



# Audio Engineering Society Conference Paper

Presented at the 2022 International Conference on  
Audio for Virtual and Augmented Reality  
2022 August 15–17, Redmond, WA, USA

*This paper was peer-reviewed as a complete manuscript for presentation at this conference. This paper is available in the AES E-Library (<http://www.aes.org/e-lib>) all rights reserved. Reproduction of this paper, or any portion thereof, is not permitted without direct permission from the Journal of the Audio Engineering Society.*

---

## A Variational Y-Autoencoder for Disentangling Gesture and Material of Interaction Sounds

Simon Schwär<sup>1</sup>, Meinard Müller<sup>1</sup>, and Sebastian J. Schlecht<sup>2</sup>

<sup>1</sup>*International Audio Laboratories, Erlangen, Germany*

<sup>2</sup>*Dept. of Signal Processing and Acoustics and Dept. of Art and Media, Aalto University, Espoo, Finland*

Correspondence should be addressed to Simon Schwär ([simon.schwaer@audiolabs-erlangen.de](mailto:simon.schwaer@audiolabs-erlangen.de))

### ABSTRACT

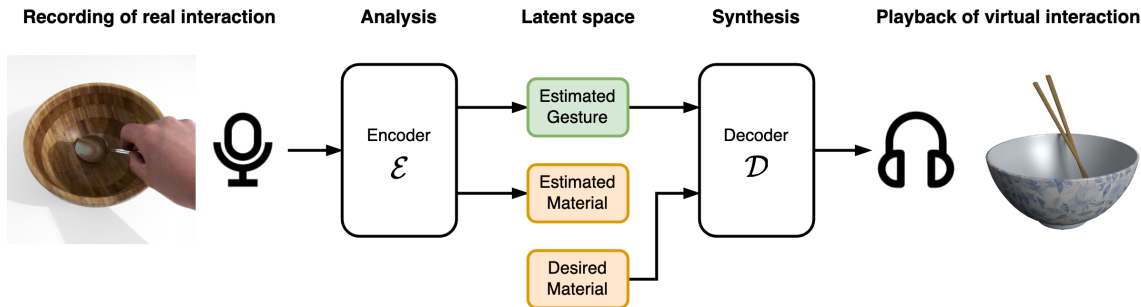
Appropriate sound effects are an important aspect of immersive virtual experiences. Particularly in mixed reality scenarios it may be desirable to change the acoustic properties of a naturally occurring interaction sound (e.g., the sound of a metal spoon scraping a wooden bowl) to a sound matching the characteristics of the corresponding interaction in the virtual environment (e.g., using wooden tools in a porcelain bowl). In this paper, we adapt the concept of a Y-Autoencoder (YAE) to the domain of sound effect analysis and synthesis. The YAE model makes it possible to disentangle the gesture and material properties of sound effects with a weakly supervised training strategy where only an identifier label for the material in each training example is given. We show that such a model makes it possible to resynthesize sound effects after exchanging the material label of an encoded example and obtain perceptually meaningful synthesis results with relatively low computational effort. By introducing a variational regularization for the encoded gesture, as well as an adversarial loss, we can further use the model to generate new and varying sound effects with the material characteristics of the training data, while the analyzed audio signal can originate from interactions with unknown materials.

### 1 Introduction

One goal of audio analysis and synthesis tasks is to identify and allow control of interpretable underlying factors of variation in the analyzed sound before it is resynthesized. A prominent example of this *disentanglement* task is timbre transfer, where some factors of variation (the instrument) are modified, while others (pitch, expression, etc.) should remain constant. In this paper, we combine methods for learning a disentangled representation and adapt them to a particular use case similar to

timbre transfer: modifying the acoustic properties of an object (the *material*) while the interaction itself, i.e., the strength, duration, etc. of the impacts that produced the sound (the *gesture*) remains constant. This can be particularly useful in mixed reality scenarios, where an interaction may be performed with a real object that is replaced by a virtual representation of the object with different characteristics.

As an example, Figure 1 shows how a user may use a metal spoon to interact with a wooden bowl.



**Fig. 1:** Overview of the proposed system: An audio signal is recorded and analyzed by the encoder  $\mathcal{E}$ . The disentangled latent representation consisting of an estimate for gesture and material can be modified (e.g. by changing the material). A signal is resynthesized using the decoder  $\mathcal{D}$  and played back over headphones.

The sound of this real interaction is analyzed, split up into a part that represents the material characteristics and a part that represents the gesture (i.e., how the bowl was struck or scraped). Before resynthesizing, the material representation is replaced to better match the characteristics of the virtual object – for example, chopsticks in a porcelain bowl. The resynthesized sound is then played back to the user, contributing to the immersive experience that the interaction actually took place with the virtual bowl.

To achieve this disentangling analysis and synthesis, we adapt the Y-Autoencoder [1] neural network training strategy to learn a representation of short audio clips where gesture and material are separated and can be individually modified. Using a single class label for each material (e.g., different combinations of spoons and bowls), we can train a model with weakly labeled training data. As opposed to other generative methods that achieve a similar disentanglement, like conditional generative adversarial networks (GANs) [2], an autoencoder (AE) can learn to analyze the disentangled representation as well. By adding a variational constraint on the learned gesture representation, we show that the gesture can also be encoded from recordings of interactions with objects that have not been part of the training set of materials. Furthermore, we add an adversarial loss to improve the resynthesis after changing the material class of an encoded example. With this, we combine advantages of generative methods and a conditional autoencoder in a single model, while achieving

reasonable audio quality with small models whose computational complexity allows real-time execution.

