

DIFFERENTIABLE SIGNAL PROCESSING WITH BLACK-BOX AUDIO EFFECTS

Marco A. Martínez Ramírez^{b,*} Oliver Wang[#] Paris Smaragdis[‡] Nicholas J. Bryan[#]

[#]Adobe Research, USA

^bCentre for Digital Music, Queen Mary University of London, UK

[‡]University of Illinois at Urbana-Champaign, USA

ABSTRACT

We present a data-driven approach to automate audio signal processing by incorporating stateful third-party, audio effects as layers within a deep neural network. We then train a deep encoder to analyze input audio and control effect parameters to perform the desired signal manipulation, requiring only input-target paired audio data as supervision. To train our network with non-differentiable black-box effects layers, we use a fast, parallel stochastic gradient approximation scheme within a standard auto differentiation graph, yielding efficient end-to-end backpropagation. We demonstrate the power of our approach with three separate automatic audio production applications: tube amplifier emulation, automatic removal of breaths and pops from voice recordings, and automatic music mastering. We validate our results with a subjective listening test, showing our approach not only can enable new automatic audio effects tasks, but can yield results comparable to a specialized, state-of-the-art commercial solution for music mastering.

Index Terms— audio effects, deep learning, differentiable signal processing, black-box optimization, gradient approximation.

1. INTRODUCTION

Audio signal processing effects (Fx) are ubiquitous and used to manipulate different sound characteristics such as loudness, dynamics, frequency, and timbre across a variety of media. Many Fx, however, can be difficult to use or are simply not powerful enough to achieve a desired task, motivating new work. Past methods to address this include audio effects circuit modeling [1], analytical methods [2], and intelligent audio effects that dynamically change their parameter settings by exploiting sound engineering best practices [3]. The most common approach for the latter is adaptive audio effects or signal processing systems based on the modeling and automation of traditional processors [4]. More recent deep learning methods for audio effects modeling and intelligent audio effects include 1) end-to-end direct transformation methods [5, 6, 7], where a neural proxy learns and applies the transformation of an audio effect target 2) parameter estimators, where a deep neural network (DNN) predicts the parameter settings of an audio effect [8, 9, 10] and 3) differentiable digital signal processing (DDSP) [11, 12], where signal processing structures are implemented within a deep learning auto-differentiation framework and trained via backpropagation.

While these past methods are very promising, they are limited in several ways. First, direct transform approaches can require special, custom modeling strategies per effect (e.g. distortion), are often based on large and expensive networks, and/or use models with limited or no editable parameter control [5, 6]. Second, methods with

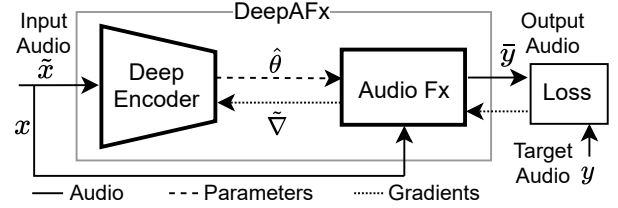


Fig. 1. Our *DeepAFx* method consists of a deep encoder that analyzes audio and predicts the parameters of one or more black-box audio effects (Fx) to achieve a desired audio processing task. At training time, gradients for black-box audio effects are approximated via a stochastic gradient method.

parameter control commonly require expensive human-labeled data from experts (e.g. inputs, control parameters, and target outputs) [7]. They are also typically optimized to minimize parameter prediction error and not audio quality directly [10], which can reduce performance. Third, DDSP approaches require a differentiable implementation for learning with backpropagation, re-implementation of each Fx, and in-depth knowledge, limiting use to known differentiable effects and causing high engineering effort [11, 12].

In this work, we introduce *DeepAFx*, a new approach to differentiable signal processing that enables the use of *arbitrary* stateful, third-party black-box processing effects (plugins) within any deep neural network as shown in Figure 1. Using our approach, we can mix-and-match traditional deep learning operators and black-box effects and learn to perform a variety of intelligent audio effects tasks without the need of neural proxies, re-implementation of Fx plugins, or expensive parameter labels (only paired input-target audio data is required). To do so, we present 1) a deep learning architecture in which an encoder analyzes input audio and learns to control Fx black-boxes that themselves perform signal manipulation 2) an end-to-end backpropagation method that allows differentiation through non-differentiable black-box audio effects layers via a fast, parallel stochastic gradient approximation scheme used within a standard auto differentiation graph 3) a training scheme that can support stateful black-box processors, 4) a delay-invariant loss to make our model more robust to effects layer group delay, 5) application of our approach to three audio production tasks including tube amplifier emulation, automatic non-speech sounds removal, and automatic music post-production or mastering, and 6) listening test results that shows our approach can not only enable new automatic audio production tasks, but is high-quality and comparable to a state-of-the-art proprietary commercial solution for automatic music mastering. Audio samples and code are online at <https://mchijmma.github.io/DeepAFx/>.

