# Introducing Latent Timbre Synthesis

Kıvanç Tatar
Simon Fraser University
Vancouver, B.C., Canada
ktatar@sfu.ca

Daniel Bisig
Zurich University of the Arts
Zurich, Switzerland
daniel.bisig@zhdk.ch

Philippe Pasquier
Simon Fraser University
Vancouver, B.C., Canada
pasquier@sfu.ca

June 2, 2020

**Abstract**

We present the Latent Timbre Synthesis (LTS), a new audio synthesis method using Deep Learning. The synthesis method allows composers and sound designers to interpolate and extrapolate between the timbre of multiple sounds using the latent space of audio frames. We provide the details of two Variational Autoencoder architectures for LTS, and compare their advantages and drawbacks. The implementation includes a fully working application with graphical user interface, called *interpolate_two*, which enables practitioners to explore the timbre between two audio excerpts of their selection using interpolation and extrapolation in the latent space of audio frames. Our implementation is open-source, and we aim to improve the accessibility of this technology by providing a guide for users with any technical background.

## 1 Introduction

Promising new research of Deep Learning (DL) for the musical applications of audio transformation and sound synthesis has recently emerged in academia; in conjunction with the increasing popularity of Deep Learning architectures [2]. Musical applications of these technologies have yet to become accessible for composers and musicians who lack expertise in Machine Learning (ML) and Artificial Intelligence (AI). Although Deep Learning has been applied to many musical tasks [2], the research on incorporating Deep Learning architectures for sound design applications in experimental electronic music is still in its early stages.

This project focuses on the integration of modern ML and AI techniques into tools for computer-assisted sound design and their applications within composition practices. We concentrate on the fields of experimental electronic music and sound art, which takes sound qualities, listening modalities and experiences [38], and organized sound theories [28] into the center [36]. We specifically focus

on audio corpus-based sound synthesis approaches that rely on large libraries of audio excerpts.

Latent Timbre Synthesis (LTS) aims to help composers by utilizing an abstract latent timbre space that is generated by training unsupervised Deep Learning (DL) models with a set of audio recordings. These new DL tools allow composers to synthesize sounds using a latent space of audio that is constrained to the timbre space of the audio recordings in the training set. Our development of ML and AI tools for computer-assisted sound design and their subsequent evaluation within composition practice serve to highlight the benefits and shortcomings of the selected machine learning algorithms for creative ideation and discovery.

In 1940s, physicist Dennis Gabor proposed [7] that a sound is composed of acoustical quanta that is bounded by time and frequency. We are inspired by this idea while applying this approach to digital audio, and asking,

- What would be the latent space of audio frames?

- Can we regenerate original audio recordings using that latent space?

- Can we create new audio synthesis methods using the latent space for sound design applications and composition practices?

- Can a DL-based model be used to provide professional grade sound synthesis tools?

- How can an audio synthesis architecture using a Deep Learning model provide the user the flexibility to generate sounds of any duration?

We are also inspired by the definition of music as "nothing but organized sound" [39] involving sound objects [29] that situate on multiple layers [32], where any sound can be used to produce music [19, 39], and strong connections exist between pitch, noise, timbre, and rhythm [32, 31, 27, 28]. In that sense, the Latent Timbre Synthesis project builds on our previous work titled Musical Agents based on Self-Organizing Maps (MASOM) [35, 37]. MASOM combines organizing sound samples in latent audio space with statistical sequence models for musical structure. The latent audio space in MASOM is generated by a Self-Organizing Map that organizes a set of audio excerpts. In LTS, we move further by aiming for an audio synthesis framework where we can synthesize sounds that do not exist in the training set.

Following the research questions and directions above, our contributions presented in this paper include two Variational Auto-encoders (VAEs) to generate a latent space of audio frames. Unlike other Deep Learning architectures such as Generative Adversarial Networks (GANs), VAEs are beneficial for our applications because these architectures can encode an existing audio frame to a latent space, as well as synthesize audio frames from latent vectors. VAEs also allow audio synthesis through interpolation and extrapolation of timbres, by using the latent vectors of audio frames.

Latent Timbre Synthesis differs from the previous works such as Granma MagNet[1] because we prioritize the flexibility to generate audio with any duration, in comparison to outputting audio excerpts of fixed-duration. We think that the flexibility of changing the duration of the generated audio is crucial

for our applications, which stands out as another contribution of LTS. Our approach focus on creating a latent space of audio frames, where we can represent an audio recording with any length as a time-series sequence of latent vectors.

The LTS framework consists of three main modules, calculation of wavelet transform based spectrogram representation, latent audio frame space generation using two specific Variational Auto-encoders and inverse synthesis using wavelet-based magnitude spectrogram generated by the decoder of the VAE. We compare the advantages and drawback of two VAE architectures for designing a synthesis tool for composition practices and sound design applications. In addition, we present and share a fully working application, called *interpolate_two* with a Graphical User Interface (GUI) that allows composers to synthesize audio using timbre interpolation and extrapolation with multiple sounds. In comparison to high computational complexity of previous works mentioned in Section 2, the low computational complexity of *interpolate_two* allows the incorporation of the sound design tool within composition practices and real-time applications. The documentation of the setup of *interpolate_two* is detailed to guide practitioners of all backgrounds. Our implementation is open-source[1], and sound examples are available[2]. We encourage our readers to dive into the code and experiment with the framework for further audio synthesis possibilities.

## 2    Related Works

We limit this section to the previous works that utilize audio spectrogram as an input for the Deep Learning architecture, with the exception of WaveNet. We situate the Latent Timbre Synthesis project within the raw-audio generation applications of Deep Learning, and WaveNet is one of the state of the art systems in the area. We also omit Deep Learning systems for speech synthesis or vocoder applications, such as MelGAN [16], while mentioning in Section 5 how we plan to incorporate them in LTS as a next step in our research.

WaveNet is a Deep Learning architecture that uses an audio corpus for the tasks of music composition, multi-speaker speech and text to speech generation, and speech recognition [24]. WaveNet applies Convolutional Neural Networks (CNNs) with two strategies to handle temporality of raw audio data: causal convolution and dilation. Causal convolutions ensure that the output only depends on the past observations. Dilated causal convolutions skips a number of inputs on each layer. The number of inputs that are skipped exponentially increases with each layer; hence, the receptive field of the network also increases exponentially [41]. Note that, the receptive field is the number of neurons that affect a single neuron in deep networks. Oord et al. [24] tested WaveNet on two audio corpora: the MagnaTagATune dataset and the YouTube piano dataset. The authors point out that "Even with a receptive field of several seconds, the models did not enforce long-range consistency which resulted in second-to-second variations in genre, instrumentation, volume and sound quality." That is, WaveNet struggled to generate long-term variations like in the case of interactive music systems that apply Markov Models [36, Section 6.1]. There has been follow-up research on the WaveNet architecture, where the authors stack mul-

---

[1]The source code is available at `https://www.gitlab.com/ktatar/latent-timbre-synthesis`.