This paper is organized as follows. In Section 2, we introduce our notation, review the concept of the Y-Autoencoder and discuss related work in the field of sound effect synthesis with neural networks. Section 3 gives an overview of our model and training strategy. In Section 4, we introduce our datasets and the experimental setups. Results are presented and discussed in Section 5. We provide audio examples and code on a supplemental website<sup>1</sup>.

## 2 Background

An autoencoder (AE) consists of an encoder  $\mathcal{E} : \mathcal{X} \rightarrow \mathbb{R}^L$  and a decoder  $\mathcal{D} : \mathbb{R}^L \rightarrow \mathcal{X}$ , so that for any  $x \in \mathcal{X}$ ,  $x \approx \mathcal{D}(\mathcal{E}(x))$ , where  $\mathcal{X}$  is the set of possible inputs (e.g., images or audio frames). The vector  $z = \mathcal{E}(x)$  is called an  $L$ -dimensional *latent representation* of  $x$ , where  $L$  is typically much smaller than the dimensionality of  $\mathcal{X}$ . The encoding and decoding can be learned when  $\mathcal{E}$  and  $\mathcal{D}$  are neural networks and the reconstruction loss function

$$\mathcal{L}_R(x) = \|x - \mathcal{D}(\mathcal{E}(x))\|_2^2 \quad (1)$$

is minimized on a training set  $\mathcal{X}_T \subset \mathcal{X}$ .

<sup>1</sup><https://audiolabs-erlangen.de/resources/2022-AVAR-InteractionSounds>

A variant of this method is the *variational* autoencoder (VAE) [3], where the encoder outputs parameters  $\mu$  and  $\sigma$  of a multivariate normal distribution with independent dimensions from which  $z$  is sampled. The objective is extended to learn a continuous and compact latent space by regularizing the encoder output distribution for any input to be close to  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ , with  $\mathbf{0}$  as the all-zero vector and  $\mathbf{I}$  as the identity matrix. This is achieved by minimizing the Kullback-Leibler (KL) divergence  $\text{KL}(\cdot \parallel \cdot)$  between the distributions, yielding a loss function with two terms:

$$\begin{aligned} \mathcal{L}_{\text{VAE}} &= \mathcal{L}_{\text{R}} + \text{KL}(\mathcal{N}(\mu, \text{diag}(\sigma)) \parallel \mathcal{N}(\mathbf{0}, \mathbf{I})) \\ &= \mathcal{L}_{\text{R}} + \frac{1}{2}(1 + \log(\sigma^2) - \sigma^2 - \mu^2). \end{aligned} \quad (2)$$

With the constraint on the latent space organization, the VAE learns a representation where the distance between two encodings reflects the similarity of the original data points.

Adversarial losses have been used successfully to improve audio quality of generative models [4] and autoencoders [5]. The loss involves a separate *discriminator* neural network  $\mathcal{A} : \mathcal{X} \rightarrow \{0, 1\}$  that is trained to discriminate a real input  $x \in \mathcal{X}_{\text{r}}$  from an output  $\hat{x} = \mathcal{E}(\mathcal{D}(x))$ . The gradient of the loss to maximize the chance of correctly identifying  $x$  and  $\hat{x}$  is used to train  $\mathcal{A}$ , whereas the gradient of the loss to maximize the chance of misclassifying  $\hat{x}$  as real is used to train  $\mathcal{E}$  and  $\mathcal{D}$ .

## 2.1 Y-Autoencoder

The Y-Autoencoder (YAE) [1] introduces a partly supervised method to disentangle a known class label  $y \in \{0, 1\}^C$  (a *one-hot* encoding with  $C$  classes) from the rest of the latent representation. We denote the class estimate as the *explicit part*  $z_e \in [0, 1]^C$  and the remaining latent representation as the implicit part  $z_i \in \mathbb{R}^L$ . The encoder and decoder of the YAE process  $z_e$  and  $z_i$  separately, so that  $\mathcal{E}_{\text{Y}} : \mathcal{X} \rightarrow \mathbb{R}^C \times \mathbb{R}^L$  with  $(z_e, z_i) = \mathcal{E}_{\text{Y}}(x)$ , and  $\mathcal{D}_{\text{Y}} : \mathbb{R}^C \times \mathbb{R}^L \rightarrow \mathcal{X}$  with  $x \approx \mathcal{D}_{\text{Y}}(z_e, z_i)$ .

To learn the disentanglement, i.e., the independence of  $z_e$  and  $z_i$ , the latent representation is decoded twice with  $\hat{x} = \mathcal{D}_{\text{Y}}(y, z_i)$  and  $\tilde{x} = \mathcal{D}_{\text{Y}}(y_{\text{r}}, z_i)$ , where  $y_{\text{r}} \in \mathbb{R}^C$  is the one-hot label of a random class. The YAE loss

$$\mathcal{L}_{\text{YAE}} = \mathcal{L}_{\text{R}} + \mathcal{L}_{\text{L}} + \mathcal{L}_{\text{E}} + \mathcal{L}_{\text{I}} \quad (3)$$