*This work was performed while interning at Adobe Research.

2. METHOD

2.1. Architecture

Our model is a sequence of two parts: a deep encoder and a black-box layer consisting of one or more audio effects. The encoder analyzes the input audio and predicts a set of parameters, which is fed into the audio effect(s) along with the input audio for transformation.

The **deep encoder** can be any deep learning architecture for learning representations of audio. To allow the deep encoder to learn long temporal dependencies, the input \tilde{x} consists of the current audio frame x centered within a larger audio frame containing previous and subsequent context samples. The input to the encoder consists of a log-scaled mel-spectrogram non-trainable layer followed by a batch normalization layer. The last layer of the encoder is a dense layer with P units and *sigmoid* activation, where P is the total number of parameters. The output is the predicted parameters $\hat{\theta}$ for the current input frame x .

The **audio Fx** layer is a stateful black-box consisting on one or more connected audio effects. This layer uses the input audio x and $\hat{\theta}$ to produce an output waveform $\bar{y} = f(x, \hat{\theta})$. We calculate a loss between the transformed audio and targeted output audio.

2.2. Gradient approximation

In order to train our model via backpropagation, we calculate the gradients of the loss function with respect to the audio Fx layer and the deep encoder. For the latter, we compute the gradient using standard automatic differentiation. For the former, gradient approximation has been addressed in various ways such as finite differences (FD) [13], FD together with evolution strategies [14], reinforcement learning [15], and by employing neural proxies [16]. For our work, we use a stochastic gradient approximation method called simultaneous permutation stochastic approximation (SPSA) [17], which is rooted in standard FD gradient methods and has low computational cost [13]. In our formulation, we ignore the gradient of the input signal x and only approximate the gradient of the parameters $\hat{\theta}$.

FD can be used to approximate the gradient of the audio Fx layer at a current set of parameters $\hat{\theta}_0$, where the gradient $\tilde{\nabla}$ is defined by the vector of partial derivatives characterized by $\tilde{\nabla} f(\hat{\theta}_0)_i = \partial f(\hat{\theta}_0) / \partial \theta_i$ for $i = 1, \dots, P$. The two-side FD approximation $\tilde{\nabla}^{FD}$ [17] is based on the backward and forward measurements of $f()$ perturbed by a constant ϵ . The i th component of $\tilde{\nabla}^{FD}$ is

$$\tilde{\nabla}^{FD} f(\hat{\theta}_0)_i = \frac{f(\hat{\theta}_0 + \epsilon \hat{d}_i^P) - f(\hat{\theta}_0 - \epsilon \hat{d}_i^P)}{2\epsilon}, \quad (1)$$

where $0 < \epsilon \ll 1$ and \hat{d}_i^P denotes a standard basis vector with 1 in the i th place and dimension P . From eq. (1), each parameter θ_i is perturbed one at a time and requires $2P$ function evaluations. This can be computationally expensive for even small P , given that $f()$ represents one or more audio effects. Even worse, when $f()$ is stateful, FD perturbation requires $2P$ instantiations of $f()$ to avoid state corruption, resulting in high memory costs.