[2]We provide sound examples at `https://kivanctatar.com/Latent-Timbre-Synthesis`.

tiple WaveNet architecture on top of each other [25], or they combine WaveNet with Vector Quantized Variational Autoencoders [4]. The main drawback of all WaveNet systems are their computational complexity and high-usage of GPU memory [3]. The technology requirements of WaveNet compromise its usage in compositional practices, where composers do not necessarily have access to computers with the state of the art GPUs.

Differentiable Digital Signal Processing (DDSP) is a toolbox made by the Google for researching Digital Signal Processing (DSP) applications of Deep Learning [5]. The authors describe the DDSP Autoencoder, which is a VAE architecture where the input are the Mel-Frequency Cepstral Coefficients (MFCCs) of an audio excerpt. We mention a comparison of using MFCCs and other audio features as the representation of timbre for audio synthesis with VAEs in Section 3.1. The architecture employs three autoencoders for fundamental frequency (*f-encoder*), loudness (*l-encoder*), and the latent space of timbre (*z-encoder*). The fundamental frequency and the loudness encoders use the CREPE architecture that is originally presented as a pitch detector [10]. The *z-encoder* architecture is inspired by the ResNet architecture in Computer Vision research [11]. The decoder, on the other hand, controls the input parameters of an additive synthesis module, a subtractive synthesis module, and a reverb. These three synthesis modules generate the final audio. The loss function compares the generated audio with the original one, using a specific function called Multi-Scale Spectrogram Loss, which is similar to comparing the spectrograms of original and generated audio.

Generative Timbre Spaces project [6] is perhaps one of the most similar previous study to the Latent Timbre Synthesis project. The application of Generative Timbre Synthesis focuses on generating a latent timbre space of conventional musical instruments. This model uses a VAE where the encoder is a 3-layer feed-forward network with 2000 units in each layer. The latent space has 64 dimensions. The authors introduce a new regularization item in the cost function. The additional regularization loss tries to force the network to satisfy perceptual similarity ratings of conventional musical instruments in Western Classical Music. These perceptual ratings are proposed in previous studies [8, 15, 12, 20, 17]. The training dataset of Generative Timbre Spaces is audio recordings of conventional musical instruments where each file is an instrument playing a note. The authors takes one frame from each audio file to train the VAE model. Hence, the architecture aims to capture the generalized timbre of a conventional musical instrument instead of the regeneration of an arbitrary audio excerpt. Likewise, the cost function is not suitable to regenerate a dataset with arbitrary audio recordings because there are no perceptual ratings available. We further discuss the issues related to the hyper-parameters of VAE in Generative Timbre Spaces in Section 3.2.

The DDSP Autoencoder as well as the Generative Timbre Spaces aims for the synthesis applications of conventional music where the model is conditioned to output an audio with a fundamental frequency constraint. In sound design, experimental electronic music, and Sound Art applications, having a fundamental frequency of a sound gesture is rather limiting. The music theory of the contemporary electronic music emphasizes the continuum between noise, pitch, and rhythm [32, 19, 31, 27, 28]. Hence, in LTS, we aim for a model that could generate any audio so that the composers and practitioners are free to explore the full potential of digital audio synthesis.

4

Table 1: A previous study provided a comparison of audio frame reconstruction losses with Variational Autoencoders using spectrograms with fixed-length windows and wavelet transform based spectrograms [6].

|  | Spectrogram | $logp(x)$ | $\|\|x - \tilde{x}\|\|^2$ |
|---|---|---|---|
| Fixed Window | STFT | -1.9237 | 0.2412 |
|  | DCT | 4.3415 | 2.2629 |
| Wavelet Transform | CQT | -2.8723 | 0.1610 |
|  | NSGT-MEL | -2.9184 | 0.1602 |
|  | NSGT-ERB | -2.9212 | 0.1511 |

## 3    System Design

The Latent Timbre Synthesis framework consists of three main parts, spectrogram calculation using a type of wavelet transform, audio frame latent space generation using Variational Autoencoders, and inverse synthesis of audio using the magnitude spectrogram generated by the decoder of VAE (Figure 1).

### 3.1    Wavelet Transform based Spectrograms

The audio feature extraction module generates spectrogram frames using a type of wavelet transform for audio, called Constant-Q Transform (CQT) [30], which gained popularity in Music Information Retrieval (MIR) research in the recent years. CQT, as well as its variant Non-stationary Gabor Transform (NSGT) [40], have been compared to the other audio features such as Mel-Frequency Cepstral Coefficients (MFCC); and previous studies showed that CQT and NSGT could perform better in MIR applications such as segmentation and musical structure analysis [23]. Naturally, segmentation and musical structure analysis tasks require computing the audio similarity [22]; thus, they are suited to create a latent space of audio frames.

A previous work [6] compared spectrograms computed with fixed windows and wavelet transform based spectrograms for applications of latent audio frame space generation using Deep Learning (DL). This comparison included Short-Time Fourier Transform, Discrete Cosine Transform, Constant-Q Transform (CQT), and Non-Stationary Gabor Transform (NSGT) variations using different frequency scales. The wavelet based transforms in this previous study were CQT and NSGT variants. The authors [6] found that wavelet transform based spectrogram representations perform better than the spectrograms calculated using fixed-length windows in regards to the log-likelihood and mean quality of the audio frame reconstructions, as shown in Table 1; while the audio frame reconstructions of wavelet transform based spectrograms gave similar results. We utilize CQT in comparison to other wavelet-based spectrograms in Table 1 because a python library for audio analysis, titled Librosa[3] [21], includes a CQT and inverse CQT implementation [30] combined with a Fast-Griffin-Lim phase estimation [26] that we explain in Section 3.3. We aimed that the LTS framework would be available for composers and sound designers of all technical
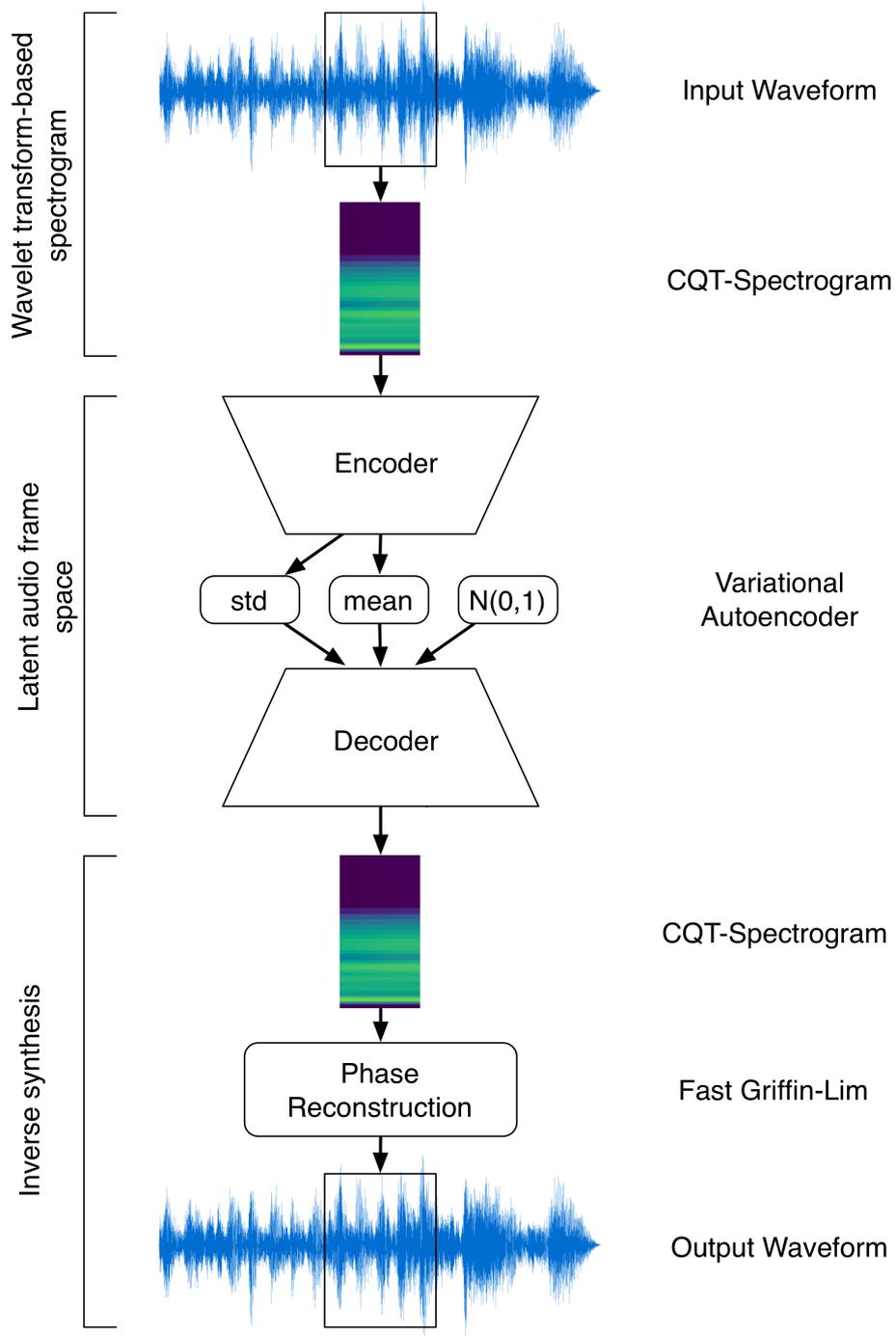
---

