# AUTOMATIC MULTITRACK MIXING WITH A DIFFERENTIABLE MIXING CONSOLE OF NEURAL AUDIO EFFECTS

*Christian J. Steinmetz* [1,2,*]    *Jordi Pons* [1]    *Santiago Pascual* [1]    *Joan Serrà* [1]

[1] Dolby Laboratories
[2] Music Technology Group, Universitat Pompeu Fabra

## ABSTRACT

Applications of deep learning to automatic multitrack mixing are largely unexplored. This is partly due to the limited available data, coupled with the fact that such data is relatively unstructured and variable. To address these challenges, we propose a domain-inspired model with a strong inductive bias for the mixing task. We achieve this with the application of pre-trained sub-networks and weight sharing, as well as with a sum/difference stereo loss function. The proposed model can be trained with a limited number of examples, is permutation invariant with respect to the input ordering, and places no limit on the number of input sources. Furthermore, it produces human-readable mixing parameters, allowing users to manually adjust or refine the generated mix. Results from a perceptual evaluation involving audio engineers indicate that our approach generates mixes that outperform baseline approaches. To the best of our knowledge, this work demonstrates the first approach in learning multitrack mixing conventions from real-world data at the waveform level, without knowledge of the underlying mixing parameters.

***Index Terms***— Intelligent music production, automatic mixing, deep learning, temporal convolutional network.

## 1. INTRODUCTION

In the post-production process, the audio engineer is tasked with creating a cohesive mixture of the recorded elements. This process involves a number of technical challenges [1], such as reducing masking, ensuring balance between the sources, and addressing noise or bleed, as well as artistic considerations, such as selecting the timbre and level of the artificial reverberation. Producing a mix is especially challenging due to the interplay between the aforementioned tasks, which are often dependant on each other, and not easily separated.

Over the past decade, the accessibility of recording and music production equipment has rapidly increased. This, along with the growth and accessibility of digital distribution platforms, has brought music production to a new, diverse demographic [2]. Nevertheless, while these tools have become affordable and readily available, the skill and expertise required for their operation has remained relatively constant. Intelligent music production (IMP) is a research field focused on the development of algorithms that provide feedback, assistance, or automation in the process of recording, mixing, or mastering music [3]. These methods often aim to address the high level of skill required in the music production process, lowering the barrier of entry, but their applications do not end there. Work in IMP may also help expedite the workflow of professional engineers, potentially uncover new understanding about current mixing conventions, or even discover new techniques for multitrack mixing.

---

* Work done during an internship at Dolby Laboratories.

IMP systems generally implement either rule-based or classical machine learning approaches [4]. Rule-based approaches rely upon establishing a set of rules and logic surrounding best practices [5–7]. While they generate convincing results for some cases, they do not provide a level of expressivity that matches human audio engineers [8]. In comparison, classical machine learning approaches allow for greater model flexibility, but have typically suffered from the lack of parametric mixing data (i.e., the exact settings of each processor in the mix). For this reason, they have been of low-complexity, limiting their practical application [9–11]. While both approaches have seen some success in addressing particular aspects of the mixing process, they ultimately have failed to capture the entire process and generalize at the scale of real-world projects.

The previous shortcomings, along with the promise of deep learning methods in multiple audio signal processing tasks, motivate the application of those within IMP. Nevertheless, there are a number of unique challenges in the application of deep learning methodologies to automatic mixing that have yet to be addressed:

1. Large variation in the types and number of sources.
2. Expectation for high-fidelity, requiring a low tolerance for artifacts and high sampling rates (at least 44.1 kHz).
3. Ability for audio engineers to view and adjust the resulting mix parameters in an intuitive manner.
4. Presence of many acceptable mixes and the difficulty in the objective evaluation of their quality.

In this paper, we address these challenges and demonstrate how we can successfully apply these methods. Our major contributions are:

- We demonstrate that temporal convolutional networks can model a series connection of signal processing devices, across their parameter spaces.
- We propose a differentiable mixing console enabling interpretability and the ability to learn from limited and unstructured data.
- We introduce a loss function based on sum and difference signals, critical in enabling learning from real-world mixes.
- In a perceptual evaluation, we demonstrate that our model can learn mixing conventions from raw audio waveforms of real-world mixes, which we believe to be the first of its kind.