SPSA offers an alternative and is an algorithm for stochastic optimization of multivariate systems. It has proven to be an efficient gradient approximation method that can be used with gradient descent for optimization [17]. The SPSA gradient estimator $\tilde{\nabla}^{SPSA}$ is based on the random perturbation of all the parameters $\hat{\theta}$ at the same time. The i th element of $\tilde{\nabla}^{SPSA}$ is

$$\tilde{\nabla}^{SPSA} f(\hat{\theta}_0)_i = \frac{f(\hat{\theta}_0 + \epsilon \hat{\Delta}_i^P) - f(\hat{\theta}_0 - \epsilon \hat{\Delta}_i^P)}{2\epsilon \hat{\Delta}_i^P}, \quad (2)$$

where $\hat{\Delta}_P$ is a P -dimensional random perturbation vector sampled from a symmetric Bernoulli distribution, i.e. $\Delta_i^P = \pm 1$ with probability of 0.5 [18]. In each iteration, the total number of function evaluations $f()$ is two, since the numerator in eq. (2) is identical for all the elements of $\tilde{\nabla}^{SPSA}$. This is especially valuable when P is large and/or when $f()$ is stateful. While the random search directions do not follow the steepest descent, over many iterations, errors tend to average to the sample mean of the stochastic process [17].

2.3. Training with stateful black-boxes

While training samples typically should be generated from an independent and identically distributed (i.i.d.) process, we note that audio effects are *stateful* systems. This means output samples depend on previous input samples or internal states [19]. Therefore, audio Fx layers should not be fed i.i.d.-sampled audio during training as this will corrupt state and yield unrealistic, transient outputs. As such, we feed consecutive non-overlapping frames of size N to the audio Fx layer, i.e. audio frames with a hop size of N samples, where the internal block size of each audio effects is set to a divisor of N . This ensures that there are no discrepancies between the behavior of the audio effects during training and inference time.

To train the model using mini-batches and maintain state, we use a separate Fx processor for each item in the batch. Therefore, for a batch size of M , we instantiate M independent audio effects for the forward pass of backpropagation. Each of the M plugins receives training samples from random audio clips, however across minibatches, each instance is fed samples from the same input until swapped with a new sample. When we consider the two-side function evaluations for SPSA, we meet the same stateful constraints by using two additional Fx per batch item. This gives a total of $3M$ audio effects when optimizing with SPSA gradients, whereas FD requires an unmanageable $(2P + 1)M$ effects for a large number of parameters P or batch size M . Finally, we parallelize our gradient operator per item in a minibatch, similar to recent distributed training strategies [20], but using a single one-GPU, multi-CPU setup.

2.4. Delay-invariant loss

Multiband audio effects correspond to audio processors that split the input signal into various frequency bands via different types of filters [19]. These filters can introduce group delay, which is a frequency-dependent time delay of the sinusoidal components of the input [21]. These types of effects and similar can also apply a 180° phase shift, i.e. to invert the sign of the input. While such properties are critical to the functioning of Fx, it means that difficulties can arise when directly applying a loss function with inexact inputs, either in time or frequency domain. Thus, we propose a delay-invariant loss function designed to mitigate these issues.

We first compute the time delay $\tau = \arg\max(\bar{y} \star y)$ between the output \bar{y} and target y audio frames via cross-correlation (\star). Then, the loss in the time domain

$$L_{time} = \min(\|\bar{y}_\tau - y_\tau\|_1, \|\bar{y}_\tau + y_\tau\|_1), \quad (3)$$

corresponds to the minimum $L1$ distance between the time-aligned target y_τ and both a 0° phase shift and 180° phase shift time-aligned output \bar{y}_τ . We compute \bar{Y}_τ and Y_τ : a 1024-point Fast Fourier Transform (FFT) magnitude of \bar{y}_τ and y_τ , respectively. The loss in the frequency domain L_{freq} is defined as

$$L_{freq} = \|\bar{Y}_\tau - Y_\tau\|_2 + \|\log \bar{Y}_\tau - \log Y_\tau\|_2. \quad (4)$$

The final loss function is $L = \alpha_1 L_{time} + \alpha_2 L_{freq}$, where we empirically tune the weighting and use $\alpha_1 = 10$ and $\alpha_2 = 1$.