[3]https://librosa.github.io/librosa/

Figure 1: Latent Timbre Synthesis framework

levels. Hence, we prioritized the options that simplify the installation of the LTS framework.

We delve into the details of Constant-Q Transform in the following. We can calculate the CQT of an audio recording [30], a discrete time domain signal $x(n)$, using the following formula:

$$X^{CQ}(k, n) = \sum_{j=n-\lfloor N_k/2 \rfloor}^{n+\lfloor N_k/2 \rfloor} x(j) a_k^*(j - n + N_k/2) \tag{1}$$

where $k$ represents the CQT frequency bins with a range of $[1, K]$, and $X^{CQ}(k, n)$ is the CQT transform. $N_k$ is the window length of a CQT bin, that is inversely proportional to $f_k$ that we define in equation 3 Notice that, $\lfloor \cdot \rfloor$ is the rounding towards negative infinity. $a_k^*$ is the negative conjugate of the basis function $a_k(n)$ and,

$$a_k(n) = \frac{1}{N_k} w(\frac{n}{N_k}) exp[-i2\pi n \frac{f_k}{f_s}] \tag{2}$$

where $w(t)$ is the window function, $f_k$ is the center frequency of bin $k$, and $f_s$ is the sampling rate. CQT requires a fundamental frequency parameter $f_1$, which is the center frequency of the lowest bin. The center frequencies of remaining bins are calculated using,

$$f_k = f_1 2^{\frac{k-1}{B}} \tag{3}$$
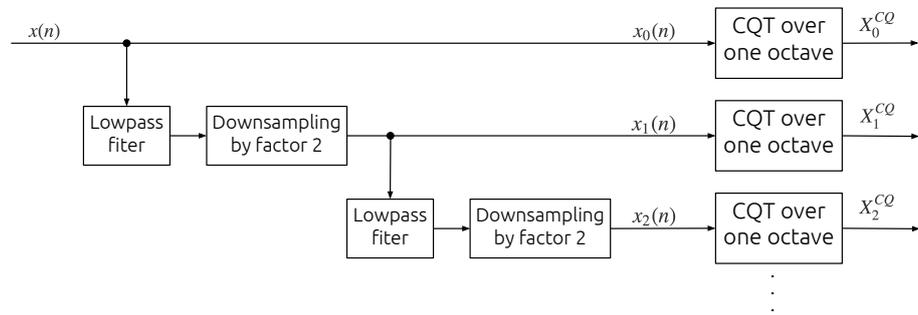
where $B$ is the number of bins per octave.

CQT is a wavelet-based transform because the window size is inversely proportional to the $f_k$ while ensuring the same Q-factor for all bins $k$. We can calculate the Q-factor using,
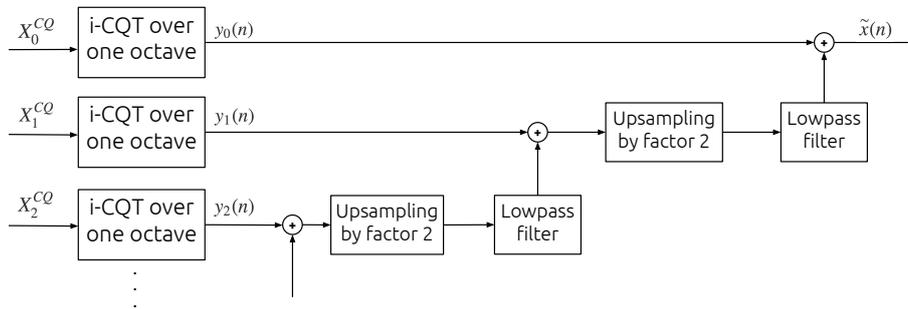
$$Q = \frac{q f_s}{f_k (2^{\frac{1}{B}} - 1)} \tag{4}$$

where $q$ is scaling factor with the range [0,1] and equals to 1 as the default setting. We direct our readers to the original publication for the specific details of the CQT [30], which also proposed a fast algorithm to compute CQT and inverse CQT (i-CQT), given in Figure 2.

All experiments in this paper used the same audio feature extraction configurations. We tried several parameters to find the configuration that could give the least amount of audio artifacts with the pipeline of calculating the CQT spectrogram and then reconstructing the audio back using the magnitude spectrogram of CQT combined with the phase estimation algorithm. Notice that these artifacts would appear even with an ideal DL model because the inverse audio synthesis introduces these artifacts to the LTS. We used 16-bit and 44.1 kHz stereo or mono audio recordings. We converted the stereo files to mono first, and then calculated the CQT spectrograms using a hop-size of 128 samples. The $f_1$ parameter was 32.7 Hz that corresponds to the musical note, C1. $q$ value in equation 4 was equal to 1, and the window function was "hann"[4]. CQT included 48 bins per octave for a total range of 8 octaves; which sums up

---

[4]https://docs.scipy.org/doc/scipy/reference/generated/scipy.signal.windows.hann.html

(a) Constant-Q Transform



(b) Inverse Constant-Q Transform

Figure 2: A fast algorithm to compute CQT and i-CQT, described in [30] and implemented in Librosa [21].

to a 384 bins in total. Hence, the input of the DL models are vectors with 384 dimensions. These parameters resulted in the least amount of artifacts in our experiments. Our source code is flexible to change these parameters, and we encourage our readers to try different parameter configurations.