## 2. DIFFERENTIABLE MIXING CONSOLE

We aim to achieve a strong inductive bias for the mixing task by incorporating knowledge from the signal processing chain of a traditional mixing console. We consider a neural network that analyzes a set of audio inputs, and then predicts a set of parameters for each
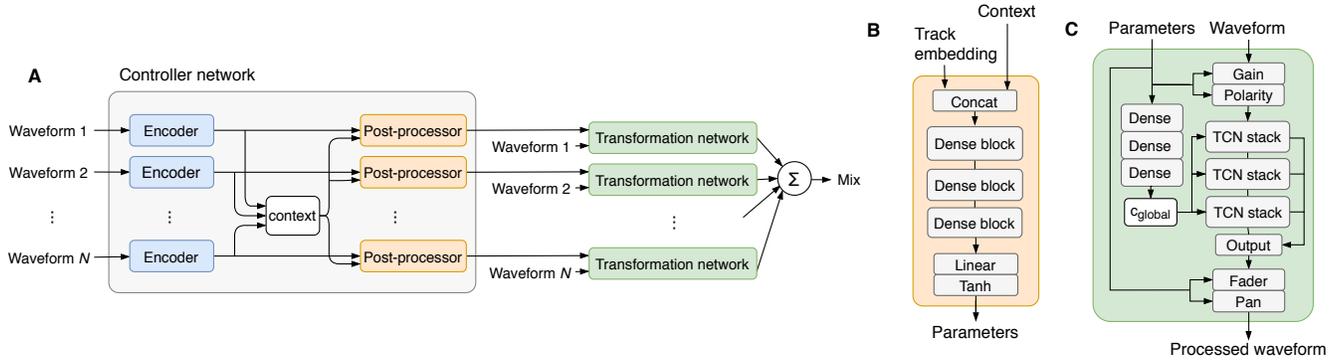
**Fig. 1**. Block diagrams of the DMC (A), its post-processor (B), and its transformation network (C).

channel in the mixing console. In order to train this network, we need the ability to compute gradients through the processors in the channel (e.g., equalization, compression, and reverb). This may be challenging, as implementations of these processors can be complex and varied. Additionally, they may not have easily-tractable or well-behaved gradients. To overcome this, we propose to replace each channel in the mixing console with a neural network that aims to emulate, as closely as possible, the processing of the original channel, which we will call the transformation network. We train this network by utilizing existing digital audio effects to generate training examples. We construct a differentiable mixing console (DMC), as shown in Fig. 1, where each channel has been replaced by an instance of the pre-trained transformation network. This enables us to train the controller network for the mixing task in an end-to-end fashion, without the need for the parameters used to create the ground truth mixes.

### 2.1. Transformation network

Various deep learning approaches have already been proposed for the task of modeling audio effects [12–17]. While previous approaches have focused on training a single model for each effect, we believe our work is the first to consider building a model that emulates a series connection of effects and their parameters, jointly. Most approaches do not consider modeling the different configurations of these devices, and those that do, only consider a sparse sampling of the parameters [14, 16]. This is due to the fact that they aim to emulate an analog device, and the process of collecting data at many configurations is often impractical. However, we are interested in modeling the behavior of digital signal processors. As a result, we can generate effectively endless examples during the training of the transformation network. To this end, we developed a Python package, pymixconsole[1], which implements a framework for controlling a chain of audio effects found in a typical mixing console, as shown in Fig. 3. Using audio recordings of various musical sources at 44.1 kHz from the SignalTrain dataset [14], we generate training examples across all configurations of the chain, on-the-fly, by processing these recordings with uniformly sampled configurations.

For the design of the transformation network, we follow a similar implementation to Damskägg et al. [18] for modeling distortion effects, which adapts a non-causal WaveNet-like model [19], formalized by Bai et al. [20] as the temporal convolutional network (TCN). This network is composed of blocks of 1-dimensional convolutions, as shown in Fig. 2. Exponentially increasing dilation factors and a kernel size of 15 are used to achieve a larger receptive field.
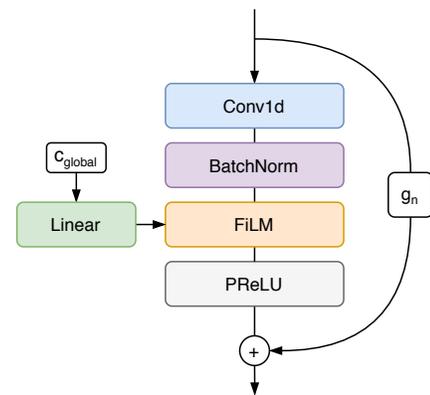
**Fig. 2**. Block diagram of the TCN block.

