a tensor and subsequently into a dataset. The following training, was performed on 80% of this subset, whereas the other 20% were used for evaluation. Regarding the duration of the training, this was held for a short time (~20 epochs).

To also mention the model configuration, this consisted of 4 convolutional layers in the encoder including ReLU activation and batch normalization. The decoder part had therefore also 4 layers, which contain respective convolutional-transpose layers with also ReLU and batch normalization except the very last layer. Another important detail was, because of the shape of the input, that the number of input channels had to be 1. Throughout the network, this parameter got varied, whereas in the first layer an expansion was set to 225, subsequently in the second to 256 output channels. To the innermost layer it was again reduced to 100. In the decoder part, again an expansion was made, whereas at the end it was reduced to 1 channel again, in order to match the input size. The next graphic (fig 4.1) showes the autoencoder that was used for the initial experiments. To be mentioned the kernel-size (e.g. 5x5), strides (s) and padding (p) here is already the configuration for the succeeding experiment. The channels, the general structure as well as amount of layers was the same for all initial experiments.

Evaluation

Regarding the evaluation of this first model, it had been tested out on the remaining 20% of the small subset. These first experiments, consisted just of evaluating the outcome of the decoder, by passing single spectrograms through the network. By this, the ability of the network to recreate audio spectrograms was proven, which further on serves as

Figure 4.1: Initial 2D convolutional autoencoder.

a base for the next experiments. Regarding the final outcoming sounds, the original preserved phase information was reused to recreate audios.

4.3.2 Spectrograms of signalframes as Input

In the previous initial experiment, the samples from the NSynth dataset have been taken as a whole for the experiment. It has been mentioned, that all samples are 4 seconds long, but some contain padding in order to come to the 4 seconds. As those zero-paddings are then also part of the trained, this could affect the behaviour but also the outcome of the model. If those zero-paddings would be left out, this would lead to unequal long samples, which brings the problem with it, to not be able to be used for the model. This is because the network has a fixed size of neurons at the input, which implies to have a fixed input shape. Furthermore this also means, that having a fixed size of 4 seconds, the input always has to be of 4 seconds, which is not desirable. Not at least, if the system should be used in real-time applications for audio synthesis, one cannot wait to have 4 seconds of a signal, to perform audio synthesis with it. Therefore it would be desirable to perform audio synthesis on smaller “frames or chunks” of an audio signal, respective spectrogram.

Pre-processing

This leads, to the idea to take chunks or frames of audio data that get transformed into distinct spectrograms of same size. For a start the length of those frames, was set to 500 ms. Furthermore it was chosen, that those frames are not consecutive, but have an overlap of 50%. Additionally those frames were multiplied with a window function, like it is used in the STFT. Similar to the window function used in the STFT, here a Hann-window is used. Again as those trimmed signals were all differently long, they had to be padded to a multiple of the frame size respective hop-length. This ensures to have equally long chunks of the signal. As of the windowed frames, the first and last

frame didn't have overlapping parts at the beginning respective end. To overcome this issue, additional zeros were padded to form one frame on each side, to get there also an overlap. In combination, having the framed and windowed signal chunks, by overlapping each again with 50% and adding the values, this would yield the original signal again. This would especially helpful when reconstructing the final signal in the end.

Having those framed and windowed signal chunks, the STFT gets applied on those, having the same configuration as in the first experiment. This then lead to have multiple spectrograms for the length of 500 ms with again a frequency resolution of 31,25 Hz. Again for the training and testing, additional to the log-mag data, the phase information, reference value and name of the sample including a number to identify the frame got preserved for later use.

Model and training

The configuration of the model is rather similar to the one used in the first setting. It has the same amount of layers on each side, despite different strides, but also different kernel-sizes got applied (see figure 4.1). Again this one has been trained on 80% of the *keyboard_synthetic* test dataset samples. Of course the significant difference here is, that now the input data are not whole spectrograms but overlapping frames. This also means, that the amount of data has increased. When collecting the spectrograms for the dataset, the names get shuffled, in order to not have the same order. In contrast, the single frames, don't get shuffled as well not during training. Again the training has been performed over 20 epochs on the local machine.

Testing/Evaluation

For the purpose of evaluating the test score this has been done with the remaining 20% of the samples. To evaluate the ability of the autoencoder to reconstruct spectrograms, the whole samples were used including those from the training. Here the spectrograms were provided in the order as they appear. Having all reconstructed spectrograms, the inverse STFT were applied with the preserved phase information. Resulting in the frames corresponding to every single input note, those got overlapped and added (in their right order). By this procedure the signal with the original length could be obtained and further on evaluated auditorily. For results see chapter 5.

A step for interpolating two different sources has not been investigated here, up to these experiments. Also the embedded space also has not been evaluated so far but will be subject of further experiments.

4.4 Experiments single frequency vectors

The above mentioned experiments, were a prove regarding the ability of convolutional autoencoders to recreate audio spectrograms. From this point on the experiments were done using the whole training dataset and also include audio synthesis. In the previous experiment the model was trained on frames of audio data which were 500ms long. Having the idea to synthesise audio on real-time input, this would mean that always a signal frame of 500ms has to be present, in order to have an input for the model.

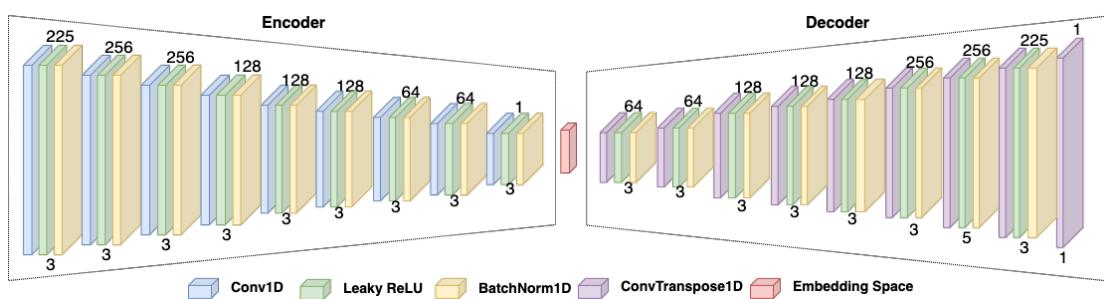
Therefore it is desirable to have an input that is as short as possible. An important note at this point is therefore, that spectrograms consist of frequency by time data. Each vector of the spectrogram along the time axis represents a short frame of time.

As mentioned before, when pre-processing was done on 16 kHz sampled data with an `n_fft` of 512, respective hop-length of 256, a time resolution of 16 ms was therefore present. Therefore the idea was to take those single frequency vectors, as input for the model. The shape therefore is frequency by channels which in the case of the previous parameters is 256x1. Having this idea, also the silence at the end of the signals could be omitted as just the frequency domain was used for the models input. This therefore enabled to take samples of different length in time. Throughout the experiment the value for the `n_fft` was increased to 1024, as an increase of the models performance could be achieved (more on that in chapter 5 and 6). Increasing this parameter therefore meant to decrease the time resolution, but increased the amount of frequency bins and thus having a better frequency resolution. Resulting in a number of 513 frequency bins (15,625 Hz/bin) and time resolution of 32ms per vector.

Neural network

As a consequence a different model had to be used, as 2D convolutions were no more suited. Therefore a model was designed that uses 1D convolutions. 1D convolutions are performing the same calculations, with the difference of having just a 1-dimensional kernel. This one dimensional kernel just operates on the frequency axis and tries to extract important features. Equally to the 2D convolutions they also apply the principle of channels. Regarding the channels, those also get expanded but then subsequently until the innermost layer reduced to 1 in order to form a single dimensional vector. Until the end again those channels get expanded but reduced again to 1 at the end to have the same shape as the input vector. In the following graphic (figure 4.2), the structure and configuration of this autoencoder is depicted.

Figure 4.2: Deep 1D convolutional autoencoder.



Training

The training process for this experiment, differs in some major points from the previous ones. The main difference lies in the shape of the data. This can be concluded as here just the single frequency vectors, were used for this network and therefore the network itself has 1D convolutions as discussed before. Furthermore, this experiment

was originally trained locally with the small subset like before, but later on, access to a GPU-accelerated instance was made available. This GPU-instance, as discussed at the beginning of this chapter, enabled to train complex ML-models with a large amount of data, over a long time. With this possibility, the training dataset could be utilized as this consists of several thousands of audio samples and would be too big to use locally. This does not mean, that the whole dataset was used all the time for training on this instance. With the progress of the project, several trainings with different amounts of data were made, that impacted the training process.

As here single frequency vectors were used, the process of creating a tensor and subsequently the dataset for the dataloader, was different. Here all spectrograms of the pre-defined and pre-processed dataset, were utilized, and the single frequency vectors get concatenated to form one big 3D array in the shape of amount by channels by frequency. To note the amount of channels here is also initially one. While experimenting, regarding the training, also different strategies regarding shuffling were used. First on just the samples were shuffled resulting in the instruments being mixed, but the frequency vectors were in the same order as in the spectrogram. This stayed the same as after each epoch the data didn't get shuffled. One more successful strategy, was to also shuffle the whole generated array, resulting in the frequency vectors of all instruments being mixed. Furthermore shuffling was also done after each epoch. The same strategies were also applied to the validation dataset, as from this point on also the validation dataset was considered. Not at least as more computational power was present from this point on. The most time, a batch-size of 32 was considered.

Introducing the validation step from this experiment on, also the scheduler to adapt the learning rate during the training was applied. Taking into account, that the network is rather complex and a huge amount of data has been used, the learning rate also had to be chosen adequately. Throughout this experiments, a small learning rate of $1e - 7$ was chosen, as this led to a more stable training and better convergence.

Testing

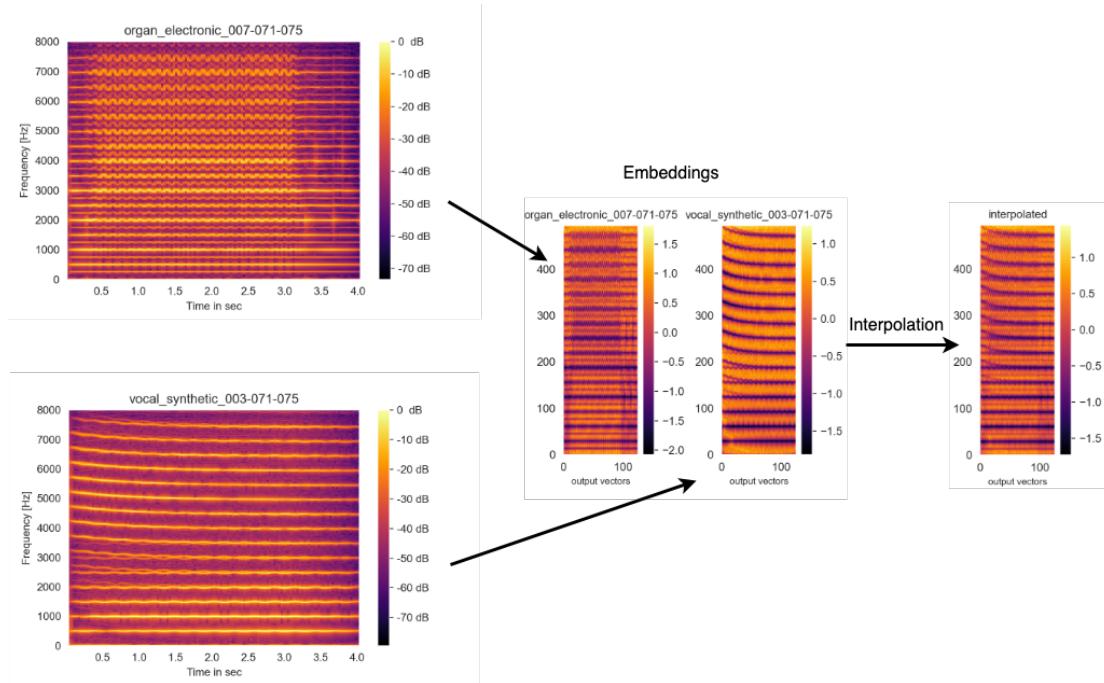
As described above, after each training, the current model state was saved, those can be used for the testing and evaluation regarding synthesis. As having numerous states that got saved, the one with the lowest error score on the validation set got considered for further steps. For the steps around testing and evaluation, the held-out test dataset was used. This stage was configured, to either take the whole test-dataset, a specific pitch, or a specific instrument source. By this technique its made possible to get the scores and therefore performance for certain instruments or pitch.