## 3.2 Deep Learning in Latent Timbre Synthesis

### 3.2.1 Autoencoders and Variational Autoencoders

Autoencoders are Deep Learning architectures for generative modelling. The architecture consists of two main modules: an encoder and a decoder (Figure 1 and 4). The encoder maps the input data $x \in R^L$ to a latent vector $z \in R^M$ where $z = encoder(x)$, and $M < L$. The decoder aims to convert a latent vector back to the original data, and ideally, $decoder(encoder(x)) = x$. The vanilla Autoencoder architecture encodes the input data vector to a single point, that is the latent vector. In comparison, Variational Autoencoder (VAE) is an improved version of the Autoencoder architecture that converts the input data vector to a stochastic distribution over the latent space. This difference is also referred as the "reparametrization trick" [13, 14, 33].

In VAE, the encoder tries to generate a latent space by approximating $p(z|x)$ while the decoder tries to capture the true posterior $p(x|z)$. The vanilla VAE approximates $p(z|x)$ using $q(z|x) \in Q$ with the assumption that $p(z|x)$ is in the form of a Gaussian distribution $N(0, I)$. This approximation is referred in the literature as *Variational Inference* [13]. Specifically, the encoder outputs the mean $\mu_M$ and the co-variance $\sigma_M$ as the inputs of the Gaussian distribution function $N(z; \mu_M, \sigma_M^2 I)$ over a latent space with $M$ number of dimensions. Hence, the encoder approximates $p(z|x)$ using $q^*(z|x) = N(z; f(x), g(x)^2 I)$ where $\mu_M = f(x)$, $f \in F$, $\sigma_M = g(x)$, and $g \in G$. The decoder's input, the latent vector $z$ is sampled from the latent distribution $q(z) = N(z; f(x), g(x)^2 I)$. Hence, the loss function consists of the reconstruction loss and the regularization term of Kullback-Leibler divergence between $q^*(z|x)$ and $p^*(z)$,

$$L_{f,g} = \mathbb{E}_{q^*(z)}[log p^*(x|z)] - \alpha \cdot D_{KL}[q^*(z|x)||p^*(z)] \tag{5}$$

We direct our readers to the original VAE publication [13] for the mathematical induction of the loss function in equation 5 [14]. Note that, some previous works introduced additional regularization terms to the loss function to condition the VAE further, such as the introduction of perceptual ratings of musical instruments in [6].

The LTS framework focuses on Variational Autoencoders in comparison to Generative Adversarial Networks because we aim for audio synthesis by interpolation and extrapolation in the latent space of audio frames (see Section 4). The input vectors of the VAE model are CQT vectors calculated from one audio frame where the window size varies with the frequency bins. We aim to generate a latent space of audio frames so that we can synthesize audio with any duration. Previous systems such as Grannma MagNet [1] utilizes audio excerpts with fixed-duration, where the training input vectors of deep learning model are 2D audio spectrograms with time along the x-axis and frequency along the y-axis. This design choice constraints these DL models to limited applications such as generating a fixed-length audio excerpt. Our approach differs from the previous systems because the training observation of DL model in LTS is one

audio-spectrum vector that is calculated from one audio-frame. This allows LTS to generate audio with any duration in sound design applications.

We focus on two Deep Learning (DL) architectures in the current version of the Latent Timbre Synthesis framework. Both models are Variational Autoencoders (VAE); however, the layers and model parameters differ. In the first VAE model, we were inspired by the previous work [6] where the authors trained a Variational Autoencoder to generate conventional musical instrument timbres with digital audio synthesis.

We initially tried the VAE architecture with the same hyper-parameter settings that were used in the Generative Timbre Spaces project. The setting of Generative Timbre Spaces [6] were unsuccessful in our experiments. The model could not learn to regenerate the audio recordings in the training dataset, and could only generate noise. Upon further investigation, we found that the Kullback-Leibler Divergence regularization term multiplier caused the issue. In Generative Timbre Spaces, the authors increase the multiplier from 0 to 2 during the first 100 epochs of the training, following the warm-up procedure [33]. We suspect that the KLD multiplier range of $[0, 2]$ is specific to the application of Generative Timbre Spaces where the training dataset consists of audio files with distinct harmonic content and low noisiness in the spectrum. Furthermore, the training dataset size of Generative Timbre Spaces is rather small, less than 100 MB, whereas we work with GBs of audio to train the LTS models. For example, the *erokia* dataset that we provide with our source code includes 2 GBs of audio that corresponds to $3,084,591$ audio frames as training data for LTS models, using a hop-size of 128 samples for the CQT calculation.
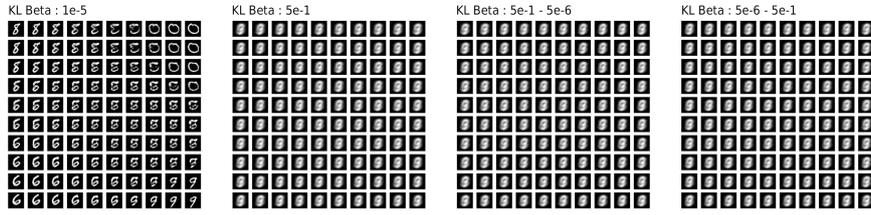
We found that the range of $[0, 2]$ for KLD multiplier was too high for our application and prohibited the VAE to learn. In addition, the warm-up procedure had adverse effects on the learning. We further investigated this issue using the MNIST dataset [18]. MNIST gave us a visual understanding of the effect of KLD multiplier. The images in MNIST are $28 * 28$ pixels, adding up to a total of 784 pixels. This is similar to LTS where the input of the VAE is a vector of 384 dimensions.