We utilize batch normalization without an affine transformation, and couple it with feature-wise linear modulation (FiLM) [21] in order to inject conditioning information from the effect parameters. The global conditioning $c_{global}$, shown in Fig. 2, is a vector generated by a small 3-layer multi-layer perceptron (MLP), which projects the signal chain parameters (e.g. 26 parameters for the complete channel) to a 128-dimensional vector. At each block, $c_{global}$ is projected via a linear layer to match the channel dimension for the FiLM operation. A residual connection with a learnable gain is also included.

We create multiple stacks of 10 TCN blocks to achieve a larger receptive field. The dilation factor $d$ of layer $l$ is given by $d_l = 2^{(l-1) \mod 10}$, where $\mod$ is the modulo operation. We consider three configurations of the TCN with 10, 20, and 30 blocks each: TCN-10, TCN-20, and TCN-30, which achieve receptive fields of 320 ms, 650 ms and 970 ms, respectively. Additionally, skip connections are included from the intermediate activations from every layer, where they are averaged before being added to the final layer.

Since the channel includes processing operations that are differentiable, like the input gain, polarity, fader, and panning, we implement these directly in the transformation network, before and after the TCN stacks, as shown in Fig. 1C. The TCN is trained to emulate only the more complex processors in the channel, shown in the dashed box in Fig. 3: the EQ, compressor, and reverb. With this configuration, we will define two different mixing tasks. First, a basic mix, where the TCN is removed and only the gain and panning are predicted by the controller, and second, a full mix, where the TCN is active, and parameters for all the processors are predicted by the controller to fully emulate the processing in the mixing console.
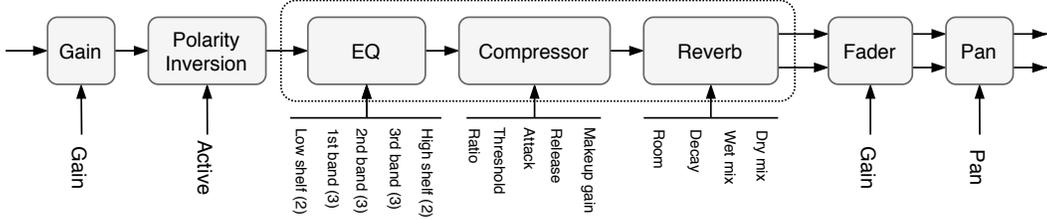
**Fig. 3**. Block diagram of the signal processing chain in the differentiable mixing console channel. The labels indicate parameters passed to each processor. The three processors in the dashed box are modeled by the transformation network, while the others are implemented directly, since they pose no challenge in backpropagation.

## 2.2. Controller network

The controller network contains a series of encoders and post-processors with shared weights, as shown in Fig. 1A. First, the encoder must learn to extract and aggregate information from the inputs that is salient for the mixing process. The encoder is constructed following the common spectrogram-based VGGish model [22], and we conduct transfer learning by using the pre-trained weights on AudioSet provided by Gemmeke et al. [23]. We found that fine-tuning these weights during training further improved performance.

The post-processor, shown in Fig. 1B, is a simple 3-layer MLP, with PReLU activations and dropout of 0.1. Recall that the weights of both the controller and transformation networks are shared across all input channels. This means the process of predicting parameters for each channel occurs independently on a per-track basis. Therefore, by default, cross-channel interactions cannot be captured, a critical consideration for creating a mix [24]. To address this issue, each instance of the post-processor is passed two inputs, the track embedding for the respective input channel, and an additional context embedding, which is computed by simply averaging the embeddings generated for all of the input sources. While this may obscure some of the information about the input sources in the mix, we found this provided sufficient context and worked in practice.