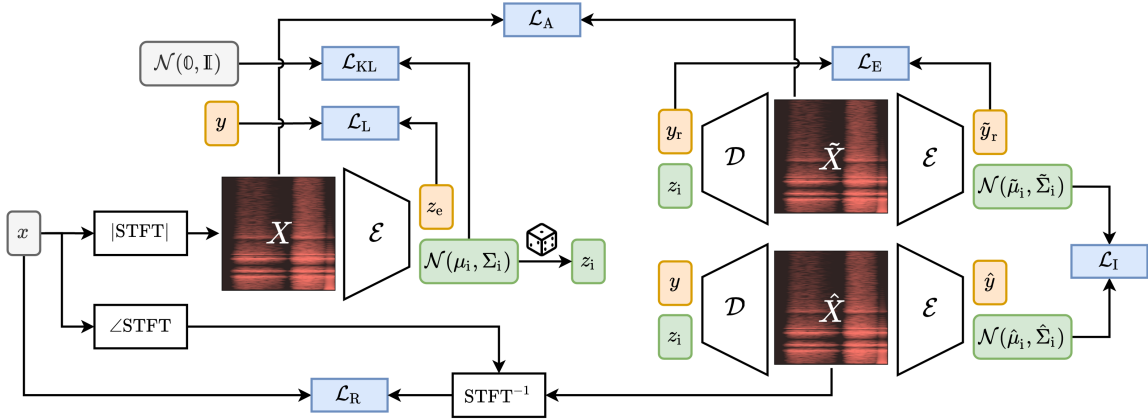
consists of four terms (nomenclature adapted by us):

1. A reconstruction loss  $\mathcal{L}_{\text{R}}$  as in (1).
2. The label classification loss  $\mathcal{L}_{\text{L}} = \text{CE}(z_e, y)$ , where the cross-entropy loss  $\text{CE}(a, b) = \sum_{j=1}^C a_j \log b_j$  between  $z_e$  and  $y$  aids the supervised classification of the inputs while enforcing that no other information is encoded in  $z_e$ .
3. The random explicit loss  $\mathcal{L}_{\text{E}} = \text{CE}(\tilde{y}_{\text{r}}, y_{\text{r}})$ : A subsequent re-encoding  $(\tilde{y}_{\text{r}}, \tilde{z}_i) = \mathcal{E}(\tilde{x})$  should correctly predict the class  $y_{\text{r}}$ .  $\mathcal{L}_{\text{E}}$  forces the decoder to make use of the information in the explicit part, as otherwise the encoder would not be able to infer this information from  $\tilde{x}$ .
4. The implicit consistency loss  $\mathcal{L}_{\text{I}} = \|\tilde{z}_i - \hat{z}_i\|_2^2$ : Re-encoding  $\hat{x}$  with  $(\hat{y}_{\text{r}}, \hat{z}_i) = \mathcal{E}_{\text{Y}}(\hat{x})$  should yield the same implicit part as re-encoding  $\tilde{x}$ . This ensures that the information in  $z_i$  is used consistently by the decoder.

## 2.2 Related Work

Interactive media like augmented and virtual reality require to adapt sound effects to the actual interaction at run-time, which makes the use of traditional recorded sound effects (“Foley”) more difficult. Physical [6] and signal [7, 8] models can be parametrized to produce appropriate sounds in real-time. However, many models are limited to a specific sound effect type, and the parametrization often requires a manual mapping from an interaction to the model configuration. Furthermore, computational cost may prohibit a very detailed modelling of the effects.

Several works have addressed the use of deep learning for sound effects. In a multimodal approach, fitting sound effects can be selected from a library by analyzing a silent video sequence of interactions between a drum stick and the environment [9]. A similar method can be used to select more general Foley sound effects [10] or even MIDI events that represent music a person may be playing in the analyzed video [11]. Direct synthesis of a wide variety of sounds can be achieved with GANs [12], which also allow to generate varying footsteps specifically [13]. However, the models allow limited control over the produced sound effects.



**Fig. 2:** Visualization of the encoding/decoding method and different losses with their respective inputs at training time. Blue color for losses, grey for inputs, orange for material class labels, green for gesture representations. Further,  $\Sigma_i = \text{diag}(\sigma_i)$ , also analogously for  $\tilde{\sigma}_i$  and  $\hat{\sigma}_i$ .

### 3 The VYAE Model

To learn the analysis and synthesis using a disentangled latent representation, we set up the *Variational Y-Autoencoder* (VYAE) model. While the training strategy roughly follows the YAE, we make some important modifications to the original YAE that we highlight in this section.

Our model input  $x$  consists of 8192 time-domain audio samples with a sampling rate of 16 kHz. The input of the encoder and output of the decoder is a log-magnitude spectrogram denoted by  $X$  and  $\hat{X}$ , respectively. The spectrogram  $X = \log |\mathcal{S}_{N,H}(x) + \epsilon|$  is calculated<sup>2</sup> using the short-time Fourier transform  $\mathcal{S}_{N,H}$  with window size  $N = 1024$  and hop size  $H = 256$ . In the following, we omit the subscripts  $N$  and  $H$  when this window and hop size are used. By discarding the first Fourier coefficient (the DC component of the signal) of the time-frequency representation, we yield  $512 \times 32$  values for  $X$ .

The encoder outputs are an estimate of the material class  $z_e \in [0, 1]^C$  for  $C$  classes, as well as parameters  $\mu_i, \sigma_i$  of a normal distribution representing the estimated gesture  $z_i \sim \mathcal{N}(\mu_i, \text{diag}(\sigma_i))$ . We use variational regularization for the gesture

<sup>2</sup>The small constant  $\epsilon$  is used to avoid taking a logarithm of 0. We omit it hereinafter.

space to facilitate a continuous and compact representation of different gestures. This allows interpolation between gestures, as well as sampling new and similar gestures from the latent space while preserving a perceptually meaningful output of the decoder. Furthermore, even though this is not generally the case [14], the variational regularization can improve generalization (cf. Fig. 6).