4.4.1 Experiments for Synthesis

Having this kind of network that was trained on the training dataset, to reconstruct single frequency vectors, some experiments have been conducted for generating audio. As during the training the network learned how to reconstruct frequency vectors, one experiment is to examine the quality of the output for a single non-modified audio. Those are similar to the ones conducted, in the first experiments above. Also by reconstructing just single notes, the preserved phase information could be reused, which enabled to apply the ISTFT. For a comparison also the Griffin-Lim algorithm was utilized.

As the main objective of this work was to examine the capability of creating novel sounds, from this point on in the project the interpolation step was introduced. With the interpolation step the encoded features in the embedded space of two instruments, were taken and value-wise interpolated. As it is known at this point, the encoder part of this network takes as input a vector of 513x1 and creates a lower-dimensional representation of it with the size of 495x1. Those representations, can be seen as the essential features, and were considered for the interpolation task. For this task, the frequency vectors of two instruments of probably the same pitch were passed through the network. Here the data did not get shuffled, as it was important to produce the values in the same order as they come from the spectrograms. The output of the encoder for each instrument, was concatenated to a 2D array 495xN where N is the number of encoded frequency vectors. To get a better idea of this concept, the next graphic (figure 4.3) shows two spectrograms with the corresponding accumulated output of the encoder.

Figure 4.3: Input spectrograms with embeddings and interpolated embedding.



Having those two representations, interpolation was performed. Like discussed in chapter 3, along the x-axis each output vector of one sample was interpolated with the corresponding output vector of the second sample. This process is also shown in figure 4.3 where the result of the interpolation process is shown. This process got applied equally for each subsequent experiments as the encoder output is always of the same shape, but more on that later on. This interpolated vector was subsequently passed through the decoder part which forms again frequency vectors of 513x1 that were accumulated to a spectrogram. This spectrogram then depicts the result of audio synthesis combining features of two different instruments.

Further on this spectrogram was converted back to audio domain. As in this case no phase information is present, the Griffin-Lim algorithm [8] for phase estimation was applied. By listening to the final sound, it should be possible to hear the characteristics of both instruments combined in this sound.

4.5 Experiments with slices of spectrograms

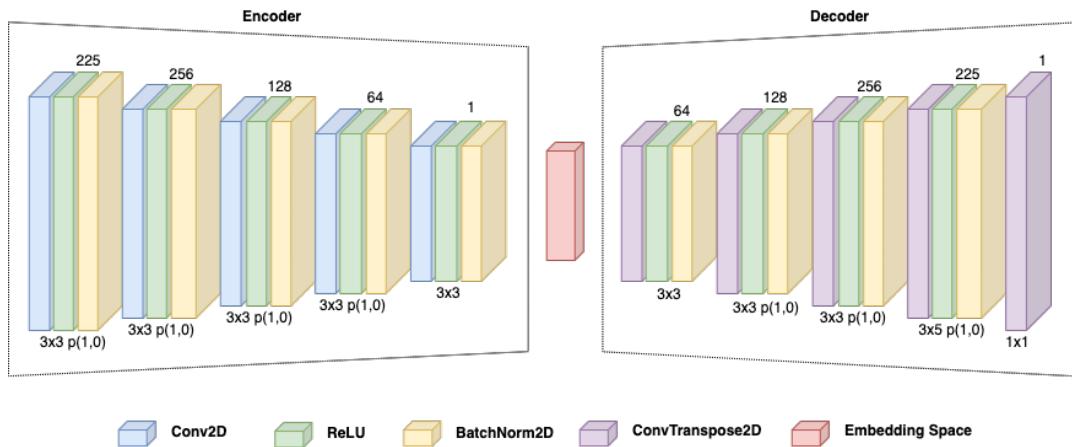
Based on the results and insights that could be gained with the previous experiments (see chapter 5) some more experiments were made. Not at least, to examine how different representations of the input data and as a consequence different model configurations influence the task of neural audio synthesis. For this case the following experiments in contrast worked with 2D representations as input. 2D data has already been used in the initial experiments with the whole spectrograms, but also spectrograms based on frames of the audio signal. As a difference for the following experiments, the spectrograms have been taken, but sliced into overlapping chunks. The spectrograms taken in this experiment, were generated with the same configurations, as in the previous experiment (`n_fft=1024, hop_length=512`). The idea to use chunks of spectrograms as input, arose to improve the models performance with regard to synthesis and recreation of spectrograms. As by using single frequency vectors, just the frequency information is present and no information about its change over time. As a theory to also incorporate the temporal axis as input, important information about the temporal frequency changes could be captured. Those frequency changes could deliver more important characteristics of the samples that could be extracted. As not the whole spectrograms should be used as input, chunks of them are favorable, to keep the input as small as possible. Therefore the choice has been made to take three consecutive frequency vectors to form a frame with the shape of (513x3). Throughout the experiments it has been tried out to either take the frames subsequently or overlapping, leading to the preference in taking the latter. The overlap has chosen to be always 2 vectors from the preceding frame.

For demonstration, let there be a spectrogram $spec = [f_0, f_1, f_2, f_3, f_4, f_5, \dots, f_n]$ of length n where f_i is the frequency vector at index i . This results in an array of $specs_slices = [[f_1, f_2, f_3], [f_2, f_3, f_4], [f_3, f_4, f_5], \dots, [f_{n-3}, f_{n-2}, f_{n-1}], [f_{n-2}, f_{n-1}, f_n]]$. The idea behind this overlap was to capture every change or time-pattern in the source spectrogram. This technique also had the advantage that all the silence of the source spectrograms could be cut away, allowing to use differently long spectrograms as input. A positive side-effect here was also, to gain more amount of input data, on which the model could be trained on. As the models output are also overlapping frames, to reconstruct the final spectrogram and audio, this also had to be considered, but more on that later on. Again when creating the tensor and dataset, the spectrograms were taken in a random order. At first the resulting frames of the training set, didn't get shuffled. Later on during the training, the frames in the batches, which again are of size 32, were shuffled like in previous examples. The same happened to the validation dataset, while it also was shuffled initially and after each epoch. An important note is that for the training and validation the whole dataset were used all the time.

Neural Network

As again input data in a 2D shape was used for these experiments, the model again had to be adapted. Instead of 1D convolutions and 1D batch normalization again, 2D convolutions and 2D batch normalization has been used. Also in contrast to the previous network, a normal ReLU activation has been used. These experiments incorporate different model configurations whereas a special focus has been given to the amount of striding and thus input compression. Here it should be evaluated, how more compression and thus a smaller latent space, influences the quality of the decoder output and further on of the generated audio. Both in terms of single note reconstruction and interpolation based synthesis. The following figure (4.4) shows, the basic network structure for this kind of experiments.

Figure 4.4: 2D convolutional autoencoder.



This network is not deep as the one used in the previous experiment because it consists of 5 layers on each side (10 in total). As the input is of size 513x3, the convolutions had to use a padding along the time-axis. If no padding would be applied, just one time a kernel of 3x3 could be applied. Thus 5 times a 3x3 kernel could be applied with the last layer having no padding. This results in a single 1D vector in the embedded space. In the decoder part, padding also had to be applied to regain the same dimensionality in the end. As previously said, these experiments should give an insight on how the striding and thus the size of the embedding influences the performance and quality of the output. Choosing more strides and a smaller embedding, the network has to learn to extract more efficient encodings, from which it can reconstruct the input data. Thus by choosing a small embedding size, the network should learn to extract the most important features of the input. Therefore the size is also of significance regarding the audio synthesis task by interpolating the embeddings. In the following table (4.1), the different configurations regarding the striding are shown.

To explain, each row represents a network configuration, that either uses one, two or three times a striding of (1,2) on each side of the network. As of the shape of the input, striding just can be applied on the frequency axis. The columns with the names “Encoder” and “Decoder” show the respective layers, where the stride was applied. For

	Encoder	Embedding-size	Decoder
Single Stride	e2	250	d4
Double Stride	e2, e4	124	d2, d4
Triple Stride	e2, e3, e5	62	d1, d3, d4

Table 4.1: Setting of stridings in network

example e2 is the second layer in the Encoder and d4 the fourth layer in the decoder. The center column called “Embedding-size” explains the size of the embedded space vector. By having these model configurations, interesting findings and results could be obtained.

4.5.1 Audio synthesis

In the previous experiment with the 1D convolutional network, the interpolation step has been introduced. This step was also applied with this network. As the embedded space vector also has the shape of a 1D vector, the interpolation procedure here was exactly the same.

Analysis of the embedded space

As with neural audio synthesis novel interesting sounds should get generated, it is also of interest to find interesting combinations for the interpolation. Combined with the fact, that the embeddings contain the extracted features/characteristics of a sound, those can be utilized for this task. It therefore has been implemented to take the output of the encoder of several notes (e.g. from the same pitch) and compute the correlation coefficients between those. The result got depicted in a correlation matrix to see which samples have the lowest correlation coefficients. Low correlation coefficients between two data samples mean that those have little similarities and thus different characteristics. By taking those embeddings with the lowest correlation coefficients and interpolate them, interesting novel sounds can be generated. More on that later in chapter 5.

4.5.2 Reconstruction and post-processing

As the shape of the output of the network is also the same as of the input of the network, this had to be considered when recreating the spectrogram. The output therefore were also again frames of 3 that had to be overlapped. During the experiments different strategies have been tried out, whereas it was preferred to average the overlapping parts of the output frames, in order to form the final spectrogram.

When having the final spectrograms, some further experiments regarding the improvement of the sound quality have been carried out. In this case this includes to correct energy in the frequency bands of the output. With this technique important properties of the sound like the transient or impulse (e.g. guitar stroke) should get preserved. For this task in advance the energy-values of the frequencies were summed up for each frequency vector of the original audio spectrograms. This was also done for the output spectrogram and its frequency vectors, after converting from db to energy.

By comparing those sums with the corresponding values of the input spectrogram, a factor could be calculated. This factor then was multiplied with the corresponding output frequency vector in order to have the same amount of energy present in the output spectrogram. If the output spectrogram was generated by an interpolated embedding, the energy sums of the two input samples were taken and averaged. This averaged values then were taken as a reference to correct the energy. By performing the inverse STFT or Griffin-Lim algorithm, the resulting audio sample should therefore also have a corrected amplitude and thus improved sound. Regarding the results, those will be discussed as well later on in the thesis.