## 2.3. Stereo loss function

A critical step in the mixing process involves panning the elements in the mix across the stereo field. This must be done in such a way to achieve proper balance between the left and right channels, while also ensuring the appropriate spatialization of the elements. When training a model to generate stereo mixes using the L1 loss or multi-resolution STFT loss [25], there is an inherent issue: these loss functions applied to multi-channel audio heavily penalize mixes that place elements on the opposite side of the stereo field compared to the ground truth mix. This poses a challenge, since these loss functions encourage the model to always place sources in the center of the stereo field. Consider a ground truth mix with an electric guitar panned to the left. While the model may know that the guitar ought to be panned to the left or the right side and not the center, it has no way of predicting the absolute orientation in the ground truth mix. Therefore, to minimize the error, the model is incentivized to place it in the center, creating a more perceptually inaccurate mix.

To address this issue, we design a stereo loss function that aims to achieve left-right invariance, so only the overall stereo balance is considered. We first compute the sum and difference signals,

$$y_{\text{sum}} = y_{\text{left}} + y_{\text{right}}$$
$$y_{\text{diff}} = y_{\text{left}} - y_{\text{right}},$$

directly on the left and right channels of the time-domain signals. While this does transform the stereo information to another representation, the absolute stereo orientation is now represented in the phase of the difference signal. To ignore this phase information, we apply a multi-resolution STFT loss [25], which is composed of a spectral convergence (SC) and a spectral log-magnitude (SM) term:

$$\ell_{\text{SC}}(\hat{y}, y) = \frac{\| \, |\text{STFT}(y)| - |\text{STFT}(\hat{y})| \, \|_{\text{F}}}{\| \, |\text{STFT}(y)| \, \|_{\text{F}}}$$

$$\ell_{\text{SM}}(\hat{y}, y) = \frac{1}{N} \left\| \log \left( |\text{STFT}(y)| \right) - \log \left( |\text{STFT}(\hat{y})| \right) \right\|_1,$$

where $|\text{STFT}(\cdot)|$ is the short-time Fourier transform magnitude, $||\cdot||_{\text{F}}$ is the Frobenius norm, $||\cdot||_1$ is the L1 norm, $N$ is the number of STFT frames, and $\hat{y}$ and $y$ denote the predicted and target signals, respectively. To compute the final multi-resolution STFT loss, $M$ different STFT configurations are chosen with varying window, hop, and frame sizes, and the error at these resolutions is averaged:

$$\ell_{\text{MR}}(\hat{y}, y) = \frac{1}{M} \sum_{m=1}^{M} \left( \ell_{\text{SC}}(\hat{y}, y) + \ell_{\text{SM}}(\hat{y}, y) \right).$$

In all our experiments, we follow the original multi-resolution STFT implementation proposed by Yamamoto et al. [25], which includes three different STFT frame sizes of 512, 1024, and 2048. Applying $\ell_{\text{MR}}$ as defined above to both the sum and difference signals, we define the stereo loss function

$$\ell(\hat{y}, y) = \ell_{\text{MR}}(\hat{y}_{\text{sum}}, y_{\text{sum}}) + \ell_{\text{MR}}(\hat{y}_{\text{diff}}, y_{\text{diff}}).$$

## 2.4. Training

We begin by training the transformation network by minimizing the mean absolute error (MAE), or L1 loss, between the predicted and ground truth processed waveforms. Since the predicted waveforms are smaller than the input waveform due to the lack of padding in the model, we take a central crop from the ground truth that matches the size of the predicted waveform. We use Adam, a learning rate of $3 \cdot 10^{-4}$, and a batch size of 32, along with plateau learning rate scheduling, halving the learning rate after the validation loss has not decreased for 20 epochs, defining an epoch as 1,000 random 1.5-second patches from the dataset.