Like in the case of our tests with audio spectrograms, we obtained similar results with the MNIST data. Higher KL-divergence values as well as the warm-up procedure significantly deteriorated the reconstructions of the trained model.
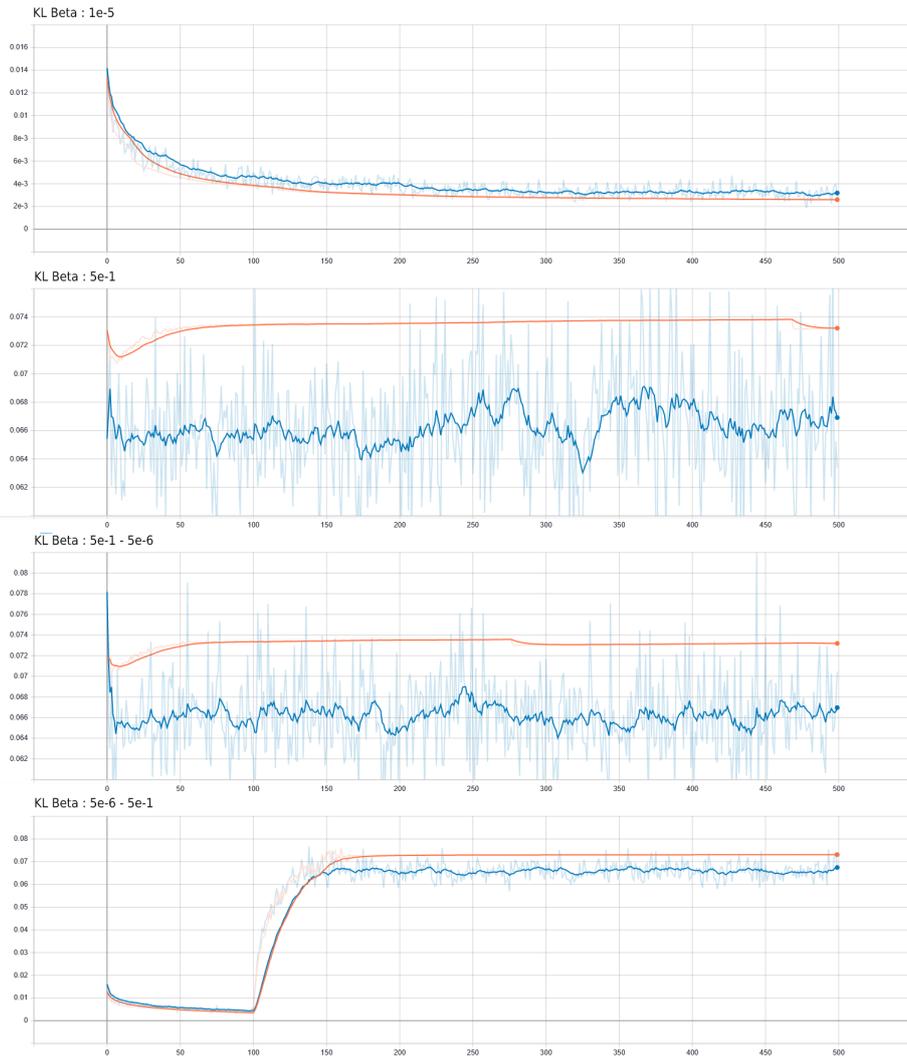
Figure 3a and 3b show the effect of KLD multiplier on the training of the Variational Autoencoder. We used the same architecture depicted in Figure 1 and only changed the input and output dimensions to 784 that corresponds the flattened vector of $28x28$ images in the MNIST dataset. In addition, we tested the warm-up procedure and the reverse settings of the warm-up procedure; however, both were detrimental to the training in our experiments. We proceeded our experiments with KLD values around $1e - 5$, given the success of this setting with the MNIST dataset.

All VAE architectures in LTS use decoder networks that are the reversed replicas of the encoder networks, as in most cases of VAEs. The first VAE model in LTS is a feed-forward network with two Dense layers including 2048 neurons, The *dense, dense_1,* and *dense_2* layers in Figure 4 apply Rectified Linear Units (ReLU) as the neuron activation functions.

We train the network for 2000 epochs, while the improvements after epoch

(a) The effect of the KLD multiplier on the reconstructions



(b) The loss values during training where x-axis is the training epoch and y-axis is the loss value

Figure 3: The effect of the KLD multiplier in the loss function on Variational Autoencoder training
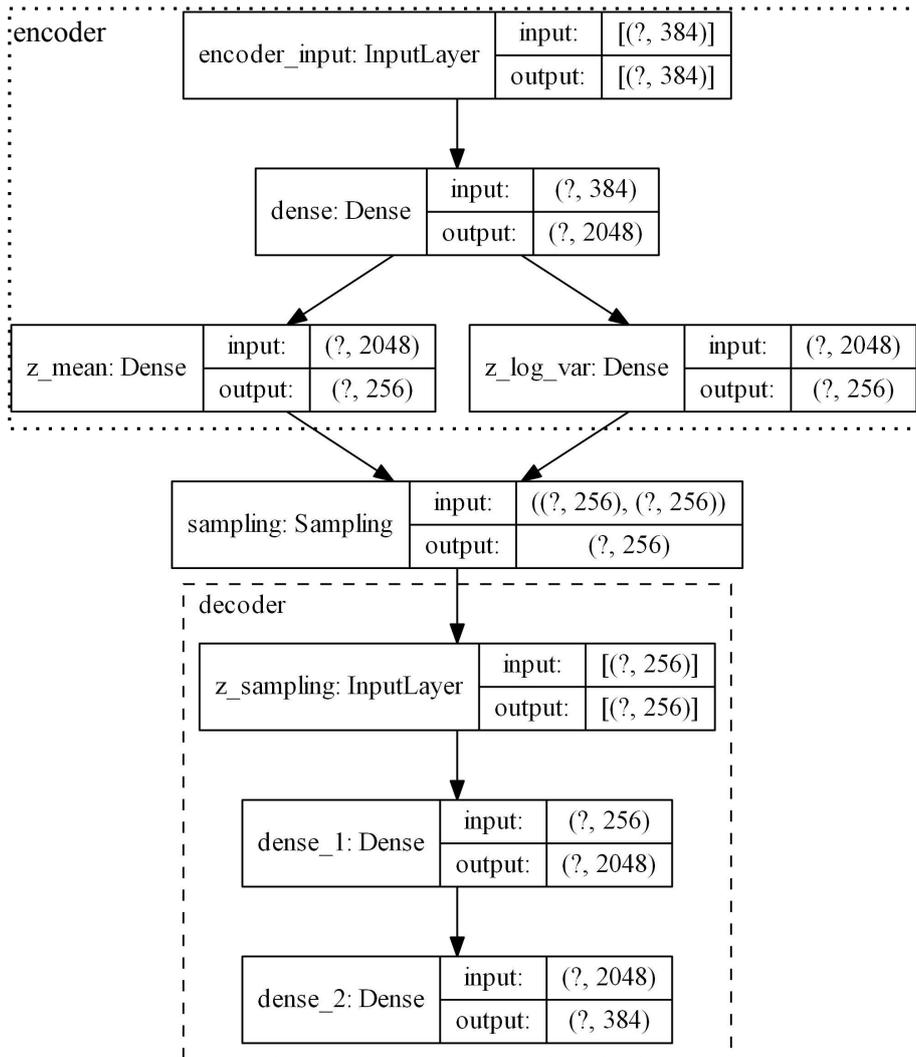
Figure 4: The first Deep Learning architecture available in Latent Timbre Synthesis
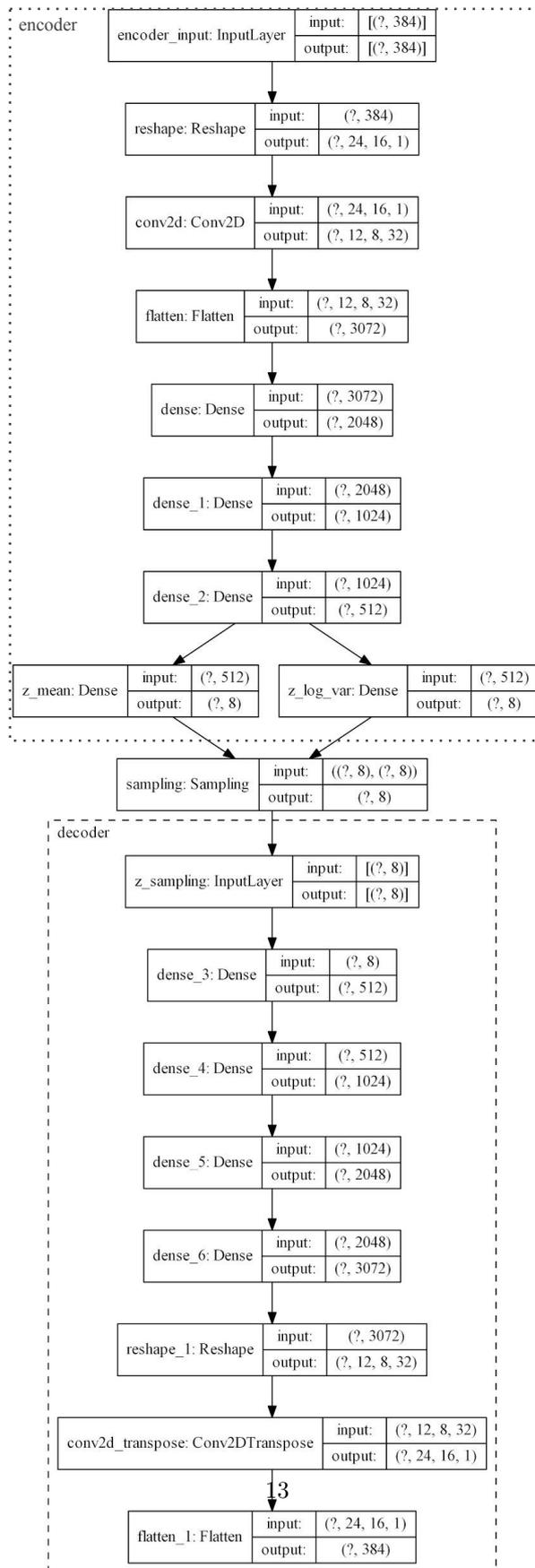
Figure 5: The second Variational Autoencoder architecture available in Latent Timbre Synthesis