4.6 Experiments with mel-scale

Up to this point, the experiments have all been conducted on log-magnitude spectrograms. Those spectrograms depict the frequency energies on a linear scale with a certain resolution (e.g. [0 Hz, 15.625 Hz, 31.25 Hz ... 8000 Hz]). Those spectrograms have been computed of a dataset consisting of musical notes. Those notes can be grouped into musical intervals describing the distance of a note to another, whereas the interval of an octave would be a doubling in frequency. As an example the note a' has a frequency of 440Hz (if perfect sine), whereas its octave on top would be an a'' with 880 Hz. If going an octave down it would result in 220Hz having an a (to be continued A 110 Hz, A' 55 Hz, ...). This also means that the higher the note gets, the larger the distance in Hz and the lower it gets, the less distance in Hz between the notes. Having this principle one can come to the conclusion that having a linear scale like in log-magnitude spectrograms, the resolution with lower notes is worse then with higher notes. With a look onto machine learning and neural networks, this also means, that there is less data available to train on for low notes. Keeping this in mind, the choice has been made to perform comparative experiments based on a different scale, which is called mel-scale [26]. This scale can be seen as a compressed form of a spectrogram. This scale is an empirical scale that is based on the human auditory perception. This means that humans percept low frequencies louder and with a greater resolution then with higher notes. Mapped on the mel-scale, in contrast to magnitude, lower notes have a greater distance between them. The opposite applies to the higher notes as there the distance gets less. This means that with this scale lower notes can be better differentiated and thus are more emphasized then with the linear scale. Even though that this scale is based on observations it was proven beneficial regarding machine learning-tasks. Not at least as several approaches, mentioned in related work, applied this scale for their audio synthesis task. Furthermore this scale could also be seen as a compressed form of a spectrogram containing the most significant properties.

4.6.1 Pre-processing

The pre-processing step does not really differ from the ones performed in previous experiments. A difference here is in order to obtain spectrograms with the mel-scale to call the *librosa* function `librosa.feature.melspectrogram(...)` which takes as input a pre-computed STFT-spectrogram. Throughout this experiment the spectrograms created with an `n_fft` of 1024 and `hop_length` of 512 have been taken to be converted to

mel-spectrograms. The resulting mel-spectrogram has a size of $128 \times t$ and thus is smaller than the mag-spectrogram ($513 \times t$). Having the mel-spectrograms the same steps as with log-magnitude were performed (db conversion, preserving the power reference,...). Also for the model-training, frames of 3 vectors of the mel-spectrogram were taken as input for the model. In the next figure (4.5) a mel-spectrogram with its corresponding log-magnitude spectrogram is displayed, to see the difference.

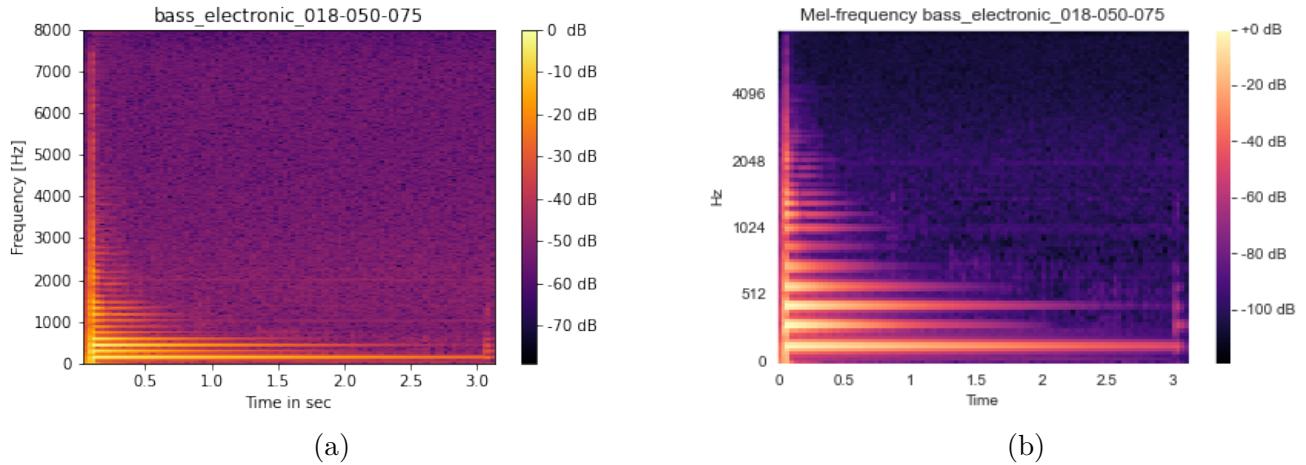


Figure 4.5: log-mag spectrogram (a), log-mel spectrogram (b)

By comparing those, it can be seen, that the scale is logarithmic and also that the distance between the low frequencies is greater than between the high-frequencies. As of this known characteristics and its already compressed form, the experiments using these kind of spectrograms thus can be expected to bring different but interesting results. Using the mel-spectrograms also the same experiments with and without interpolation got conducted to get a comparison.

Neural Network

Like in the previous experiment, also three different networks with single, double or triple strides on each side were applied and assessed further on. The basic network structure and configuration does not differ of the one used in the previous experimental setting as having the same amount of layers and same configuration of the channels. Nevertheless some differences arise, because the last layer applies a kernel of 1×2 to regain the same input shape. Also in the networks using three or two times strides on each side, padding in the layer $e4$ and $d4$ is set to 1,1. Regarding the network with three strides on each side, those have been set onto layers $e2$, $e3$, $e4$ and $d2$, $d3$, $d4$. As the input of the network is already of smaller size, the corresponding embeddings are therefore significantly smaller like in the previous experiments (56, 28, 14). Those were also altered like in the previous experiments by interpolation to synthesize novel sounds.

Reconstruction and post-processing

The steps around the reconstruction and post-processing are similar to the previous experiments. Because as input also frames of 3 were used, the output consists of 3 frames too. Those also were overlapping and averaged to form a final log-mel spectrogram. Again also the energy got corrected to have the same amount of energy present like in the input. Having the final log-mel spectrograms, those could be either directly converted to an audio using the function `librosa.feature.inverse.mel_to_audio` or first to an magnitude spectrogram. The latter is used when just single audio samples were reconstructed (without interpolation) as here the phase information could be reused, and thus the ISTFT can be applied to obtain the final sound. Otherwise the mentioned function is performing the conversion to magnitude and Griffin-Lim at once to obtain an estimated signal.

Chapter 5

Results

This chapter, shows the results that could be obtained through the experiments, mentioned in chapter 4. Those results contain either numbers, for assessing the model performance with a certain error performance, but also of graphics, displaying the outputs of the models. The outputs of the models, being spectrograms, will get compared with their corresponding input spectrograms. As the output of the encoder part (embedding) play a crucial role regarding audio synthesis, those also get displayed and assessed further on. Some additional graphics containing additional spectrograms, embeddings and signal plots will get displayed in the appendix A. Regarding the final listenable sounds, those get discussed in the next chapter.

5.1 Results regarding models with single Frequency Vectors

As described in chapter 4 experiments have been conducted by using the single frequency vectors of the spectrograms. As a model, a 18 layer deep neural network with 1D convolutions has been trained on the reconstruction of single frequency vectors. Throughout the development different settings have been tried out in order to find a model with an optimal training process and performance. The amount of data on which the model has been trained on, has been varied as different performances regarding the convergence could be observed. First trainings have been made on keyboard_synthetic mixed with guitar_acoustic. Here it could be observed, that no sufficient convergence could be reached ($MSE > 200$). Training a distinct model per instrument source, showed that the instrument source has a significant influence on the convergence. For example a model trained solely on keyboard_synthetic converged better than a model with guitar_acoustic. Comparing the performance of a network trained on guitar_acoustic with one trained on e.g. guitar_electronic, showed that with the latter a better score could be reached (~17). Therefore the decision was made, to combine keyboard_synthetic and guitar_electronic instead of acoustic, for the training. Compared to the one mixed with guitar_acoustic a score of ~78 could be reached. Subsequently more instrument sources were added, until the whole training dataset was utilized. This model showed the best performance regarding its error scores and was chosen to be considered for the ongoing experiments. The next table 5.1 shows the error scores of this model. As a short note, each trained model was chosen based on the best validation error score.

	MSE-Score
Training	4,237
Validation	4,638
Test	1,190

Table 5.1: MSE-Scores 1D convolutional autoencoder

Here it can be seen that the error scores for training and validation are rather close, while the score on the held out test set was significantly smaller. This score could be reached with a training over 8 epochs and was considered as the best score, as after those 8 epochs the validation score increased again (despite decreasing training score). The next table 5.2 is also interesting as it shows, the error scores on the test set, regarding different pitches. As displaying all pitches would go beyond the scope, pitch classes from 030 to 100 have been chosen. In this case it can be seen, that all error scores have roughly the same value and do not differ significantly. Furthermore there is no visible trend if the score increases depending the pitch.

Pitch	MSE-Score
030	1,049
035	1,393
040	0,935
045	1,096
050	1,019
055	1,108
060	1,086
065	1,111
070	1,164
075	1,041
080	1,177
085	1,753
090	1,100
095	1,796
100	2,043

Table 5.2: MSE-Scores for specific pitch classes using 1D convolutional autoencoder.

5.1.1 Experiments of single reconstruction

Having the trained network, this one got evaluated towards the ability of reconstructing audio spectrograms and further on recreating the sounds. The next graphics (5.1, 5.2 and 5.3) show the result of taking a guitar sample as input and reconstructing it.

With a special look onto the reconstructability of spectrograms, figure 5.1 shows the input spectrogram and the generated output spectrograms from the network. Here it can be seen, that the reconstructed spectrogram, differs in a few points from the input

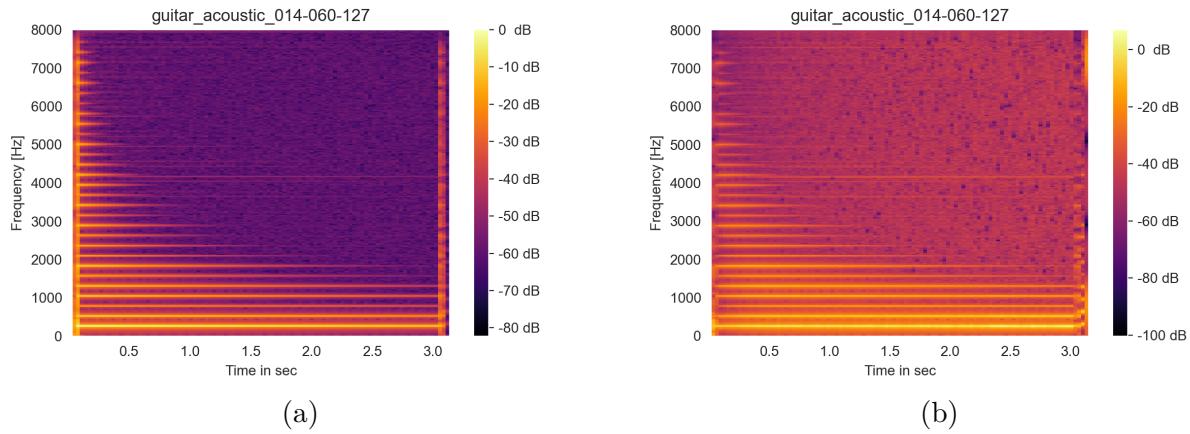


Figure 5.1: input spectrogram (a), reconstructed spectrogram. (b).

spectrogram. First of all it can be seen, that across the frequency areas that contain little energy in the input spectrogram, in the output more energy is present. This also means that between those areas and the sound-characteristic high energy areas, less difference is present. Furthermore regarding the original broad spectra at the beginning and at the end, those are hardly present in the output spectrogram. Knowing that these represent the stroke and the damp of the string, it can be said, that these are not present in the output. This gets also confirmed through looking at the time-domain signal.