To train the DMC, we create instances of the pre-trained enocder and transformation network, along with the post-processor, for each input source in the multitrack input. We again use Adam and the same learning rate, but we scale the batch size due to memory constraints: for the basic task we use a batch size of 16 and for the full task we use a batch size of 2. For all DMC models, we again use plateau learning rate scheduling, halving the learning rate after the validation loss has not decreased for 200 epochs, and defining an epoch as 100 random 5-second patches.
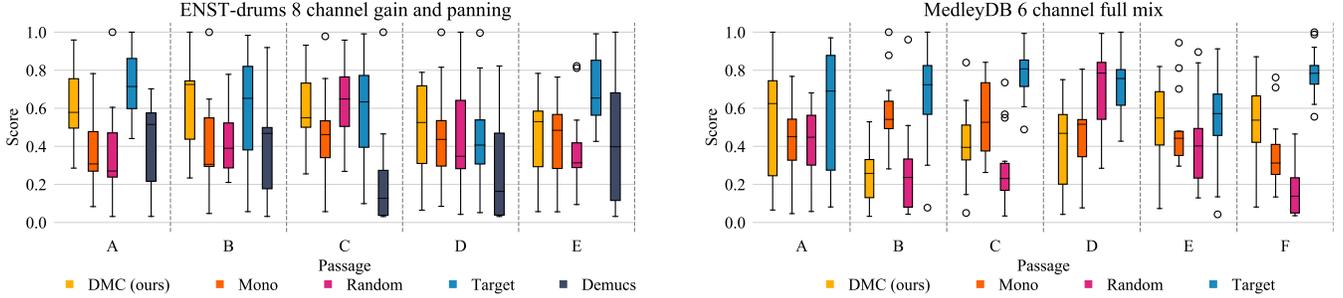
**Fig. 4**. Perceptual evaluation on the ENST-Drums gain and panning mixing task (left), and for MedleyDB full-channel mixing task (right).

## 3. EVALUATION AND RESULTS

### 3.1. Audio effect modeling

We evaluated variants of the TCN for channel emulation task by comparing the MAE, as well as the multi-resolution STFT distance, between the ground truth and predicted waveforms. We found that increasing the receptive field improved performance, with the TCN-30 (MAE: 0.024, STFT: 2.210) outperforming the shallower models, TCN-20 (MAE: 0.027, STFT: 2.315) and TCN-10 (MAE: 0.035, STFT: 2.701) Additionally, on a separate task of modeling an analog compressor, our TCN-20 (MAE: $4 \cdot 10^{-3}$, STFT: 0.606) outperformed the current state-of-the-art for this task, SignalTrain [14] (MAE: $8 \cdot 10^{-3}$, STFT: 1.657), by a substantial margin. Further details of these evaluations are presented in [26] (Ch. 4).

### 3.2. Multitrack mixing

**Datasets —** First, we consider the ENST-Drums dataset [27], which includes around 3 h of multitrack recordings of drummers. Each example in the dataset includes 8 sources from the drum kit. We generate training, validation, and test splits (80/10/10) following Moffat et al. [11]. This dataset provides an initial indication of the ability to learn mixing conventions from mixes with consistent sources and mixing techniques across the dataset. Next, to demonstrate the ability of our framework to generalize to real-world use cases, we further consider MedleyDB [28,29], which provides realistic and diverse multitrack recordings of complete songs across a number of genres. The dataset contains 196 songs with around 7 h of recordings. Due to memory constraints, we train models using songs with $\leq 16$ inputs for the basic task and $\leq 6$ inputs for the full task, resulting in 120 and 65 songs, respectively. Similarly to ENST-Drums, we create an 80/10/10 split of the data, plus we ensure that songs from the same artist do not fall into different splits.

**Baselines —** We consider three baselines. The first one is the mono mix, a sum of the inputs. The second one scales each input by a random gain ($-12$ dB to $+12$ dB), along with random panning. The third one is what we consider a canonical deep learning approach for processing time domain signals: we adapt the Demucs architecture [30], originally designed for source separation. Unlike our DMC model, this architecture does not present an inductive bias for the mixing task. To be comparable to the DMC model, we remove the LSTM layers from the center of the original network, and also scale down the number of channels, which results in around 80 M parameters. We pass a fixed number of input sources and train the model to predict the ground truth mix from the dataset. For songs with fewer inputs than we train with, we fill these inputs with zeros.

**Perceptual evaluation —** Due to the subjective nature of mixing, we conduct a perceptual evaluation with the Web Audio Evaluation

Tool [31] using the APE test design [32]. We enlisted 16 audio engineers with mixing experience. They were presented passages from the test set mixed by each approach, and were instructed to rate each on a scale from 0 to 1. More details are shown in [26] (Ch. 5).