50 are rather in the exploitation[5] phase of the learning, and help to minimize the floor noise in the generated CQT magnitude spectrogram. The learning rate is $1e^{-4}$ and the KLD multiplier in the cost function is $5e^{-4}$. The learning rate and the KLD multiplier parameters are dependent on the training dataset. We recommend the readers who would be interested to try their own dataset to start with the hyper-parameter settings above, and change the parameters when needed.

The latent space of our first VAE model consists of 256 dimensions. This number is still acceptable given that the previous work [6] used a 64 dimensional latent space to cluster a much smaller range of conventional musical instrument timbres. Yet, we explored the possibility to find a deeper network that could generate a latent space with a smaller number of dimensions.

The second VAE network that we present in this paper aimed to decrease the number of latent space dimensions of the first network using a deeper architecture, shown in Figure 5. We pursued an iterative design procedure where we tried to decrease the number of dimensions of the latent space while maintaining the reconstruction quality and trying to achieve lower final loss values in the training runs. Our experiments aimed for a 8 dimensional latent space, and we first tested increasing the number of dense layers. Our experiments found an optimum of 4 dense layers where the number of neurons is changing by two-folds each layer (Figure 5). In addition to the Dense layers, we also tried adding convolutional layers on top of the Dense layers. We imagined that the convolutional layers could grasp the relationships between frequency bins of a given CQT vector, such as harmonics. The convolutional layer generates a 2-dimensional vector with size 24x16, which is in relation to 48 bins per octave for 8 octaves while trying to remain as close as to a square. Notice that, each row in the 2D vector generated by the convolutional layer corresponds to half of an octave. We imagined that this could further help the network to capture the harmonic content in the CQT vectors. Our experiments found an optimum of one additional convolutional layer with 32 filters and a kernel size of 3 and a stride size of 2 on top of the dense layers, shown in Figure 5. The final second network could decrease the number of latent space dimensions while giving a low final loss value. However, the network introduced a consistent floor noise sound in the reconstructions. The introduction of normalization techniques such as batch, instance, or layer normalization could not eliminate the floor noise in the reconstructions. Our findings suggest a balance between the number of latent space dimensions and the floor noise in the reconstructions where increasing the number of latent space dimensions is the solution to eliminate the noise.

Given that audio quality is of great importance in composition tasks, we focused on using and disseminating the first network given in Figure 4. Additionally, the low computational complexity of the first network can be an advantage when we combine the VAE with Deep Learning models for time-series sequence generation in our future work, detailed in Section 5.

## 3.3 Inverse Synthesis and Audio Reconstruction

This version of LTS uses the Fast Griffin-Lim [26] phase estimation algorithm (GLA) for generating audio from CQT magnitude spectrograms that the VAE

---

[5]Exploration and exploitation are two search strategies in optimization applications [34, Section 5.3].

decoder outputs. Briefly, the GLA estimates the phase component of a magnitude spectrogram by iterating the inverse synthesis and the spectrogram calculation multiple times, initially proposed in [9] and shown in Algorithm 1 [26]. Given an an audio signal $x(n)$ and its frequency transform $X(i)$,

---

**Algorithm 1** Griffin-Lim Algorithm

---

1: **Set:** $\angle X_0(i)$
2: **Initialize:** $X_0(i) = |X(i)| \cdot e^{j\angle X_0(i)}$
3: **for** $n = 1, 2, \ldots, N$ **do**
4: $\quad X_n(i) = T(IT(|X(i)| \cdot e^{j\angle X_{n-1}(i)}))$
5: **end for**
6: $\hat{x}(n) = IT(X_N(i))$

---

where $N$ is the total number of GLA iterations, $T$ and $IT$ is the frequency transform and inverse frequency transform function respectively; such as Short-Fourier Transform, or Constant-Q Transform in our case. Note that, the space of audio spectrograms is a subset of the complex number space. The iterative process of Griffin-Lim moves the complex spectrogram of the estimated signal $\hat{x}(n)$ towards the complex number space of audio signals in each iteration, as proven in [9].

The Fast Griffin-Lim algorithm (F-GLA) [26] is a revision of the original Griffin-Lim algorithm.

---

**Algorithm 2** Fast Griffin-Lim Algorithm

---

1: **Set:** $\angle X_0(i)$
2: **Initialize:** $X_0(i) = |X(i)| \cdot e^{j\angle X_0(i)}$
3: **Initialize:** $Y_0(i) = T(IT(|X(i)| \cdot e^{j\angle X_0(i)}))$
4: **for** $n = 1, 2, \ldots, N$ **do**
5: $\quad Y_n(i) = T(IT(|X(i)| \cdot e^{j\angle X_{n-1}(i)}))$
6: $\quad X_n(i) = Y_n(i) + \alpha(Y_n(i) - Y_{n-1}(i))$
7: **end for**
8: $\hat{x}(n) = IT(X_N(i))$

---

where $\alpha$ is a constant. A previous study [26] showed that the F-GLA revision significantly improves signal-to-noise ratio (SNR) compared to the GLA, where the setting $\alpha = 1$ resulted in the highest SNR value.

The F-GLA module has the highest computation time in the LTS framework. We are aware that a revision of this module using another Deep Learning architecture for vocoder applications can improve the computational complexity of LTS while making LTS more lightweight within real-time applications. We further address this in the Section 5.