We modify the YAE loss function (3) to accommodate the variational regularization and the inputs/outputs in audio domain. This way, the VYAE loss  $\mathcal{L}_{\text{VYAE}}$  becomes

$$\mathcal{L}_{\text{VYAE}} = \mathcal{L}_R + \mathcal{L}_L + \mathcal{L}_E + \mathcal{L}_I + \mathcal{L}_{\text{KL}} + \mathcal{L}_A. \quad (4)$$

Fig. 2 gives an overview of the training procedure and the inputs of the loss terms. The individual terms are defined as follows:

1.  $\mathcal{L}_R$ : We use a multi-scale spectral (MSS) distance [15], where time-domain audio is transformed to multiple log-magnitude spectrograms with varying hop and window sizes and the  $L^1$  distance between the spectrograms is accumulated. We first have to reconstruct time-domain audio from the decoder output to compare it with the input signal frame  $x$ . To do this, we use the phase of the input audio's time-frequency representation and combine it with the spectrogram output of

the decoder to yield

$$\hat{x} = \mathcal{S}^{-1}\left(\hat{X} \cdot \exp(i\angle(\mathcal{S}(x)))\right), \quad (5)$$

where  $i$  denotes the imaginary unit and  $\angle(\cdot)$  is the angle of the complex Fourier coefficients. Note that only the reconstruction with the correct class label is compared in  $\mathcal{L}_R$ . This way, at training time, the phase of  $x$  is only used for a signal that contains a reconstruction with the original material, which should ultimately be very similar to  $x$  anyway.

The MSS distance is then given by

$$\mathcal{L}_R = \sum_{(N,H) \in \mathcal{M}} \left\| \log |\mathcal{S}_{N,H}(x)| - \log |\mathcal{S}_{N,H}(\hat{x})| \right\|_1, \quad (6)$$

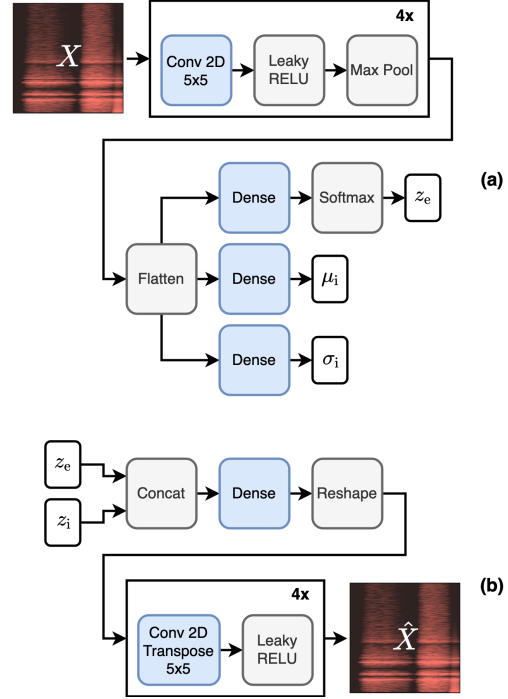
where  $\mathcal{M}$  is a set of window and hop sizes ( $\mathcal{M} = \{(2048, 1024), (1024, 256), (512, 128)\}$  in our experiments).

2.  $\mathcal{L}_L$  and  $\mathcal{L}_E$ : We use cross-entropy between  $y$  and  $z_e$  and between  $y_r$  and  $\tilde{y}_r$ , respectively, as in the original YAE.
3.  $\mathcal{L}_I$ : The re-encoding of  $\tilde{X}$  and  $\hat{X}$  yields parameters  $\tilde{\mu}_i$ ,  $\tilde{\sigma}_i$  and  $\hat{\mu}_i$ ,  $\hat{\sigma}_i$  to represent the estimated gesture (cf. Fig. 2). To fulfil the purpose of  $\mathcal{L}_I$  in the original YAE, the distributions  $\mathcal{N}(\tilde{\mu}_i, \text{diag}(\tilde{\sigma}_i))$  and  $\mathcal{N}(\hat{\mu}_i, \text{diag}(\hat{\sigma}_i))$  should be as similar as possible. We quantify the similarity using the Bhattacharyya distance [16]. For two multivariate normal distributions  $\mathcal{N}(\mu_1, \Sigma_1)$  and  $\mathcal{N}(\mu_2, \Sigma_2)$ , the loss term is given by

$$\mathcal{L}_I = \frac{1}{8}(\mu_1 - \mu_2)^T \Sigma^{-1}(\mu_1 - \mu_2) + \frac{1}{2} \log \left( \frac{\det \Sigma}{\sqrt{\det \Sigma_1 \det \Sigma_2}} \right) \quad (7)$$

where  $\Sigma = (\Sigma_1 + \Sigma_2)/2$ . As the dimensions of the VAE are assumed to be independent,  $\Sigma_1$  and  $\Sigma_2$  are diagonal and their inverse and determinant can be calculated efficiently.

4.  $\mathcal{L}_{KL}$ : We use the KL divergence  $\text{KL}(\mathcal{N}(\mu_i, \text{diag}(\sigma_i)) \parallel \mathcal{N}(\mathbf{0}, \mathbf{I}))$  as in (2). It is sufficient to regularize the first encoding, as this ensures a loss for the encoder that favors the continuous and compact representation.



**Fig. 3:** Encoder (a) and decoder (b) architecture. Blue color for building blocks with learnable weights.

5.  $\mathcal{L}_A$ : An adversarial loss between  $X$  and  $\tilde{X}$  (the decoder output generated with a random material) improves the decoding of gestures where the material has changed. To achieve this, we use a discriminator  $\mathcal{A}_Y : \mathcal{X} \times [0, 1]^C \rightarrow \{0, 1\}$  that in addition to  $X$  and  $\tilde{X}$  receives the corresponding material class label as an input. This way,  $\mathcal{A}_Y$  learns to discriminate material-specific characteristics of the spectrogram.