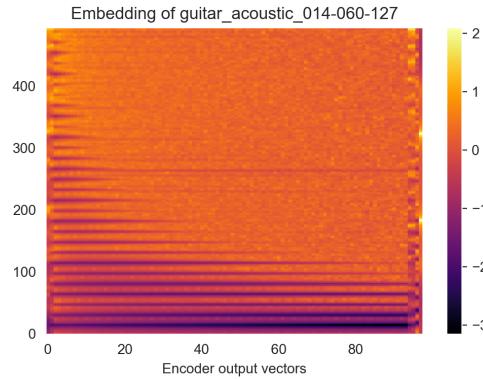


Figure 5.2: encoding of guitar acoustic.

Looking at the graphic that depicts the embedded space, there it already can be seen that those broad spectra are not preserved. Despite of that the high energy areas (harmonics) in the input spectrogram are preserved as negative values in the embedding while the original low energy areas get represented by positive values. Also as there are no strides applied in the network, there is almost no compression in the embedding.

A final look onto the time domain plots (figure 5.3 also reveal, that there's no impulse at the beginning of the signal. Furthermore it also can be said, that the amplitude in general differs in its course but also in how strong it is.

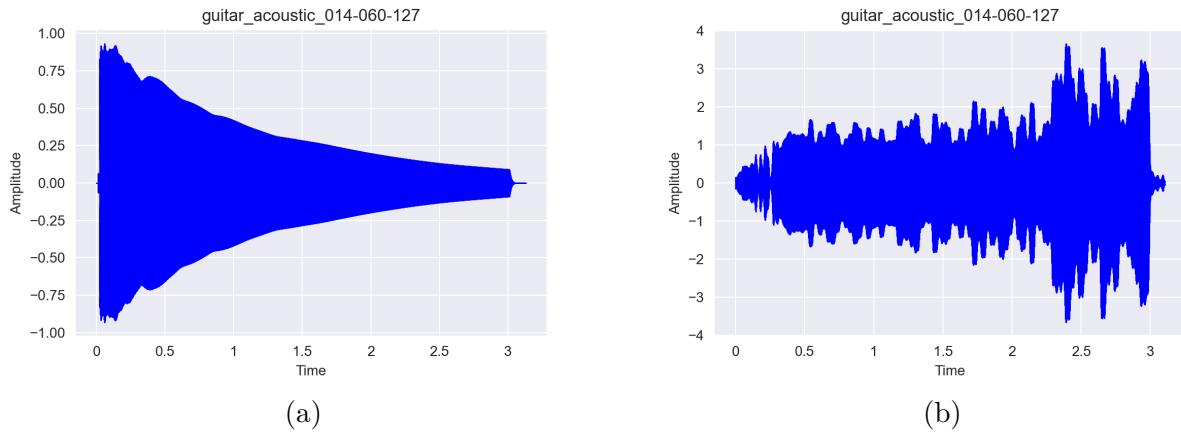


Figure 5.3: input signal (a), output signal (b).

5.1.2 Experiments with interpolation in embedding

With this kind of 1D convolutional autoencoder first experiments were made, that use interpolation in embedded space to synthesize novel sounds. In chapter 3 and 4 the concept has already been described in detail. Having this concept, instrument sources were chosen, which in order got encoded and interpolated together to in order form a new embedding. The next graphic 5.4 shows the spectrograms of the two input instruments that were chosen to create a novel sound. For this representation an acoustic guitar and acoustic brass sample were taken and encoded. By comparing those two spectrograms, it can be seen, that they differ significantly in their structure. While it can be seen, that the guitar looses energy in its higher frequency ranges, the brass sample keeps its frequency energy over the whole time. Furthermore as the harmonics stay constant in the guitar sample, those alternate in the brass. Those properties make it interesting to use them for the interpolation and thus synthesise a novel sound. For comparative reasons, in further experiments also those two instruments have been used for the synthesis task which can be seen later on.

After encoding, the outputs were taken and interpolated as can be seen in the next figure 5.5a. As it can be seen, those embeddings are a compressed form of the input spectrograms. When looking onto the generated interpolated embedding, it can be said, that this one incorporates both instruments encoded features. Having this new vector, this one was fed into the decoder network to in order generate an output spectrogram that can be seen in figure 5.5b.

Having the final spectrogram here as output, it can be seen, that spectral data of both input spectrograms are basically contained. Similar to the output spectrogram in figure 5.1b this one also does not contain the impulse (guitar stroke) at the beginning. To finally obtain a listenable sound, the output spectrogram has to get converted back to time-domain, which in this case gets done with the Griffin-Lim algorithm [8] as no phase information is present. To improve the quality of the output but also to examine the performance of other types of networks some further experiments have been made using 2D convolutions but also additional post-processing steps.

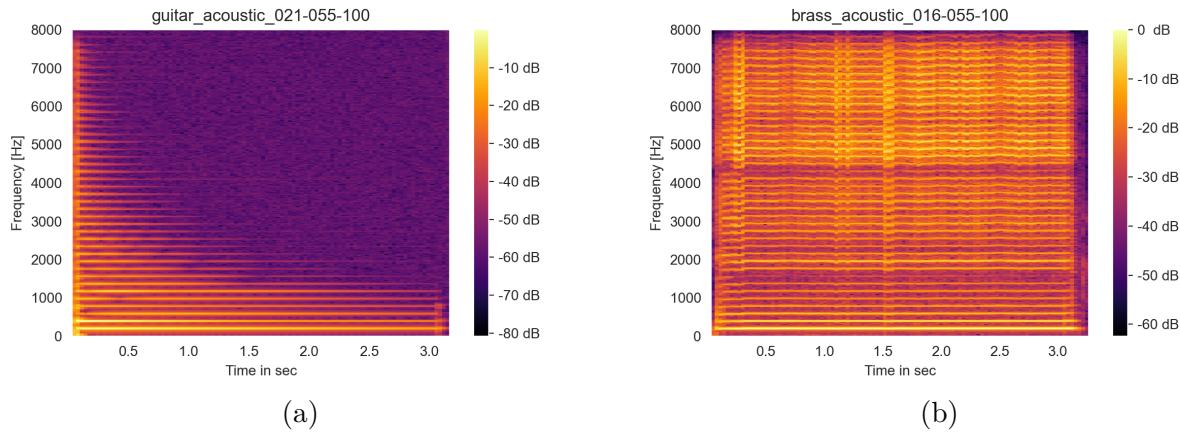


Figure 5.4: guitar acoustic (a), brass acoustic (b).

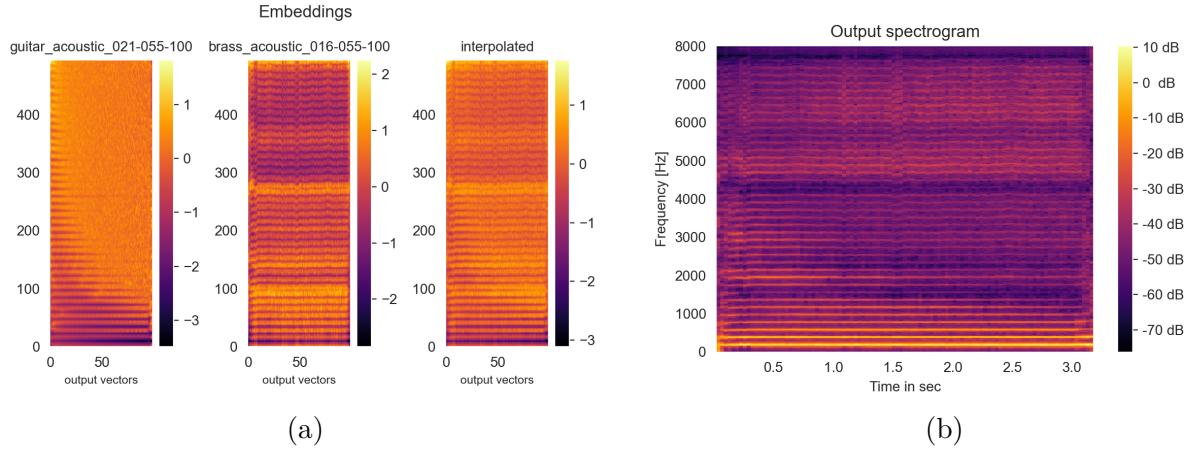


Figure 5.5: embedding interpolation (a), output signal (b).

5.2 Results of experiments with spectrogram frames

The here shown results, correspond to the described experiments in section 4.5. Here again a 2D convolutional autoencoder has been utilized to (re)synthesize audio. As input data for the neural networks, overlapping spectrograms snippets with the length of 3 along the time axis have been considered. In contrast to the previous network, this was supposed to bring better results, regarding the output but especially to preserve transients. As different networks with different stridings (compression) were trained and used, the results should also show the differences in the output but also in their quality regarding audio (re)synthesis. Additionally the encoded data gets examined by calculating the correlation coefficients between samples (mostly of the same pitch). By that it can be seen how much the encodings of different instruments differ. Additionally it also can be examined how good the different networks extract the essential features of the input data. Throughout the project these correlation coefficients got depicted in a

so-called correlation matrix. Displaying one here would not work out sufficiently, as because of its many values, those are hardly readable without zooming in. Because of this, just the findings get discussed later on in chapter Discussion. Finally also additional post-processing steps get applied, which help to correct the energy and thus the quality of the output. The results shown here are without energy correction, but the effect will also be discussed later on in the discussion.

The following table 5.3 shows the MSE error scores of the different 2D convolutional networks. As mentioned in chapter 4, three networks were trained, that differ in their configuration regarding striding and thus have different embedding sizes. These networks therefore are called single, double- and triple stride networks, regarding the amount of striding on each network side. By this they can be distinguished better in this work. In table 5.3 it can be seen that by using a different amount of strides, this has an impact onto the score. The network with just one stride on each side, has the lowest error whereas the networks with two or three strides have a significant higher error. Regarding the difference between training, validation and testing error all three networks show the same behaviour as the validation error is higher than the training error with the test score being the best.

	single-stride	double-stride	triple-stride
Training	9,779	13,717	18,292
Validation	10,094	14,056	19,056
Test	7,826	10,708	16,655

Table 5.3: MSE-Scores 2D convolutional autoencoder - single stride, double stride, triple stride.

Compared to the error score of the 1D convolutional network, it can be said, that all error scores are significantly higher. In case of this work, the score does not mean, that the quality of the output is worse, especially regarding the final output sound. The latter will get discussed later on in the discussion of this work

As during the previous experiment also the scores regarding different pitches were calculated, the next table 5.4 shows the different scores regarding pitches ranging from 30 to 100. Similar to the 1D convolutional network, the error scores regarding the pitch do not show a specific trend. Interestingly the scores of the triple-strided network has some rather strong outliers regarding samples around the classes 045-050. Having the same error scores this also does not mean, that all pitches have the same audible quality but more on that in chapter 6.