This concludes the explanation of the LTS architecture, where we covered our specific choice of Constant-Q Transform, two Variational Autoencoder architectures, and the inverse synthesis using the GLA phase estimation. In the following, we introduce the *interpolate_two* algorithm.

Figure 6: The Max GUI of the *interpolate_two* application

# 4 Interpolation and Extrapolation in Latent Timbre Space

The first sound design application of Latent Timbre Synthesis is the *interpolate_two* framework that allows composers to synthesize sounds using interpolation and extrapolation with two sounds. The framework requires a trained model to generate sounds using the latent audio frame space. The user can select the duration of the generated sound, and can choose an excerpt from two audio files. These two excerpts have the same duration. The algorithm uses these two excerpts for synthesis with interpolation and extrapolation in the timbre space. The interpolation amount sets how much of the latent vector is copied from one of the audio excerpts. For example, 30% interpolation is adding 30% of the latent vectors of the first audio and 70% of the latent vectors of the second audio. The percentages above 100 or below 0 corresponds to extrapolations. For instance, 120% is moving 20% away from the second audio in the direction the latent vector that points from audio 1 to audio 2. The algorithm synthesizes the audio by calculating every audio frame using inverse synthesis from a generated spectrogram. Hence, the user sets interpolation amounts for each latent vector that corresponds to one audio frame. The user can draw an interpolation curve to change the interpolation percentage in time using the LTS framework.

The application, *interpolate_two* consists of two components: Max GUI and the python engine. The Max GUI handles user interactions while the python engine reacts to OSC messages coming from Max. The python engine runs the deep learning model and audio feature extraction (CQT calculations), as well as inverse synthesis that generates audio from CQT magnitude spectrogram combined with Fast Griffin-Lim phase estimation.

The framework of *interpolate_two* is compatible with all variations of the VAE architectures given in Figure 4 and 5. The B an C regions in the GUI in Figure 6 load the datasets and models that are produced within a particular run with the dataset. Using the regions D and E, the user can select two audio files from a dataset to choose excerpts. The file dropdown menu allows to select an audio file from the dataset. "Zoom to selection" sets the view to the selection area. Clicking *ctrl* (or *cmd* on macos) and then dragging the mouse up & down on the waveform views applies zoom in & out. The waveform in the region F sets apply interpolation (or extrapolation) amounts per frame, where x-axis is the time and y-axis sets the interpolation percentage for a frame. When both waveform views are zoomed to the selection, the x-axis of the interpolation curve corresponds to the x-axis of the waveform. The interpolation curve view is a [waveform˜] object. The "Vertical Zoom" parameter in the inspector of this Max object sets the maximum interpolation/extrapolation amount. The default maximum is 1.3; hence, [1.0,1.3] and [-1.0, -1.3] are the extrapolation regions. It is possible to extrapolate even more by changing the vertical zoom parameter; however, the higher amounts are likely to give audible distortions. The normalize toggle in region H allows the user to normalize the generated audio to prevent audio distortions or extremely high audio volumes while exploring extrapolation possibilities.

The section H send messages to the python engine to handle output generation. "Generate & Play" initiates the python engine to synthesize a sound using the current interpolation curve and the audio selections. "Play Again" plays the previous generated sound, without going through the deep learning calculation. "STOP" immediately stops the audio coming out of the python engine. Phase iterations sets the number of iterations of the Fast Griffin-Lim algorithm. Higher number of iterations (max. 64) gives better results; however, the calculation takes significantly longer. The phase estimation algorithm is the bottleneck of computational complexity of this framework. Still, the calculation of the audio takes 50% of the audio duration with phase iteration set to 1. That is, calculating a 2-second sound takes around 1-second on a laptop with NVIDIA RTX 2080 Max-Q GPU and 2.20 GHz Intel i7-8750H CPU.

# 5   Future Work

We are currently finalizing our VAE model analysis where we visualize the latent audio frame space. In addition, we are conducting a qualitative study where composers and practitioners used the LTS framework for sound design applications. We aim to obtain a better understanding of the creative potential of these algorithms by evaluating them in real-world application scenarios within composition practices. Our future work includes a publication where we plan to include the DL model analysis and the qualitative study, as well as a compilation album release.

# Acknowledgements

# References

[1] Akten, M.: Grannma MagNet (2018). URL https://www.memo.tv/works/grannma-magnet/. Library Catalog: www.memo.tv

[2] Briot, J.P., Pachet, F.: Deep learning for music generation: challenges and directions. Neural Computing and Applications **32**(4), 981–993 (2020). DOI 10.1007/s00521-018-3813-6. URL http://link.springer.com/10.1007/s00521-018-3813-6

[3] Dieleman, S.: Sander Dieleman: Generating music in the raw audio domain. URL https://www.youtube.com/watch?v=y8mOZSJA7Bc

[4] Dieleman, S., Oord, A.v.d., Simonyan, K.: The challenge of realistic music generation: modelling raw audio at scale. In: Proceedings of the 32nd Conference on Neural Information Processing Systems (NeurIPS 2018), p. 11. Montreal QC, Canada (2018)

[5] Engel, J., Hantrakul, L.H., Gu, C., Roberts, A.: DDSP: Differentiable Digital Signal Processing (2019). URL https://openreview.net/forum?id=B1x1ma4tDr

[6] Esling, P., Chemla-Romeu-Santos, A., Bitton, A.: Generative timbre spaces: regularizing variational auto-encoders with perceptual metrics. arXiv:1805.08501 [cs, eess] (2018). URL http://arxiv.org/abs/1805.08501. ArXiv: 1805.08501

[7] Gabor, D.: Acoustical Quanta and the Theory of Hearing. Nature **159**(4044), 591–594 (1947). DOI 10.1038/159591a0. URL http://www.nature.com/articles/159591a0

[8] Grey, J.M.: Multidimensional perceptual scaling of musical timbres. The Journal of the Acoustical Society of America **61**(5), 1270–1277 (1977). DOI 10.1121/1.381428. URL http://asa.scitation.org/doi/10.1121/1.381428

[9] Griffin, D.W., Lim, J.S.: Signal estimation from modified short-time Fourier transform. IEEE Transactions on Acoustics, Speech, and Signal Processing **32**(2), 236–243 (1984). DOI 10.1109/TASSP.1984.1164317. URL http://ieeexplore.ieee.org/document/1164317/

[10] Hantrakul, L., Engel, J., Roberts, A., Gu, C.: Fast and Flexible Neural Audio Synthesis. In: Proceedings of the 20th International Society for Music Information Retrieval Conference (ISMIR 2019), p. 7 (2019)

[11] He, K., Zhang, X., Ren, S., Sun, J.: Deep Residual Learning for Image Recognition. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 770–778. IEEE, Las Vegas, NV, USA (2016). DOI 10.1109/CVPR.2016.90. URL http://ieeexplore.ieee.org/document/7780459/

[12] Iverson, P., Krumhansl, C.L.: Isolating the dynamic attributes of musical timbrea. The Journal of the Acoustical Society of America **94**(5), 2595–2603 (1993). Publisher: Acoustical Society of America