The loss for  $\mathcal{A}_Y$  is separate from  $\mathcal{L}_{VYAE}$  and given by

$$\mathcal{L}_{\mathcal{A}_Y} = \text{CE}(\mathbf{1}, \mathcal{A}_Y(X, y)) + \text{CE}(\mathbf{0}, \mathcal{A}_Y(\tilde{X}, y_r)), \quad (8)$$

where  $\mathbf{1}$  and  $\mathbf{0}$  are the vectors of all ones and zeros of appropriate dimension. The loss term for  $\mathcal{E}$  and  $\mathcal{D}$  is

$$\mathcal{L}_A = \text{CE}(\mathbf{1}, \mathcal{A}_Y(\tilde{X}, y_r)). \quad (9)$$

Figure 3 shows the neural network architectures used for encoder and decoder. The encoder consists of a 2D convolutional neural network (CNN) with

four layers with an increasing number of kernels (16, 32, 48, 64) and a  $5 \times 5$  kernel size. Each layer is followed by a *LeakyRELU* [17] non-linearity and a max pooling layer to downsample the layer output on the time- and frequency axes by selecting the maximum value in each  $1 \times 4$  block in the first layer and  $2 \times 2$  block in each subsequent layer. The last convolutional layer outputs 4096 values that are flattened to one dimension and form the input of three fully-connected (dense) layers which output  $z_e, \mu_i \in \mathbb{R}^L$  and  $\sigma_i \in \mathbb{R}^L$  with  $L = 32$ , respectively. As  $z_e$  represents a probability distribution over the material classes, the softmax nonlinearity ensures that the values sum to 1. The discriminator  $\mathcal{A}_Y$  uses the same architecture as the encoder, where the additional class label input is concatenated to the flattened convolutional layer output and the final dense layer outputs a single value between 0 and 1.

The decoder concatenates the inputs  $z_e$  and  $z_i$ , resulting in  $L + C$  values that form the input of a dense layer with 4096 outputs, which are reshaped to match the output dimensions of the encoding convolutional layers. Four subsequent upsampling convolutional layers invert the dimensionality reduction of the max pooling to yield the output log-magnitude spectrogram  $\hat{X}$ .

## 4 Experiments

### 4.1 Datasets

We compile two datasets to evaluate the proposed approach of disentangling material and gesture in recordings of real interaction sounds.

The *Spoon/Bowl* (SB) dataset is a collection of two-minute recordings of interactions using spoons in a bowl (scraping, tapping, etc.). In each recording, a different combination of plastic/metal/wooden spoon and plastic/metal/wooden/glass/porcelain bowl is used to perform random interactions, yielding  $C = 15$  classes and 30 minutes of audio material in total. All interactions were recorded in a silent environment with a close-up large-diaphragm condenser microphone. We generate a total of 95000 random excerpts of length 8192 samples (0.512 seconds at 16 kHz sampling rate) from the audio files to form a training dataset. As a test set, we

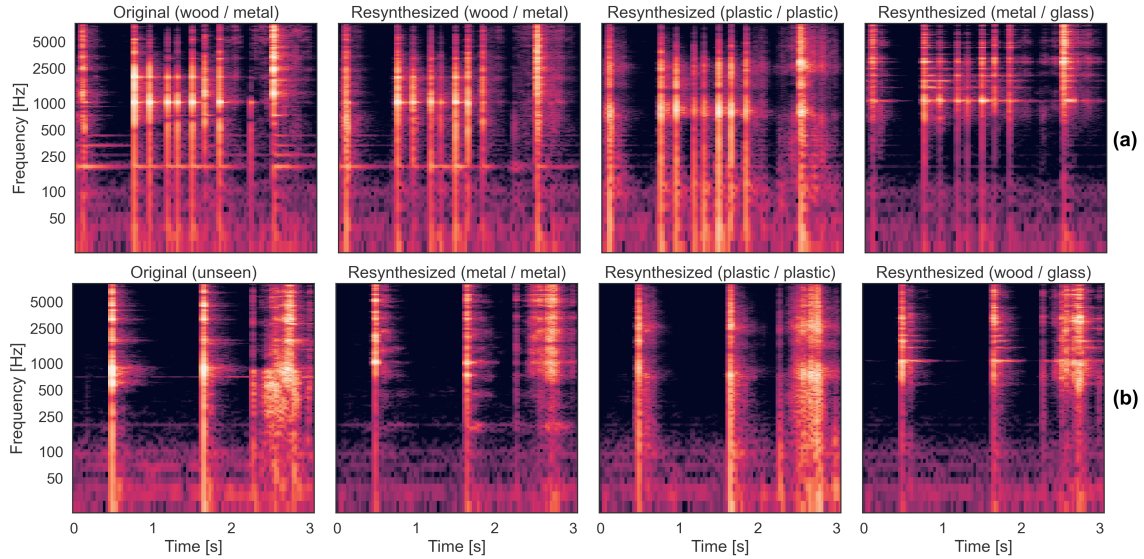
recorded additional audio clips of 20 seconds duration for each class that are not part of the training set. The test set additionally contains interactions with different spoons and bowls recorded in varying acoustic conditions that do not occur in the training set at all. We use them as an “unseen” class to evaluate the generalization of gesture estimations. We denote the classes by “[spoon material] / [bowl material]” and “unseen” in the following.