5.2.1 Experiments of single reconstruction

Having mentioned the error scores regarding reconstructing spectral audio data, those cannot be taken solely to assess the performance of the network. Therefore experiments, like with the previous network, were conducted in recreating single spectrograms. With those the ability towards reconstructing audio spectrograms becomes assessed visually and auditorily for the purpose of audio (re)synthesis. For comparative reasons, the same spectrogram source has been chosen like with the previous network. As the original input

Pitch	single-stride	double-stride	triple-stride
030	7,084	9,175	12,241
035	7,912	12,294	15,426
040	7,292	11,205	18,639
045	7,251	12,434	22,151
050	6,840	9,903	20,287
055	7,535	10,434	18,102
060	7,444	10,128	17,670
065	8,313	10,286	16,497
070	7,410	9,868	14,894
075	7,850	9,949	14,273
080	8,624	10,529	15,486
085	9,299	11,090	16,598
090	8,798	11,435	16,415
095	9,948	12,032	17,049
100	12,346	13,179	19,812

Table 5.4: MSE-Scores for specific pitch classes using 2D convolutional autoencoder.

spectrogram has already been shown in figure 5.1a, here just the output of the encoder part (embedding) but also the total output gets depicted (see figures 5.6, 5.7 and 5.8). The spectrograms in each of these graphics were generated with single-, double- and triple-stride networks and do not contain the energy correcting post-processing mechanism. Despite of this fact, it can be seen when looking onto all the output spectrograms (5.6a, 5.7a and 5.8a), that all preserve the broad spectra at the beginning and ending of the spectrograms, especially by looking onto the first one. When comparing again the low energy areas of the input spectrogram, to the output spectrograms, it can be said, that they also contain similar little energy. Again it can also be seen, that the embedding looks like a spectrogram but in a compressed form of the input, as it also contains similar structures. Contrary to the embedding in figure 5.2 the high energy areas have positive numbers while original low energy area have negative numbers.

With a look on the output of the double-stride network (figure 5.7a), it also can be said, that the broad spectra are preserved but with less energy. This can be also noticed when looking on the output spectrogram of the triple stride network (figure 5.8). The latter also shows less energy in the high energy areas and less “precise harmonics” (washed out).

Looking onto the embeddings of the two networks it can be said, that those contain significantly less values across the y-axis. Comparing it to the input but also the output, despite of the significant compression, the significant harmonic features and general structures are present as positive numbers (regarding the colorscale). As concerning the double strided network, the embedding still has a fine granularity contrary to the one obtained by the triple strided network.

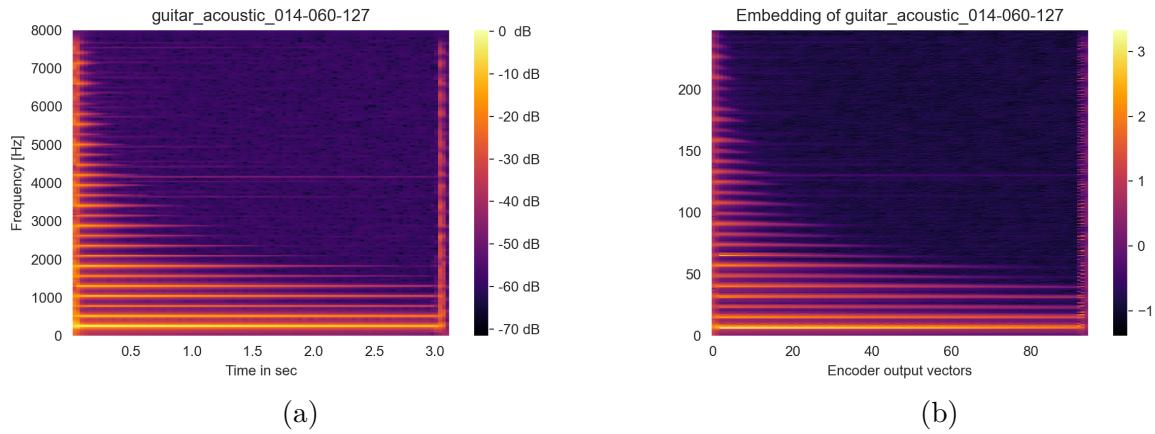


Figure 5.6: reconstruction of guitar acoustic (a), embedding of guitar acoustic (b) single stride model.

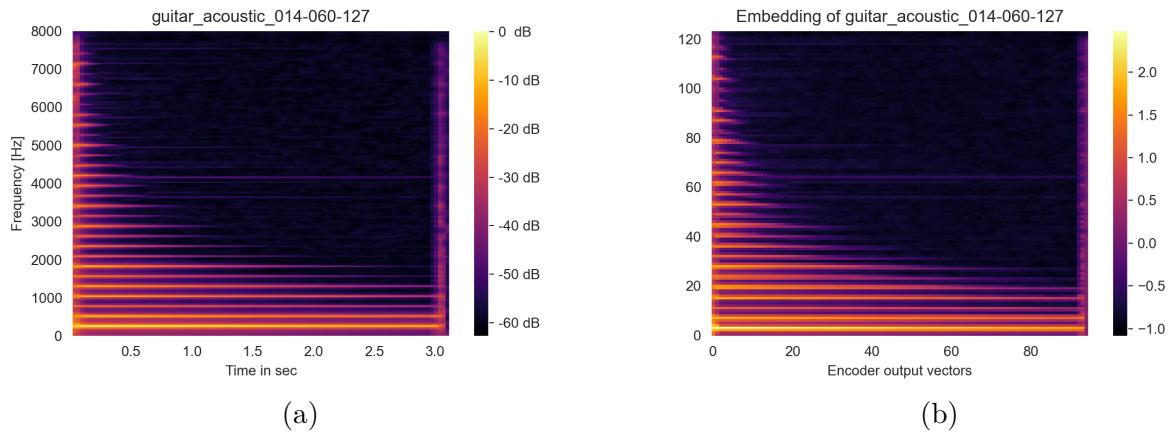


Figure 5.7: reconstruction of guitar acoustic (a), embedding of guitar acoustic (b) double stride model.

5.2.2 Experiments with interpolation in embedding

The results in this section show the resulting outputspectrograms, that got generated by interpolating of the embedding space vectors. This has already been done with the 1D convolutional network where novel sounds could be generated. As with the 2D convolutional networks used in this experiment, promising results in reconstructing single spectrograms could be obtained, this experiments yield interesting results. Not at least, as the embedding got more compressed, this also effects synthesizing new audio, as this is done by interpolating the embeddings. The following graphics (5.9, 5.10 and 5.11) show again the interpolation between the same two instrumental sources, for comparative reasons. Again also the resulting output spectrograms are displayed to see the final result. The auditory quality again gets assessed and discussed in the next chapter. When looking at the embeddings in the three graphics, that after interpolating, the features

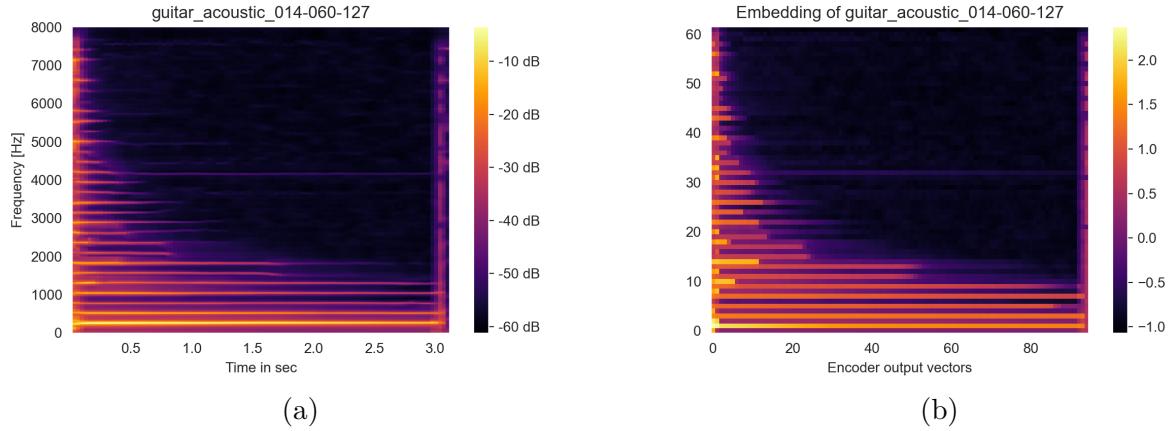


Figure 5.8: reconstruction of guitar acoustic (a), embedding of guitar acoustic (b) triple stride model.

of both instruments are visible in the result. Here it can be seen that the “harmonic” features of the brass sample do not fade contrary to the guitar samples. Nevertheless because of the interpolation the features are present but having lower values as it interpolates mostly between negative and positive numbers. The areas where the both samples have common features (lower harmonics), rather stay equally valued as those have close values in both embeddings.

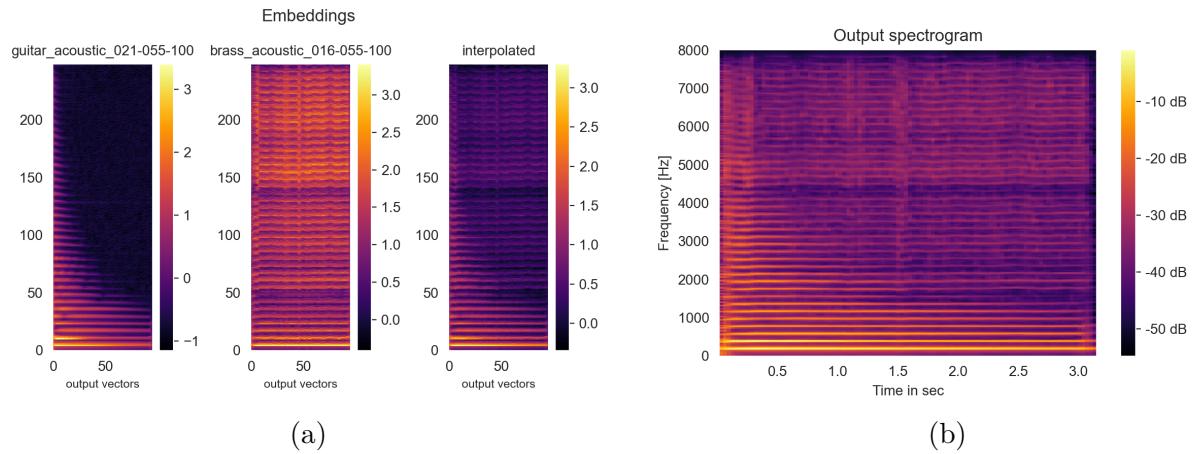


Figure 5.9: embedding interpolation (a), output signal (b).

By looking again on the decoded output it can be seen, that again both instruments are present. Compared to the result in the first experiments, the guitar sample is more present with these kind of network. Contrary to the experiment with single sample reconstruction, here the difference between the different strided networks, can be noticed even more. With special notice to the higher harmonics of the original brass samples. Those harmonics appear more precise with the single strided network (see figure 5.9b).

In the output of the double-stride network in figure 5.10b the “contours” coming from the brass sample, are less sharp than with the single stride.