[13] Kingma, D.P., Welling, M.: Auto-Encoding Variational Bayes. arXiv:1312.6114 [cs, stat] (2014). URL `http://arxiv.org/abs/1312.6114`. ArXiv: 1312.6114

[14] Kingma, D.P., Welling, M.: An Introduction to Variational Autoencoders. Foundations and Trends in Machine Learning **12**(4), 307–392 (2019). DOI 10.1561/2200000056. URL `http://arxiv.org/abs/1906.02691`. ArXiv: 1906.02691

[15] Krumhansl, C.L.: Why is musical timbre so hard to understand. Structure and perception of electroacoustic sound and music **9**, 43–53 (1989)

[16] Kumar, K., Kumar, R., de Boissiere, T., Gestin, L., Teoh, W.Z., Sotelo, J., de Brebisson, A., Bengio, Y., Courville, A.: MelGAN: Generative Adversarial Networks for Conditional Waveform Synthesis. In: Proceedings of the 33rd Conference on Neural Information Processing Systems (NeurIPS 2019), p. 12. Vancouver, BC, Canada (2019)

[17] Lakatos, S.: A common perceptual space for harmonic and percussive timbres. Perception & psychophysics **62**(7), 1426–1439 (2000). Publisher: Springer

[18] LeCun, Y., Cortes, C., Burges, C.: MNIST handwritten digit database. URL `http://yann.lecun.com/exdb/mnist/`

[19] Luigi, R.: The Art of Noise. A Great Bear Pamphlet (1967)

[20] McAdams, S., Winsberg, S., Donnadieu, S., De Soete, G., Krimphoff, J.: Perceptual scaling of synthesized musical timbres: Common dimensions, specificities, and latent subject classes. Psychological research **58**(3), 177–192 (1995). Publisher: Springer

[21] McFee, B., Raffel, C., Liang, D., Ellis, D.P., McVicar, M., Battenberg, E., Nieto, O.: librosa: Audio and Music Signal Analysis in Python. In: Proceedings of The 14th Python in Science Conference (SCIPY 2015) (2015)

[22] Müller, M.: Fundamentals of Music Processing. Springer International Publishing, Cham (2015). URL `http://link.springer.com/10.1007/978-3-319-21945-5`

[23] Nieto, O., Bello, J.P.: Systematic Exploration Of Computational Music Structure Research. In: Proceedings of the 17th International Society for Music Information Retrieval Conference (ISMIR 2016), p. 7. New York, NY, USA (2016)

[24] Oord, A.v.d., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves, A., Kalchbrenner, N., Senior, A., Kavukcuoglu, K.: Wavenet: A generative model for raw audio. arXiv preprint arXiv:1609.03499 (2016)

[25] Oord, A.v.d., Li, Y., Babuschkin, I., Simonyan, K., Vinyals, O., Kavukcuoglu, K., Driessche, G.v.d., Lockhart, E., Cobo, L.C., Stimberg, F., Casagrande, N., Grewe, D., Noury, S., Dieleman, S., Elsen, E., Kalchbrenner, N., Zen, H., Graves, A., King, H., Walters, T., Belov, D., Hassabis, D.: Parallel WaveNet: Fast High-Fidelity Speech Synthesis. arXiv:1711.10433 [cs] (2017). URL http://arxiv.org/abs/1711.10433. ArXiv: 1711.10433

[26] Perraudin, N., Balazs, P., Sondergaard, P.L.: A fast Griffin-Lim algorithm. In: 2013 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics, pp. 1–4. IEEE, New Paltz, NY (2013). DOI 10.1109/WASPAA. 2013.6701851. URL http://ieeexplore.ieee.org/document/6701851/

[27] Roads, C.: Microsound. The MIT Press, Cambridge, Mass. (2004)

[28] Roads, C.: Composing electronic music: a new aesthetic. Oxford University Press, Oxford (2015)

[29] Schaeffer, P.: Trait des objets musicaux, nouv. edn. Seuil (1964)

[30] Schrkhuber, C., Klapuri, A.: Constant-Q Transform Toolbox For Music Processing. In: Proceedings of the 7th Sound and Music Computing Conference (SMC 2010), p. 8. Barcelona, Spain (2010)

[31] Smalley, D.: Spectromorphology: explaining sound-shapes. Organised Sound **2**(02), 107–126 (1997). DOI 10.1017/S1355771897009059. URL http://journals.cambridge.org/article_S1355771897009059

[32] Stockhausen, K.: Four Criteria of Electronic Music with Examples from Kontakte (1972). URL https://www.youtube.com/watch?v=7xyGtI7KKIY&list=PLRBdTyZ76lvAFOtZvocPjpRVTL6htJzoP

[33] Snderby, C.K., Raiko, T., Maale, L., Snderby, S.K., Winther, O.: How to Train Deep Variational Autoencoders and Probabilistic Ladder Networks. In: Proceedings of the 23rd international conference on Machine learning (ICML 2016). ACM Press, Pittsburgh, Pennsylvania (2016)

[34] Tatar, K., Macret, M., Pasquier, P.: Automatic Synthesizer Preset Generation with PresetGen. Journal of New Music Research **45**(2), 124–144 (2016). DOI 10.1080/09298215.2016.1175481. URL http://dx.doi.org/10.1080/09298215.2016.1175481

[35] Tatar, K., Pasquier, P.: MASOM: A Musical Agent Architecture based on Self Organizing Maps, Affective Computing, and Variable Markov Models. In: Proceedings of the 5th International Workshop on Musical Metacreation (MUME 2017). Atlanta, Georgia, USA (2017)

[36] Tatar, K., Pasquier, P.: Musical agents: A typology and state of the art towards Musical Metacreation. Journal of New Music Research **48**(1), 56–105 (2019). DOI 10.1080/09298215.2018.1511736. URL https://www.tandfonline.com/doi/full/10.1080/09298215.2018.1511736

[37] Tatar, K., Pasquier, P., Siu, R.: Audio-based Musical Artificial Intelligence and Audio-Reactive Visual Agents in Revive. In: Proceedings of the joint International Computer Music Conference and New York City Electroacoustic Music Festival 2019 (ICMC-NYCEMF 2019), p. 8. International Computer Music Association, New York City, NY, USA (2019)

[38] Tuuri, K., Eerola, T.: Formulating a Revised Taxonomy for Modes of Listening. Journal of New Music Research **41**(2), 137–152 (2012). DOI 10.1080/09298215.2011.614951. URL http://www.tandfonline.com/doi/abs/10.1080/09298215.2011.614951

[39] Varese, E., Wen-chung, C.: The liberation of Sound. Perspectives of New Music **5**(1), 11–19 (1966). URL https://www.jstor.org/stable/832385?origin=JSTOR-pdf&seq=1#page_scan_tab_contents

[40] Velasco, G.A., Holighaus, N., Drer, M., Grill, T.: Constructing An Invertible Constant-Q Transform With Nonstationary Gabor Frames. In: Proceedings of the 14th International Conference on Digital Audio Effects (DAFx-11)), p. 7. Paris, France (2011)

[41] Yu, F., Koltun, V.: Multi-scale context aggregation by dilated convolutions. arXiv preprint arXiv:1511.07122 (2015)