The *Freesound Footsteps* (FF) dataset consists of 55 recordings of footsteps obtained from the Freesound [18] database with a total of around 17 minutes of audio material. The recordings were selected to contain a variety of footstep sounds, including dry, wet, indoor and outdoor surfaces, with a total of  $C = 55$  classes. Background noise, reverb and other degradations are present on some recordings, and all were converted to a single channel signal with 16 kHz sampling rate. As opposed to the SB data, the amount of data for each class varies, with recording lengths between 2 and 118 seconds. We generate 10 random excerpts of the same length as the SB excerpts for each second of audio to form a training set, resulting in a total of 10018 excerpts (between 22 and 1187 per class). We reserve longer clips from two classes (“high heels” and “carpet”) that are not part of the training set for evaluation purposes.

### 4.2 VYAE Model

We train the full model as described in Section 3 separately with the SB and FF datasets for 6.5M and 0.7M steps (70 epochs), respectively, using a batch size of 8 and a learning rate of  $10^{-4}$ . The training can be completed in less than 60 minutes on a single Nvidia Geforce RTX 2080 Ti GPU using Tensorflow. We multiply  $\mathcal{L}_R$  by a constant factor of  $\frac{1}{300}$  to facilitate a comparable magnitude of all loss terms.

To evaluate the disentanglement in a real-time setting, we use the test sets of SB and FF and process the audio signals frame-wise with the corresponding trained VYAE model. By using an overlap of 128 samples between frames, we can cross-fade between successive frames to avoid potential discontinuities in the reconstructed audio and are able to process a continuous audio stream with our model. Successful disentanglement should allow to



**Fig. 4:** Results with Model I trained on the Spoon/Bowl dataset. (a) Original class “metal / metal”. (b) Original class “unseen” that is not part of the training dataset.

exchange the material while preserving the gesture in the output signal. We discard the estimated class  $z_e$  of each frame and manually specify the material class  $y$  in the input of the decoder.

### 4.3 Variational regularization and adversarial loss

To investigate the effectiveness of variational regularization and adversarial loss, we train three different versions of the VYAE model (see Table 1) with the same parameters as above. Model I is the full model using the  $\mathcal{L}_{\text{VYAE}}$  as in (4), including  $\mathcal{L}_A$  and  $\mathcal{L}_{\text{KL}}$ . Models II and III are trained without adversarial loss  $\mathcal{L}_A$ . The variational regularization  $\mathcal{L}_{\text{KL}}$  is additionally removed in Model III, where the encoder outputs the  $L$ -dimensional gesture representation directly.

Model	I	II	III
Variational reg. $\mathcal{L}_{\text{KL}}$	Yes	Yes	No
Adversarial loss $\mathcal{L}_A$	Yes	No	No

**Table 1:** Model configurations to investigate the effects of the loss terms  $\mathcal{L}_{\text{KL}}$  and  $\mathcal{L}_A$ .

To compare the different model outputs to a real recorded example from the target class, we compile

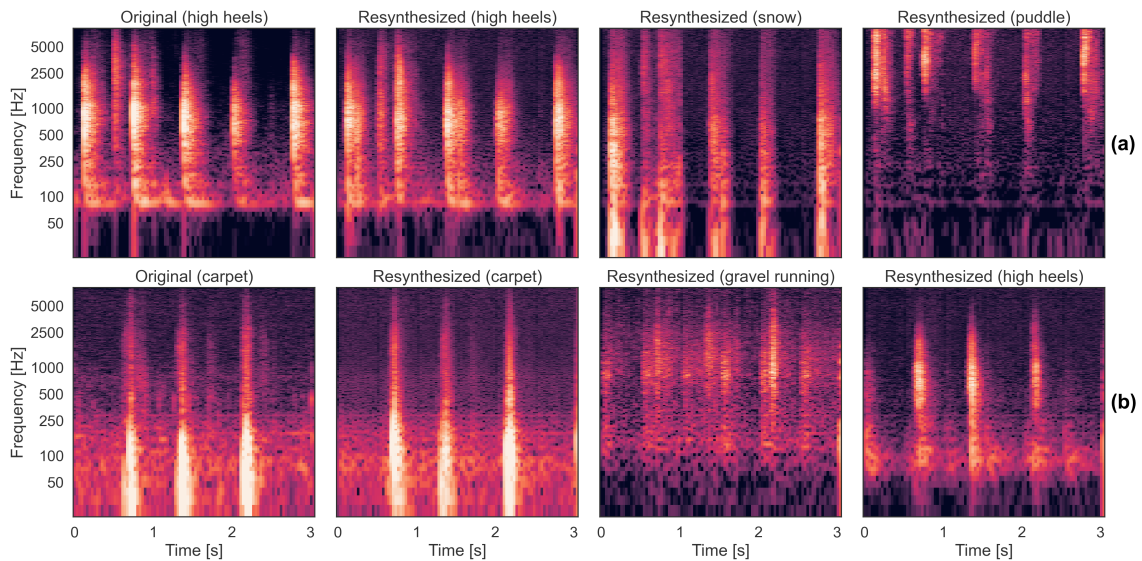
a set of single short impact interactions from the SB test set and align them in time to allow a direct qualitative comparison.

## 5 Results and Discussion

Figure 4 shows excerpts from the longer sequences of the SB test set. We encourage the reader to listen to the examples on the accompanying website<sup>3</sup>. The resynthesis of the examples with different material classes demonstrates that the material characteristics (visually best represented by the strong modes at different frequencies) can be changed independent of the gesture. This includes the “unseen” class in Fig. 4 (b), which shows that the learned gesture representation generalizes to encode the actual underlying factors of variation in the spoon/bowl interaction. However, artifacts of improper disentanglement are sometimes visible and audible. As an example, Fig. 4 (a) exhibits a peak below 250 Hz in the spectrum of the “wood / metal” recording that is preserved to some degree in the resynthesized versions with the “plastic / plastic” material class.