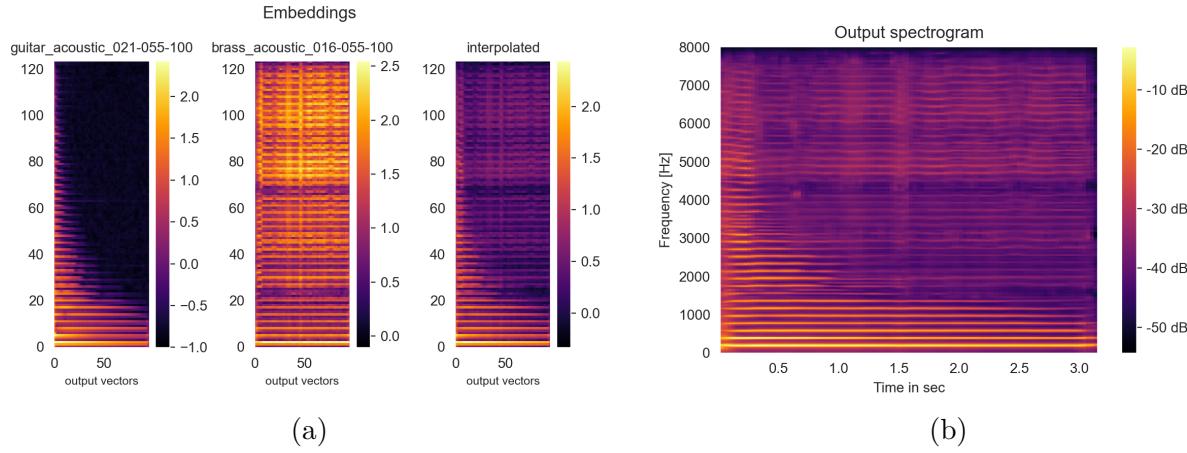


Figure 5.10: embedding interpolation (a), output signal (b).

Taking the output of the triple strided network (figure 5.11) into account, the difference to the other networks can be seen clearly. The encoded features, again contain the same structure as the input spectrograms, but due to the three-fold striding, just the most significant information is present. Comparing the interpolated embeddings it also can be said that in the latter, the guitar features can be recognized better. Looking at the final output spectrogram in 5.11b it can be seen, that the original fine harmonics aren't present anymore. The energy values of the different frequencies over time don't remain constant and show “washed out” contours. Also the fine changes over time are not as precisely present as with previous network configurations. Nevertheless both instrumental features can be recognized in the final output. Again in the discussion part, the auditive quality of all outputs presented here, with and without energy correction, gets assessed and brought into relation with the here displayed spectrograms.

5.3 Results with mel-spectrograms

The final experiments are again done with a 2D convolutional network, but using mel-scale spectrograms instead of log-magnitude. In chapter 4 under section 4.6 an introduction to the mel-scale has been given. As those log-mel spectrograms are a compressed version of the log-magnitude spectrograms and emphasizes the lower frequencies, these experiments are expected to bring different but interesting results. This is again meant regarding the model performance but also towards the quality of reconstructing spectrograms with or without modification of the embeddings. The results again get depicted in tables filled with the MSE scores of training, validation and testing but also the scores regarding specific pitch classes of the test set. Later on again the graphics contain the output spectrograms and embeddings with and without interpolation. Again three different networks with single, double and triple strides were used. Those are similar in their structure to the ones used with log-magnitude spectrograms, which was already

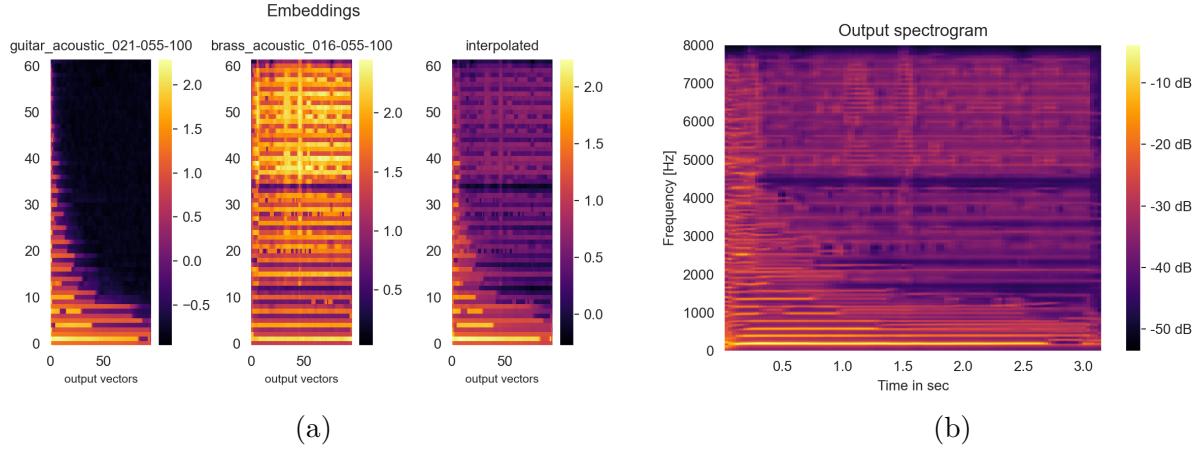


Figure 5.11: embedding interpolation (a), output signal (b)

explained in the previous chapter.

The following table 5.5 shows the training, validation and testing scores of all three different networks. First of all the MSE scores of all networks are significantly higher than the ones of the log-magnitude networks. An exception makes the testing score of the single-stride network as it is equally to the training score of the triple-strided log-magnitude network. The scores between the networks, also increase with the amount of compression, making the single-stride network again the best performing network regarding its score. Interestingly the validation score of the single-stride network, is the only one that is higher as its corresponding training score. The double- and triple-stride networks therefore show that they have a lower score on reconstructing unseen data. Again the auditory quality will get assessed in the next chapter 6.

	single-stride	double-stride	triple-stride
Training	30,579	49,685	62,833
Validation	31,084	48,439	60,471
Test	18,322	33,689	43,300

Table 5.5: MSE-Scores 2D convolutional autoencoder with log-mel spectrograms - single stride, double stride, triple stride.

Like in the previous experiments, the network performance has been evaluated towards certain pitch classes. The following table 5.6 shows the MSE scores of different pitch classes. For comparative reasons, again they got categorized on a scale from 30 to 100 in 5 steps. Contrary to those scores in the previous experiments, here the values show, that the lower pitched samples have the best score. With some exceptions/outliers, the error scores increase with a clear trend, the higher pitched the used samples were. Interestingly the difference in the score, between the highest and lowest used samples, increases, the more compression is used.

pitch	single-stride	double-stride	triple-stride
030	17,178	24,567	30,302
035	18,380	26,142	36,509
040	19,382	34,125	44,062
045	22,036	55,770	51,097
050	22,398	41,101	49,810
055	22,159	43,215	56,664
060	21,472	40,493	54,690
065	22,303	40,140	51,926
070	23,875	41,950	50,904
075	24,571	37,424	49,297
080	29,053	43,572	60,926
085	29,395	46,822	62,750
090	33,316	51,408	72,206
095	31,385	51,593	70,607
100	32,129	51,426	70,890

Table 5.6: MSE-Scores for specific pitch classes using 2D convolutional autoencoder using mel-scale.

5.3.1 Experiments of single reconstruction

Similar to the other experiments to assess the ability of reconstructing spectrograms and therefore recreate audio samples, in this section the embeddings and corresponding output spectrograms are depicted. There again the properties and respective differences to the other experiments get assessed. Regarding the auditory quality the results get discussed in chapter 6. The here shown results again correspond just to the reconstruction of single samples without interpolation. As in these last experiments the mel-scale as different scale was used, the corresponding input spectrograms are shown in advance. Figure 5.12 shows the input spectrogram of the single reconstruction experiments.

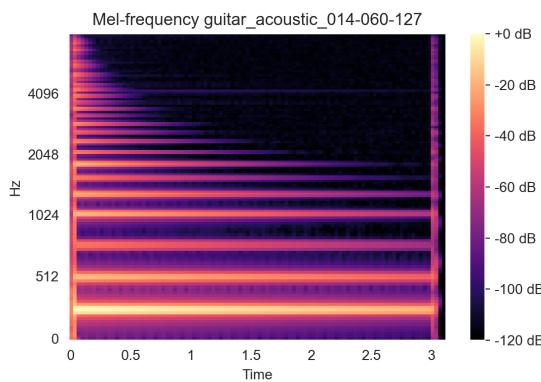


Figure 5.12: original mel-spectrogram of guitar acoustic.

The next three figures (5.13, 5.14 and 5.15) show the output spectrograms but also the embeddings of the guitar sample depicted in figure 5.12. As its known at this point, the mel-spectrograms are already a compressed version log-magnitude spectrograms, as they contain 128 values instead 513. Thus it can be seen on the y-axis of the embeddings, that those embeddings are even more compressed than in the last experiment (size of 56, 28 and 14). Furthermore it can be noticed regarding the embeddings of the first two networks (figures 5.13b and 5.14b), that contrary to the previous experiments with log-magnitude spectrograms, the low energy areas are positive valued. The high-energy areas and harmonics are therefore negative valued. Concerning the scale it can be also said that those embeddings have a narrower value range then those in the previous experiments.

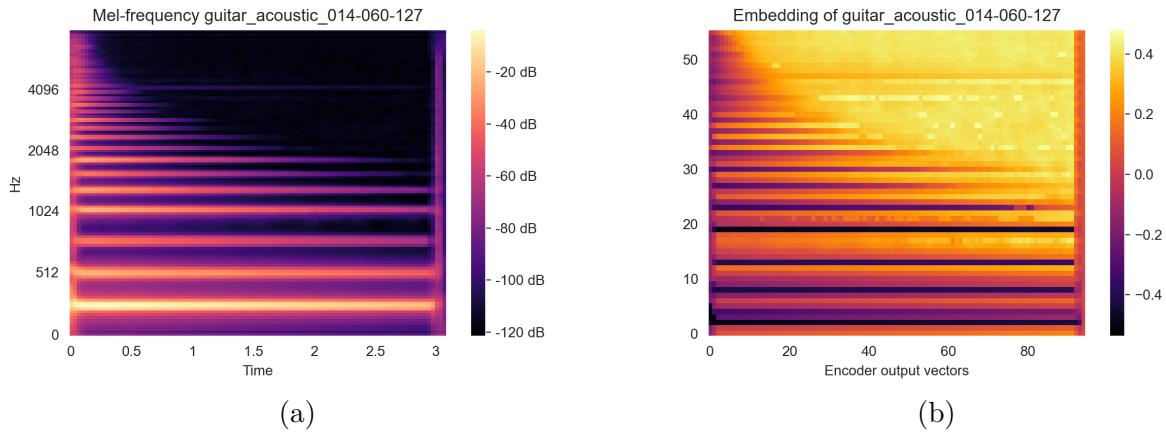


Figure 5.13: reconstruction of guitar acoustic (a), embedding of guitar acoustic (b)
single stride model.

In figure 5.13a the output of the single strided network can be seen. Comparing to the input spectrogram, its reconstruction is the best out of those three networks, as it contains the most details. Despite of this, it can be seen, that the broad band spectra at the beginning and end, are not as intense as with the input. Further on it can be seen, that changes over time in the frequencies between the lower harmonics, are not present in the output. The embeddings in figures 5.13b and 5.14b show that the features representing the lower harmonics, are recognizable as negative values. The latter already shows a lack of features in the higher harmonics. This is caused, because the input log-mel spectrograms have a shorter distance between the higher harmonics then in the lower ones. Regarding the output of the double-stride network it can therefore be said, that the higher harmonics can be less distinguished from each other then with the single-stride network (see area between >2048 Hz).