**Results —** We first report results on the ENST-Drums dataset, on the basic mixing task in Fig. 4 (left). On average, the target mixes tend to be rated the highest, with mixes from our DMC following close behind, and even surpassing the target mixes in the case of passages B and D. On average, the mono and random mixes were rated lower, with the Demucs-like model being constantly rated the lowest (listeners indicated there were artifacts, likely from the transposed convolutions used by the model). To formalize these results, we perform the Kruskal-Wallis H-test, which points to a difference between the approaches ($F = 64.01$, $p = 8 \cdot 10^{-14}$). Our further ad-hoc analysis with Conover's test reveals there is not a significant difference between the target and DMC mixes ($p_{\text{adj}} = 0.08$).

We continue with MedleyDB and the full mixing task. We omit the Demucs-like model since it was unable to converge when training on MedleyDB, which we posit was a result of its lack of permutation invariance. Results are shown in Fig. 4 (right). The Kruskal-Wallis H-test again reveals that there is a difference between the approaches ($F = 48.1$, $p = 8.8 \cdot 10^{-10}$), and Conover's test reveals that all approaches perform differently from each other, with $p_{adj} = 1.1 \cdot 10^{-9}$ for the target vs. DMC comparison. Interestingly, in passages A and E, the DMC is nearly on par with the target. However, in B, C, and D, DMC performs poorly. In F, DMC performs clearly better than the baselines, but does not reach the level of the target. Note also, that random mixes include only gain and panning, and not the entire signal chain. Therefore, the DMC has a much more challenging task. Our listening suggests that failure cases arise from the over application of reverb on elements like the vocals, which listeners rated harshly. We provide listening examples in `https://csteinmetz1.github.io/dmc-icassp2021/`.

## 4. DISCUSSION

We outline and address a number of challenges in applying deep learning methods to build a model for automatic mixing trained directly on realistic multitrack data. We build a model with a strong inductive bias for this task, taking inspiration from the mixing console. By employing pre-trained sub-networks, weight sharing, and a stereo loss function, we demonstrate, to our knowledge for the first time, the ability to learn mixing conventions directly from waveforms of real-world multitracks. In the process, we demonstrate that the TCN can adequately model a series connection of effects over a dense sampling of their parameters. While the results on the complete mixing task are somewhat limited, we hypothesize that with larger models and more data, performance could be further improved.

# 5. REFERENCES

[1] Alex Case, *Mix smart: Pro audio tips for your multitrack mix*, Focal Press, 2011.

[2] Daniel A. Walzer, "Independent music production: how individuality, technology and creative entrepreneurship influence contemporary music industry practices," *Creative Industries Journal*, vol. 10, no. 1, pp. 21–39, 2017.

[3] Brecht De Man, Ryan Stables, and Joshua D. Reiss, *Intelligent Music Production*, Audio Engineering Society Presents. Taylor & Francis, 2019.

[4] David Moffat and Mark B. Sandler, "Approaches in intelligent music production," *Arts*, vol. 8, no. 4, 2019.

[5] Enrique Perez Gonzalez and Joshua D. Reiss, "Automatic gain and fader control for live mixing," in *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, 2009, pp. 1–4.

[6] Stuart Mansbridge, Saoirse Finn, and Joshua D. Reiss, "An autonomous system for multitrack stereo pan positioning," in *133rd Audio Engineering Society Convention (AES)*, 2012.

[7] Brecht De Man and Joshua D. Reiss, "A knowledge-engineered autonomous mixing system," in *135th Audio Engineering Society Convention (AES)*, 2013.

[8] Brecht De Man, *Towards a better understanding of mix engineering*, Ph.D. thesis, Queen Mary University of London, 2017.

[9] Bennett Kolasinski, "A framework for automatic mixing using timbral similarity measures and genetic optimization," in *124th Audio Engineering Society Convention (AES)*, 2008.

[10] Jeffrey Scott, Matthew Prockup, Erik M Schmidt, and Youngmoo E Kim, "Automatic multi-track mixing using linear dynamical systems," in *Sound and Music Computing Conf. (SMC)*, 2011.

[11] David Moffat and Mark Sandler, "Machine learning multitrack gain mixing of drums," in *147th Audio Engineering Society Convention (AES)*, 2019.

[12] John Covert and David L. Livingston, "A vacuum-tube guitar amplifier model using a recurrent neural network," in *Proc. of IEEE Southeastcon*, 2013, pp. 1–5.