In total, we analyze and resynthesize 1785 frames for the sequential processing of the SB test set

<sup>3</sup><https://audiolabs-erlangen.de/resources/2022-AVAR-InteractionSounds>



**Fig. 5:** Results with Model I trained on the Freesound Footsteps dataset. (a) Original class “high heels”. (b) Original class “carpet”.

excluding the “unseen” class. The computational complexity of the model easily allows the real-time execution on a recent computer. With an Apple M1 Max CPU, we can analyze and resynthesize a frame of 0.512 seconds at 16 kHz sampling rate in 0.027 seconds on average. Furthermore, for 96% of the frames, the material of the input signal is correctly classified, while most misclassified frames contain very little activity.

Results with the FF dataset are shown in Figure 5. With the FF test set, we can find some combinations where the synthesis with changed material fails. In Fig. 5(b), the slow steps in the “carpet” class are resynthesized with the “gravel running” class label. As the “gravel running” class recording only contains faster steps, the model produces fast steps from the “carpet” encoded gesture as well. This shows how the learned disentanglement also depends on a balanced training set that contains all expected kinds of gestures for all classes.

Another limitation of our method manifests in the conversion from “carpet” to “high heels”. As we use the phase from the input signal for the output, the more transient “high heels” steps are not accurately reproduced, even though the magnitude spectrogram is similar to an original spectrogram from the class (cf. Fig. 5). As the conversion from

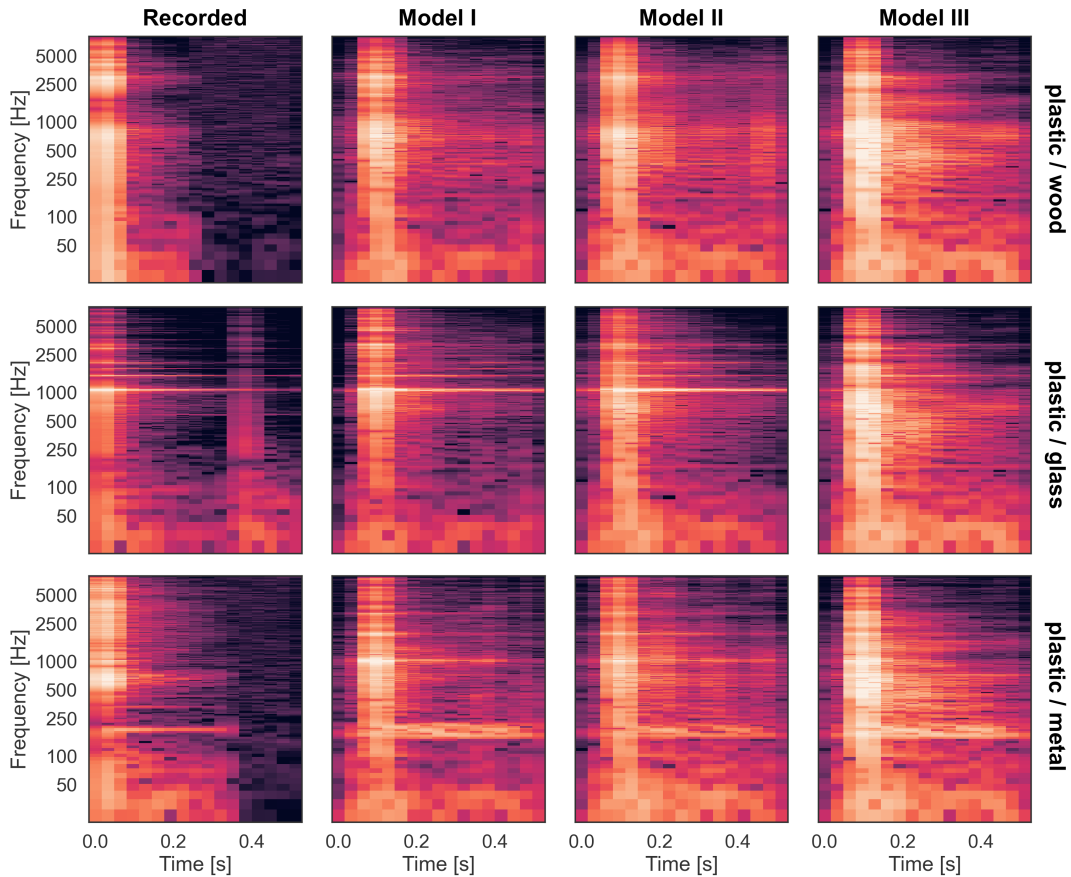
magnitude spectrogram to time domain at runtime is independent of the training method, a different phase reconstruction method, for example based on deep learning [4], could further improve our results.

Figure 6 shows the effect of variational regularization and adversarial loss. We take a frame with a single impact interaction from the “unseen” class recording, estimate the gesture with the encoder, and resynthesize the sound with the classes “plastic / wood”, “plastic / glass”, and “plastic / metal”. The first column shows a recorded example of the respective class for comparison, while the other three columns show spectrograms of the resynthesized signal. Model III (without  $\mathcal{L}_{KL}$ , column 4) clearly fails to generalize to the unseen class input and the resynthesis with different classes exhibits only small variations in the output signal. Model II (without  $\mathcal{L}_A$ , column 3) produces some artifacts related to the more reverberant recording conditions of the “unseen” class example, whereas Model I performs best in this comparison.