Comparing the result of the double-strided network with the one of the triple strided network, significant differences can be discovered. Especially when looking onto the embedding, it can be discovered, that due to the small amount of values on the y-axis the features are less distinguishable. Thus it is also more difficult to describe them. Nevertheless it can be said for this input spectrogram, that the colored areas still represent the most significant features that describe the input. Interestingly when looking onto

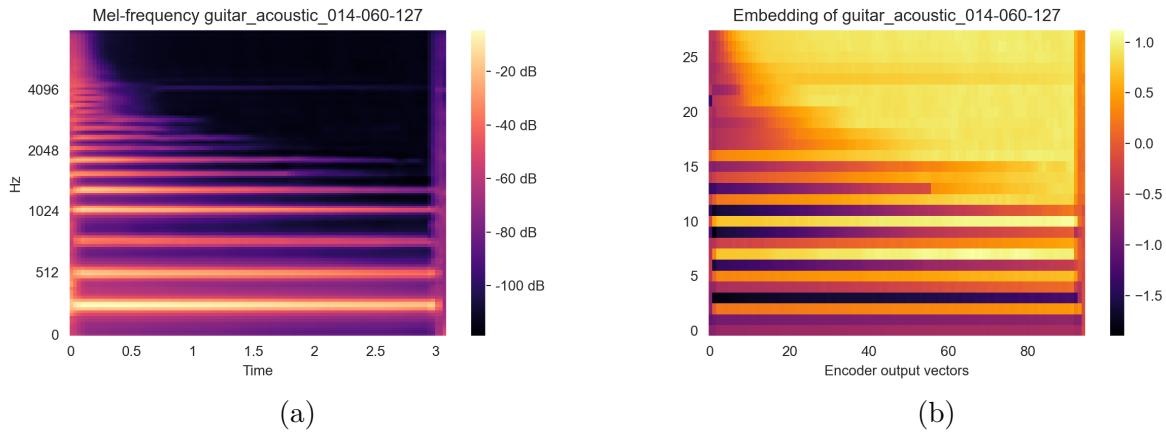


Figure 5.14: reconstruction of guitar acoustic (a), embedding of guitar acoustic (b) double stride model.

the output spectrogram in 5.15a, depending the lower harmonics, no differences can be seen to the previous output spectrograms. The part around 4096 Hz also shows some distinct spots that contain more energy, which are not present in the input spectrogram (artefacts).

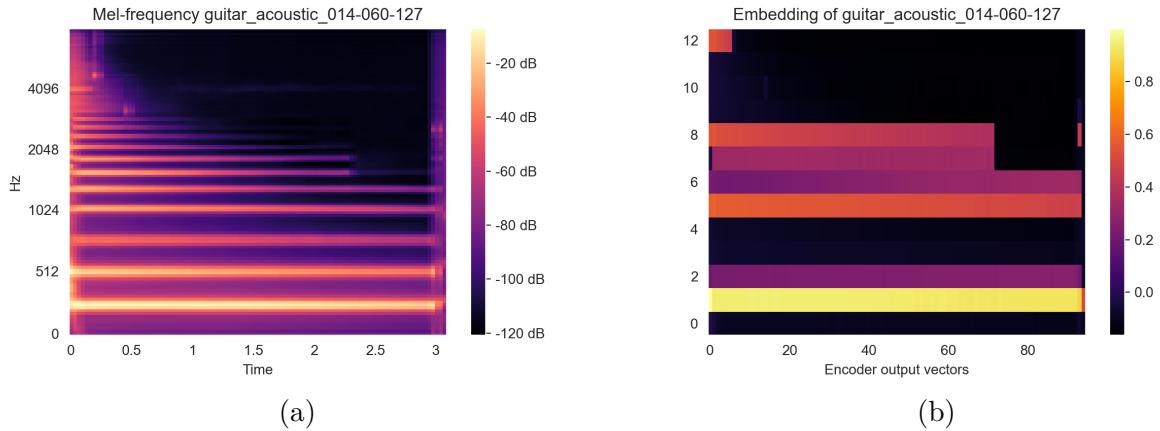


Figure 5.15: reconstruction of guitar acoustic (a), Embedding of guitar acoustic (b) triple stride model.

5.3.2 Experiments with interpolation in embedding

The final experiments in this work also deal with the introduction of the interpolation step to synthesize novel sounds. With the experiments of reconstructing single audios, it could be seen that despite of the small compressed embeddings, the input spectrograms could be reconstructed adequately. Having those smaller embeddings, less values are available to reconstruct spectrograms which in conclusion means that altering those

brings also significant changes. This fact makes these experiments even more interesting, as here again, the embeddings of two input samples, get interpolated, to generate a novel sounds. As an example, the same guitar and brass samples, that were used in the previous experiments, are utilized to synthesize audio. The following graphics 5.16 show the input log-mel spectrograms of the two source signals. To mention at this point again, the higher harmonics are closer in the scale, while the lower harmonics have a greater distance. Especially when looking at the brass sample this can be noticed as the energy does not fade.

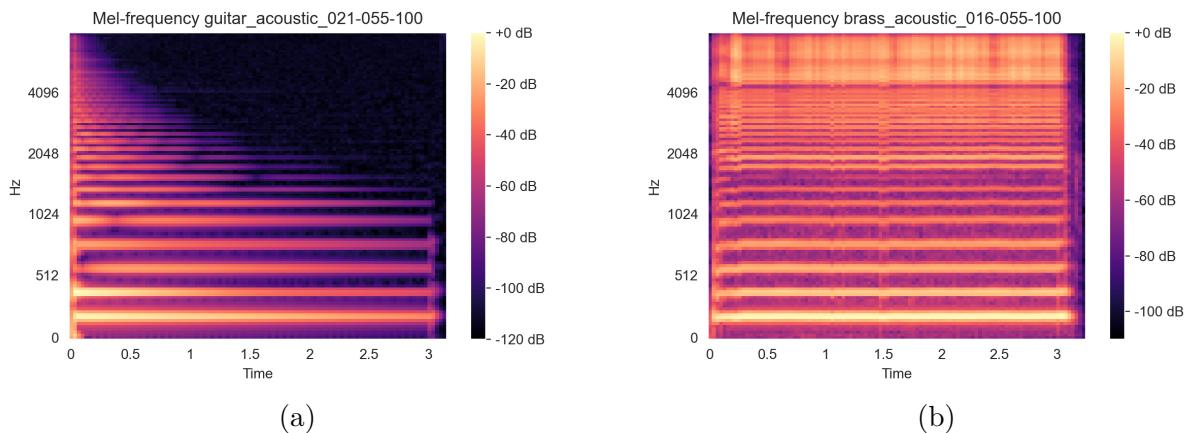


Figure 5.16: original guitar acoustic (a), original brass acoustic (b)

Knowing the source log-mel spectrograms, the following graphics show the embeddings of those spectrograms, the interpolated embedding as well as the resulting output spectrogram using again the three differently strided networks (single-, double- and triple-stride). By comparing all the output spectrograms to the two input spectrograms, it again can be said, that the output of the single stride network in figure 5.17b is the most precise and contains the most of both instruments. The upper harmonics of the brass sample can still be recognized in the output spectrogram despite the “fading” influence of the guitar sample. Fine changes in the harmonics over time of the trumpet, like at second 1.5, cannot be recognized anymore. Furthermore the guitar sample can also be recognized in the output spectrogram. With a look onto the embeddings those again show the high energy areas with low or negative values. The original low energy areas in turn have positive values. Furthermore due to the little distance between the higher frequency ranges in the input log-mel spectrograms, the embeddings in this area just depict the high-energy areas. This leads the output spectrogram to not have a high precision in the higher harmonics. Nevertheless through interpolation the resulting embedding contains features of both instruments as well as the resulting spectrogram.

The next figure 5.18, shows the output, of the double-stride network where already significant differences to the output of the single-strided network can be seen. First of all regarding the embeddings, the features representing the higher harmonics of both instruments ($\text{range } >= 15$) are less intense while the features of the lower harmonics, are still distinguishable. Having the interpolated embedding the features of both embeddings can be still recognized, despite the less precision in the y-axis above 15. As a result the

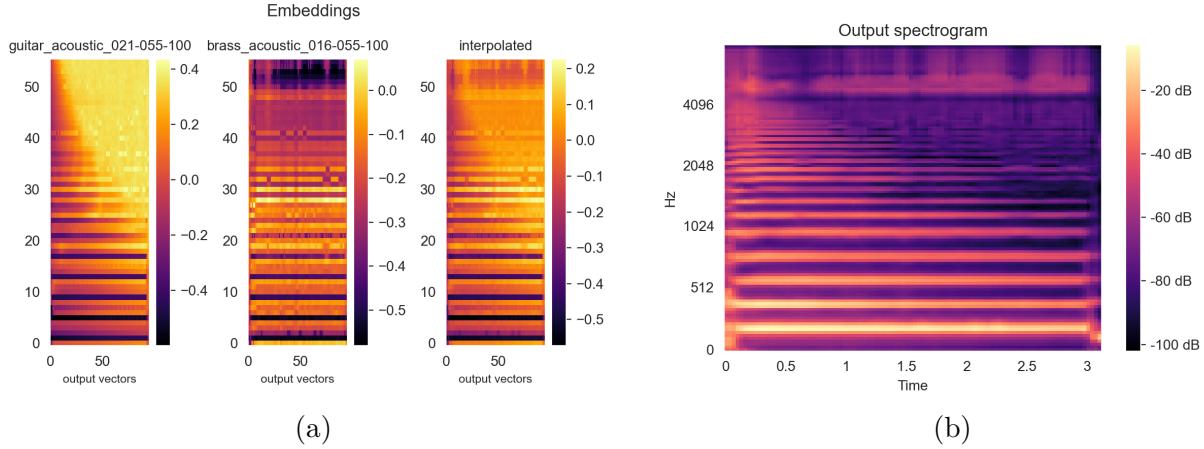


Figure 5.17: embedding interpolation (a), output signal (b).

output spectrogram, the higher harmonics of both instruments are rather “washed out” and appear blurry in the spectrogram.

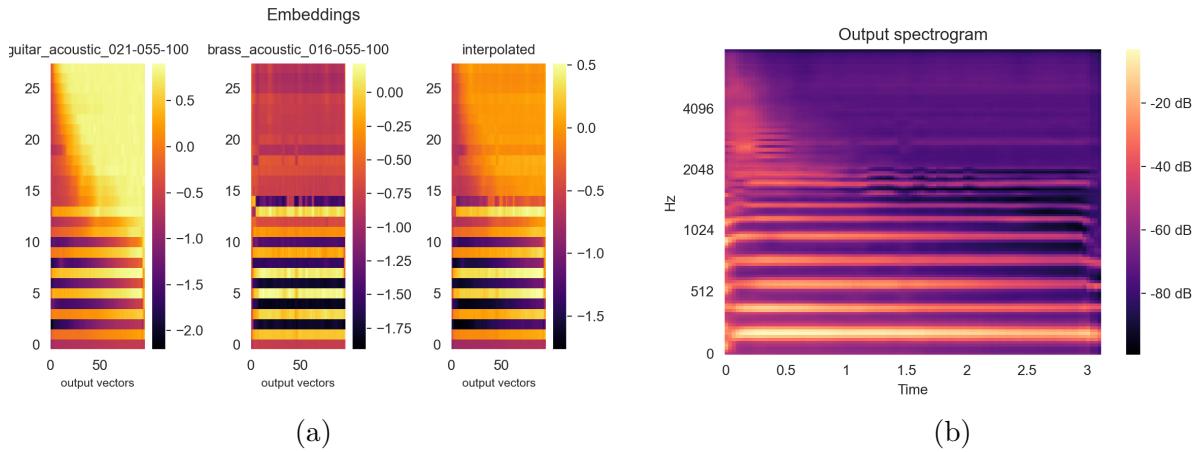


Figure 5.18: embedding interpolation (a), output signal (b).