[13] Zhichen Zhang, Edward Olbrych, Joseph Bruchalski, Thomas J. McCormick, and David L. Livingston, "A vacuum-tube guitar amplifier model using long/short-term memory networks," in *Proc. of IEEE SoutheastCon*, 2018, pp. 1–5.

[14] Scott H. Hawley, Benjamin L. Colburn, and Stylianos Ioannis Mimilakis, "Profiling audio compressors with deep neural networks," in *147th Audio Engineering Society Convention (AES*, 2019.

[15] Marco A Martínez Ramírez, Emmanouil Benetos, and Joshua D. Reiss, "Deep learning for black-box modeling of audio effects," *Applied Sciences*, vol. 10, pp. 638, 2020.

[16] Alec Wright, Eero-Pekka Damskägg, Lauri Juvela, and Vesa Välimäki, "Real-time guitar amplifier emulation with deep learning," *Applied Sciences*, vol. 10, no. 3, pp. 766, 2020.

[17] Alec Wright and Vesa Välimäki, "Perceptual loss function for neural modeling of audio systems," in *IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 251–255.

[18] Eero-Pekka Damskägg, Lauri Juvela, Vesa Välimäki, et al., "Real-time modeling of audio distortion circuits with deep learning," in *Sound and Music Computing Conf. (SMC)*, 2019.

[19] Dario Rethage, Jordi Pons, and Xavier Serra, "A WaveNet for speech denoising," in *IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, 2018, pp. 5069–5073.

[20] Shaojie Bai, J Zico Kolter, and Vladlen Koltun, "An empirical evaluation of generic convolutional and recurrent networks for sequence modeling," *arXiv:1803.01271*, 2018.

[21] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville, "FiLM: Visual reasoning with a general conditioning layer," in *32nd AAAI Conf. on Artificial Intelligence*, 2018.

[22] Shawn Hershey, Sourish Chaudhuri, Daniel PW Ellis, Jort F Gemmeke, Aren Jansen, R Channing Moore, Manoj Plakal, Devin Platt, Rif A Saurous, Bryan Seybold, et al., "CNN architectures for large-scale audio classification," in *IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, 2017, pp. 131–135.

[23] Jort F. Gemmeke, Daniel P. W. Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R. Channing Moore, Manoj Plakal, and Marvin Ritter, "AudioSet: An ontology and human-labeled dataset for audio events," in *IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, 2017.

[24] Sebastian Vega and Jordi Janer, "Quantifying masking in multi-track recordings," in *Sound and Music Computing Conf. (SMC)*, 2010.

[25] Ryuichi Yamamoto, Eunwoo Song, and Jae-Min Kim, "Parallel WaveGAN: A fast waveform generation model based on generative adversarial networks with multi-resolution spectrogram," in *IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 6199–6203.

[26] Christian J. Steinmetz, "Learning to mix with neural audio effects in the waveform domain," M.S. thesis, Universitat Pompeu Fabra, 2020, https://doi.org/10.5281/zenodo.4091203.

[27] Olivier Gillet and Gaël Richard, "ENST-Drums: An extensive audio-visual database for drum signals processing," in *Int. Society for Music Information Retrieval Conf. (ISMIR)*, 2006.

[28] Rachel Bittner, Justin Salamon, Mike Tierney, Matthias Mauch, Chris Cannam, and Juan Bello, "MedleyDB: A multitrack dataset for annotation-intensive MIR research," in *Int. Society for Music Information Retrieval Conf. (ISMIR)*, 2014.

[29] Rachel Bittner, Julia Wilkins, Hanna Yip, and Juan Pablo Bello, "MedleyDB 2.0: New data and a system for sustainable data collection," in *Int. Society for Music Information Retrieval Conf. (ISMIR)*, 2016.

[30] Alexandre Défossez, Nicolas Usunier, Léon Bottou, and Francis Bach, "Music source separation in the waveform domain," *arXiv:1911.13254*, 2019.

[31] Nicholas Jillings, Brecht De Man, David Moffat, and Joshua D. Reiss, "Web audio evaluation tool: A framework for subjective assessment of audio," in *Proc. of the 2nd Web Audio Conference (WAC)*, 2016.

[32] Brecht De Man and Joshua D. Reiss, "APE: Audio perceptual evaluation toolbox for MATLAB," in *136th Audio Engineering Society Convention (AES)*, 2014.