## 6 Conclusions

In this paper, we showed how a Variational Y-Autoencoder makes it possible to learn disentan-





**Fig. 6:** Results with different models trained on the Spoon/Bowl dataset. The first column shows a recorded example from the material classes “plastic / wood”, “plastic / glass”, and “plastic / metal”. Columns 2-4 show a resynthesis using these classes and a gesture estimated from an “unseen” example.

gling gesture and material in recordings of interaction sound effects. We improved the generalization and resynthesis quality of the model by adding variational regularization and an adversarial loss to the training procedure. The model allows for resynthesis of the sound effect with an exchanged material, even when the encoder input does not belong to a class known to the model. To our knowledge, this is the first application of a timbre transfer method to interaction sound effects. While more complex scenarios, e.g. including nonisotropic materials, may require a multi-modal approach for convincing resynthesis, this pure *audio-to-audio* transformation may be an interesting tool for sound design in mixed reality applications. In future work, we would like to explore disentanglement methods that facilitate a learnable representation of the material

as well. This way, authoring of the desired material characteristics for a sound effect could be done by encoding recordings of a few example interactions and storing the estimated material representation for later resynthesis.

## Acknowledgements

During this project, Simon Schwär was supported by a fellowship within the IFI programme of the German Academic Exchange Service (DAAD). The International Audio Laboratories Erlangen are a joint institution of the Friedrich-Alexander-Universität Erlangen-Nürnberg (FAU) and Fraunhofer Institute for Integrated Circuits IIS.

## References

- [1] Patacchiola, M., Fox-Roberts, P., and Rosten, E., “Y-Autoencoders: Disentangling latent representations via sequential encoding,” *Pattern Recognition Letters*, 140, pp. 59–65, 2020, doi:10.1016/j.patrec.2020.09.025.
- [2] Mirza, M. and Osindero, S., “Conditional Generative Adversarial Nets,” *ArXiv preprint*, abs/1411.1784, 2014.
- [3] Kingma, D. and Welling, M., “Auto-Encoding Variational Bayes,” in *2nd International Conference on Learning Representations (ICLR)*, Banff, AB, Canada, 2014.
- [4] Kumar, K., Kumar, R., de Boissiere, T., Gestin, L., Teoh, W. Z., Sotelo, J., de Brebisson, A., Bengio, Y., and Courville, A., “MelGAN: Generative Adversarial Networks for Conditional Waveform Synthesis,” in *Proceedings of the 33rd Conference on Neural Information Processing Systems (NeurIPS)*, 2019.
- [5] Caillon, A. and Esling, P., “RAVE: A variational autoencoder for fast and high-quality neural audio synthesis,” *ArXiv preprint*, abs/2111.05011, 2021.
- [6] Bilbao, S., *Numerical Sound Synthesis: Finite Difference Schemes and Simulation in Musical Acoustics*, John Wiley & Sons, Ltd., USA, 2009.
- [7] Serra, X. and Smith, J., “Spectral Modeling Synthesis: A Sound Analysis/Synthesis System Based on a Deterministic Plus Stochastic Decomposition,” *Computer Music Journal*, 14(4), pp. 12–24, 1990.
- [8] Cook, P. R., *Real Sound Synthesis for Interactive Applications*, A. K. Peters, Ltd., USA, 2002, ISBN 1568811683.
- [9] Owens, A., Isola, P., McDermott, J. H., Torralba, A., Adelson, E., and Freeman, W., “Visually Indicated Sounds,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2405–2413, 2016, doi:10.1109/CVPR.2016.264.
- [10] Ghose, S. and Prevost, J. J., “AutoFoley: Artificial Synthesis of Synchronized Sound Tracks for Silent Videos With Deep Learning,” *IEEE Transactions on Multimedia*, 23, pp. 1895–1907, 2021, doi:10.1109/TMM.2020.3005033.
- [11] Gan, C., Huang, D., Chen, P., Tenenbaum, J. B., and Torralba, A., “Foley Music: Learning to Generate Music from Videos,” in *16th European Conference on Computer Vision (ECCV)*, online, 2020.
- [12] Donahue, C., McAuley, J., and Puckette, M., “Adversarial Audio Synthesis,” in *7th International Conference on Learning Representations (ICLR)*, 2019.
- [13] Comunità, M., Phan, H., and Reiss, J. D., “Neural Synthesis of Footsteps Sound Effects with Generative Adversarial Networks,” *ArXiv preprint*, abs/2110.09605, 2021.
- [14] Bozkurt, A., Esmaeili, B., Tristan, J.-B., Brooks, D., Dy, J., and van de Meent, J.-W., “Rate-Regularization and Generalization in Variational Autoencoders,” in *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, pp. 3880–3888, 2021.
- [15] Engel, J., Hantrakul, L. H., Gu, C., and Roberts, A., “DDSP: Differentiable Digital Signal Processing,” in *8th International Conference on Learning Representations (ICLR)*, 2020.
- [16] Bhattacharyya, A. K., “On a Measure of Divergence between Two Statistical Populations Defined by Their Probability Distributions,” *Bulletin of the Calcutta Mathematical Society*, (35), pp. 99–109, 1943.
- [17] Maas, A. L., Hannun, A. Y., and Ng, A. Y., “Rectifier Nonlinearities Improve Neural Network Acoustic Models,” in *Proceedings of the 30th International Conference on Machine Learning*, Atlanta, Georgia, USA, 2013.
- [18] Font, F., Roma, G., and Serra, X., “Freesound Technical Demo,” in *ACM International Conference on Multimedia (MM’13)*, pp. 411–412, ACM, Barcelona, Spain, 2013.