Coming to the output of the triple-stride network in figure 5.19, here it can be seen that the features contained in the embeddings have no common structure with the input spectrograms. The values therefore cannot be described sufficiently. Also those appear different to the one depicted in figure 5.15b. When looking at the output spectrogram generated of the interpolated embedding, it can be seen that in the area around 1024 to 2048 significant energy changes regarding certain frequencies are present which cannot be recognized in the input. Similar to the output of the previous network in figure 5.18b fine structures of the input are also not present anymore.

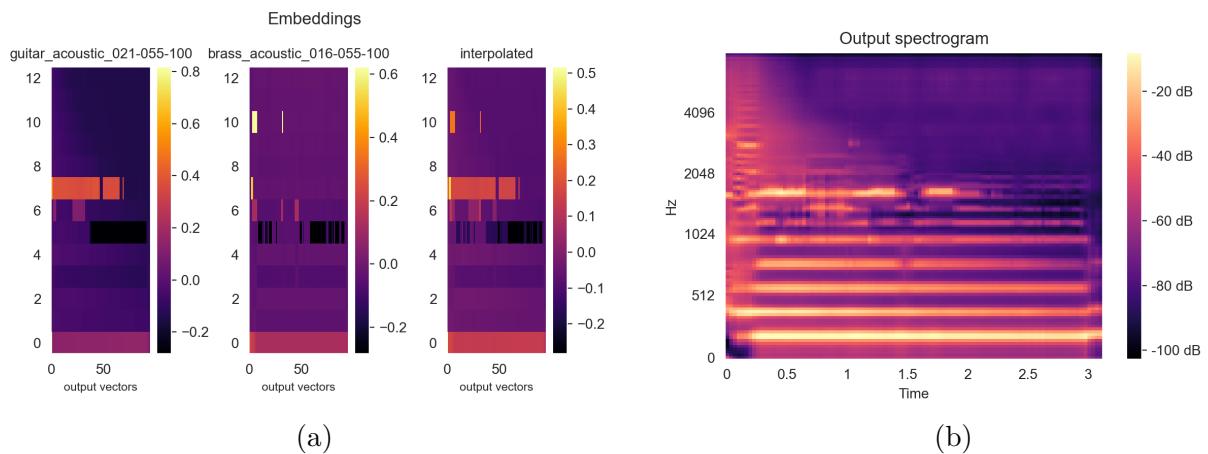


Figure 5.19: embedding interpolation (a), output signal (b).

Chapter 6

Discussion/Evaluation

6.1 Interpretation of results / observations

- 6.1.1 Quality of audio
- 6.1.2 Differences between instruments
- 6.1.3 Impact of model configurations
- 6.1.4 Impact of pre processing
- 6.1.5 Impact of post-processing

6.2 Comparison to other approaches

6.3 Limitations/difficulties

6.4 Outlook

The listenable sounds / output sounds of the networks get discussed at this point.

By listening to the final obtained sample, it can be heard, that it already contains features of both instruments despite the guitar is not as present as expected (guitar stroke not audible). Furthermore the sound contains distortion which is not desirable.

Chapter 7

Conclusion

Chapter 8

Future Work

Appendix A

Additional Graphics

Appendix B

Supplementary Materials

List of supplementary data submitted to the degree-granting institution for archival storage (in ZIP format).

B.1 PDF Files

Path: /

thesis.pdf Master/Bachelor thesis (complete document)

B.2 Media Files

Path: /media

- *.ai, *.pdf Adobe Illustrator files
- *.jpg, *.png raster images
- *.mp3 audio files
- *.mp4 video files

B.3 Online Sources (PDF Captures)

Path: /online-sources

Reliquienschrein-Wikipedia.pdf **WikiReliquienschrein2022**

Appendix C

Questionnaire

Appendix D

LaTeX Source Code

References

Literature

- [1] Keunwoo Choi et al. *A Tutorial on Deep Learning for Music Information Retrieval*. 2018. arXiv: 1709.04396 [cs.CV] (cit. on p. 14).
- [2] Joseph Colonel, Christopher Curro, and Sam Keene. *Autoencoding Neural Networks as Musical Audio Synthesizers*. 2018. eprint: 2004.13172 (eess.AS) (cit. on pp. 4, 5, 23).
- [3] Joseph T Colonel and Sam Keene. “Conditioning Autoencoder Latent Spaces for Real-Time Timbre Interpolation and Synthesis”. In: *2020 International Joint Conference on Neural Networks (IJCNN)*. 2020, pp. 1–7. DOI: 10.1109/IJCNN48605.2020.9207666 (cit. on pp. 4, 5, 23).
- [4] Vincent Dumoulin and Francesco Visin. *A guide to convolution arithmetic for deep learning*. 2018. arXiv: 1603.07285 [stat.ML] (cit. on p. 21).
- [5] Jesse H. Engel et al. “GANSynth: Adversarial Neural Audio Synthesis”. *CoRR* abs/1902.08710 (2019). arXiv: 1902.08710. URL: <http://arxiv.org/abs/1902.08710> (cit. on p. 3).
- [6] Jesse H. Engel et al. “Neural Audio Synthesis of Musical Notes with WaveNet Autoencoders”. *CoRR* abs/1704.01279 (2017). arXiv: 1704.01279. URL: <http://arxiv.org/abs/1704.01279> (cit. on pp. 2, 3, 13, 19, 23, 26, 28).
- [7] Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. “Image Style Transfer Using Convolutional Neural Networks”. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016, pp. 2414–2423. DOI: 10.1109/CVPR.2016.265 (cit. on pp. 7, 8, 11).
- [8] D. Griffin and Jae Lim. “Signal estimation from modified short-time Fourier transform”. *IEEE Transactions on Acoustics, Speech, and Signal Processing* 32.2 (1984), pp. 236–243. DOI: 10.1109/TASSP.1984.1164317 (cit. on pp. 14, 24, 36, 45).
- [9] Eric Grinstein et al. “Audio Style Transfer”. In: *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 2018, pp. 586–590. DOI: 10.1109/ICASSP.2018.8461711 (cit. on pp. 7, 9).
- [10] Lamtharn Hantrakul et al. “Fast and Flexible Neural Audio Synthesis.” In: *ISMIR*. 2019, pp. 524–530 (cit. on p. 3).

- [11] Gerhard Heinzel, Albrecht Rüdiger, and Roland Schilling. “Spectrum and spectral density estimation by the Discrete Fourier transform (DFT), including a comprehensive list of window functions and some new at-top windows” (2002) (cit. on p. 16).
- [12] G. E. Hinton and R. R. Salakhutdinov. “Reducing the Dimensionality of Data with Neural Networks”. *Science* 313.5786 (2006), pp. 504–507. DOI: 10.1126/science.1127647. eprint: <https://www.science.org/doi/pdf/10.1126/science.1127647> (cit. on pp. 13, 20).
- [13] Geoffrey E Hinton. “How neural networks learn from experience”. *Scientific American* 267.3 (1992), pp. 144–151 (cit. on p. 18).
- [14] Sergey Ioffe and Christian Szegedy. *Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift*. 2015. arXiv: 1502.03167 [cs.LG] (cit. on p. 19).
- [15] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. “Perceptual losses for real-time style transfer and super-resolution”. In: *European conference on computer vision*. Springer. 2016, pp. 694–711 (cit. on p. 11).
- [16] Michael I. Jordan and Christopher M. Bishop. “Neural Networks”. *ACM Comput. Surv.* 28.1 (Mar. 1996), pp. 73–75. DOI: 10.1145/234313.234348 (cit. on p. 18).
- [17] joseph colonel joseph, christopher curro christopher, and sam keene sam. “improving neural net auto encoders for music synthesis”. *journal of the audio engineering society* (Oct. 2017) (cit. on pp. 4, 23).
- [18] Diederik P. Kingma and Jimmy Ba. *Adam: A Method for Stochastic Optimization*. 2014. DOI: 10.48550/ARXIV.1412.6980 (cit. on p. 22).
- [19] Xuehao Liu, Sarah Delany, and Susan Mckeever. “Sound Transformation: Applying Image Neural Style Transfer Networks to Audio Spectograms”. In: Aug. 2019, pp. 330–341. DOI: 10.1007/978-3-030-29891-3_29 (cit. on pp. 8, 11).
- [20] Anastasia Natsiou, Luca Longo, and Sean O’Leary. *An investigation of the reconstruction capacity of stacked convolutional autoencoders for log-mel-spectrograms*. 2023. DOI: 10.48550/ARXIV.2301.07665 (cit. on p. 3).
- [21] Keiron O’Shea and Ryan Nash. *An Introduction to Convolutional Neural Networks*. 2015. arXiv: 1511.08458 [cs.NE] (cit. on p. 18).
- [22] Adam Paszke et al. *PyTorch: An Imperative Style, High-Performance Deep Learning Library*. 2019. arXiv: 1912.01703 [cs.LG] (cit. on pp. 20, 27).
- [23] Nathanaël Perraudin, Peter Balazs, and Peter L. Søndergaard. “A fast Griffin-Lim algorithm”. In: *2013 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*. 2013, pp. 1–4. DOI: 10.1109/WASPAA.2013.6701851 (cit. on p. 25).
- [24] Dhruv Ramani et al. “Autoencoder Based Architecture For Fast & Real Time Audio Style Transfer”. *CoRR* abs/1812.07159 (2018). arXiv: 1812.07159. URL: <http://arxiv.org/abs/1812.07159> (cit. on pp. 7, 11, 13, 19).

- [25] Fanny Roche et al. *Autoencoders for music sound modeling: a comparison of linear, shallow, deep, recurrent and variational models*. 2019. arXiv: 1806.04096 [eess.AS] (cit. on pp. 6, 23).
- [26] Stanley Smith Stevens, John Volkman, and Edwin Broomell Newman. “A scale for the measurement of the psychological magnitude pitch”. *The journal of the acoustical society of america* 8.3 (1937), pp. 185–190 (cit. on p. 39).
- [27] Prateek Verma and Julius O Smith. “Neural style transfer for audio spectrograms”. *arXiv preprint arXiv:1801.01589* (2018) (cit. on p. 8).

Software

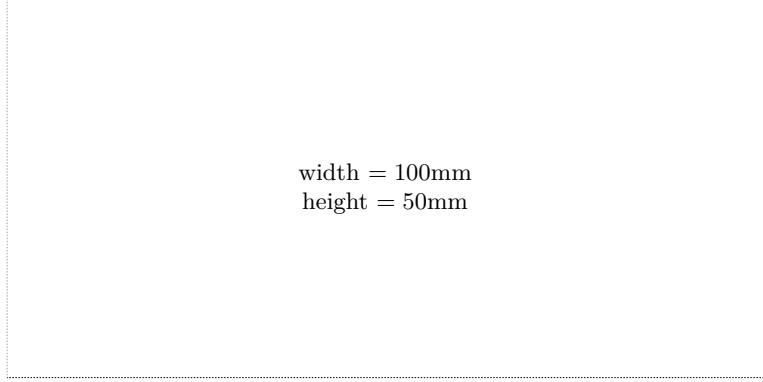
- [28] Brian McFee et al. *librosa/librosa: 0.9.1*. Version 0.9.1. Feb. 2022. doi: 10.5281/zendo.6097378 (cit. on pp. 15, 25, 27).

Online sources

- [29] Dmitry Ulyanov and Vadim Lebedev. *Audio texture synthesis and style transfer*. 2016. URL: <https://dmitryulyanov.github.io/audio-texture-synthesis-and-style-transfer> (visited on 03/14/2023) (cit. on p. 10).

Check Final Print Size

— Check final print size! —



width = 100mm
height = 50mm

— Remove this page after printing! —