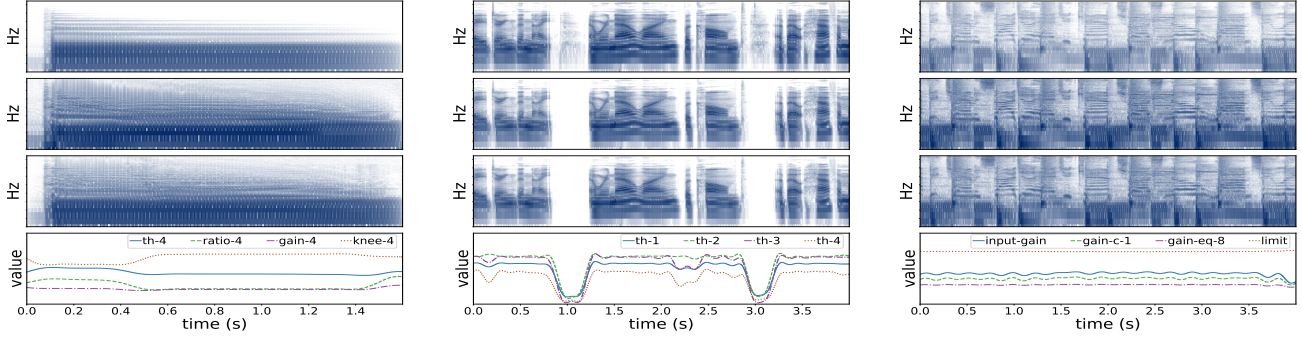


Fig. 2. (Left) Tube amplifier emulation. (Center) Automatic non-speech sounds removal. (Right) Automatic music mastering. From top to bottom: input, target, our result, and selected predicted parameters. th , and $gain-c$ mean threshold and compressor makeup gain, respectively.

3. EXPERIMENTS

We apply our approach to three applications including tube amplifier emulation, automatic non-speech sounds removal, and automatic music mastering. In each case, we use one or more different effects and different datasets to illustrate the generality of our approach. We also explore two deep encoder variants; an *Inception* network [22] and *MobileNetV2* [23]. The number of parameters for each encoder is approx. 2.8M and 2.2M, respectively. The input context \tilde{x} and current audio frame x are 40960 and 1024 samples (1.85 sec and 46 ms) at a 22,050Hz sampling rate. The log-scaled mel-spectrogram input layer has a 46 ms window size, 25% overlap, and 128 mel-bands. We use effects from the *LV2* audio plugin open standard [24] and experiment with continuous parameters scaled between 0 and 1. At inference we apply a lowpass filter to $\hat{\theta}$ to ensure smoothness. For intuition, our final optimization scheme with the Inception encoder and parametric equalizer Fx (EQ) [25] takes 3 minutes per epoch (1000 steps) to train with batch size $M = 100$ on a Tesla-V100 GPU. We use Adam optimization and early stop on validation loss.

3.1. Tube amplifier emulation

The emulation of distortion effects is a highly-researched field [1, 2] whose goal is to recreate the sound of analog reference devices (e.g. tube amplifiers or distortion circuits). Amplifier emulation using deep learning has been explored in the context of deep neural proxies [26, 27] and been shown to achieve virtually indistinguishable emulation under certain conditions [5, 6]. Such approaches, however, are *specifically* designed for this task, making it an interesting benchmark. We ask – is it possible to emulate distortion using a simple multiband dynamic range compressor?

A compressor modifies the amplitude dynamics of audio by applying a time-varying gain. Compressors are commonly used for loudness control and typically introduce little harmonic distortion in contrast to tube amplifiers [19]. Nonetheless, we apply our approach and learn 21 parameters; the *threshold*, *makeup gain*, *ratio*, and *knee* for each of the 4 frequency bands; the 3 *frequency splits*; and the *input* and *output gains* on a multiband compressor [28]. For training, we use the subset of the IDMT-SMT-Audio-Effects dataset [29] used by Martínez et al. [5] consisting of 1250 raw notes from various electric guitars and bass guitars and processed through a *Universal Audio 6176 Vintage Channel Strip* tube preamplifier. The train, validation, and test dataset sizes are 31.6, 4.0, and 4.0 minutes, respectively. For all tasks, the non-trainable Fx parameters are set to defaults with the exception of the attack and release gates, which are set to their minimum (10 ms) due to the fast control rate of our model (46 ms).

3.2. Automatic non-speech sound removal

Removing non-speech vocal sounds, such as breaths and lip smacks, is a common task performed by sound engineers [30]. This task can be done by manually editing the audio waveform or by using a *noise gate*, to reduce signal below a certain *threshold* via a *reduction gain* and *ratio* setting [19]. Both approaches can be time consuming, require expert knowledge, and, to the best of our knowledge, the automation of this task has not been investigated with the tangential exception of an intelligent noise gate for drum recordings [31] and automatic detection of breaths for speaker recognition [32, 33].

We train a model with a *multiband noise gate* [28] Fx layer with 17 parameters (*threshold*, *reduction gain* and *ratio* for each of the 4 frequency bands; the 3 *frequency splits*; and the *input* and *output gains*). We use the *DAPS* dataset [34] which contains 100 raw and clean speech recordings with manually removed breaths and lip smacks. The train, validation, and test dataset size is 213.5, 30.2, and 23.8 minutes, respectively.

3.3. Automatic music mastering

Music post-production or mastering is the process of enhancing a recording by manipulating its dynamics and frequency content [3]. This manipulation is typically done by an experienced mastering engineer and is carried out using dynamic range effects, such as a compressor and limiter; and frequency-based processors, such as EQ. Automatic mastering has been explored using adaptive audio effects [35, 36], using DNNs to predict dynamic range compression gain coefficients [37], and even commercialized as proprietary online mastering software (OMS) [38].

We train a model with multiple audio effects in series, similar to the way mastering engineers perform this task, consisting of a *multiband compressor* [28], a *graphic EQ* [25] and *mono limiter* [39]. A limiter is a dynamic range processor that more aggressively adjusts any input above the threshold, e.g. by using a ratio of ∞ [40]. We train our model to learn 50 parameters – 16 parameters for the multiband compressor (*threshold*, *makeup gain* and *ratio* for each of the 4 frequency bands, the 3 *frequency splits*, and the *input gain*), 33 parameters for the *graphic EQ* (the *gain* for each of the 32 frequency bands and the *output gain*), and 1 parameter for the limiter (*threshold*). For training data, we collected 138 unmastered and mastered music tracks from *The ‘Mixing Secrets’ Free Multitrack Download Library* [41], and as preprocessing step, we perform time-alignment using cross-correlation and loudness normalize each unmastered track to -25dBFS. The train, validation, and test dataset size is 429.3, 51.1, and 50.3 minutes, and respectively.

	Model	Epochs	Time	\tilde{d}_{MFCC}
Tube amplifier emulation	Inception	97	9.07	0.2596
	MobileNetV2	63	6.4	0.2186
	CAFx	723	5.5	0.0826
Non-speech sounds removal	Inception	89	7.4	0.0186
	MobileNetV2	60	4.8	0.0231
Music mastering	Inception	202	19.8	0.0282
	MobileNetV2	178	17.5	0.0542
	OMS	-	-	0.0157

Table 1. Training epochs & time (hours), and test MFCC distance.

4. RESULTS & ANALYSIS

4.1. Qualitative Evaluation

We show qualitative analysis by visual inspection via Figure 2, in which we depict input, target and predicted spectrograms for all tasks. For the tube amplifier task, it can be seen that the model emulates very closely the analog target. Also, the multiband compressor parameters move over time along the attack and decay of the guitar note to emulate the nonlinear hysteresis behavior of the tube amplifier. For the non-speech sound removal task, the model successfully removes breath and lip smacking sounds. It can be seen how the different thresholds vary according to the non-speech sounds present at the first and third second of the input recording. The music mastering example also successfully matches the sound engineer mastering target and the parameters of the three audio effects gradually change based on the input content of the unmastered music track.

4.2. Quantitative Evaluation

We show quantitative evaluation in Table 1, including number of training epochs until early stopping, training time in hours, and the mean cosine distance of the mel-frequency cepstral coefficients or \tilde{d}_{MFCC} as a proxy for a perceptual metric, following past work [5]. For this, we compute 13 MFCCs from a log-power mel-spectrogram using a window size of 1024 samples, 25% hop size and 128 bands. As a baseline for the tube amplifier emulation task, we use the *Convolutional Audio Effects Modeling Network (CAFx)* [5], and for the music mastering task, we use an online mastering software (OMS) [38]. We see that the tube amplifier emulation distance for our approach is higher than the other two tasks, likely caused by using a compressor to achieve distortion. We see both encoders achieved similar performance, although the Inception model tends to perform slightly better, and all training times are under a day. We also find the CAFx and OMS models have lower distance, but yield to perceptual listening tests for further conclusions.

4.3. Perceptual Evaluation

We performed a subjective listening test to evaluate the perceptual quality of our results using the multiple stimulus hidden reference (MUSHRA) protocol [42]. Seventeen participants took part in the experiment which was conducted using The Web Audio Evaluation Tool [43]. Five participants identified as female and twelve participants identified as male. Participants were among musicians, sound engineers, or experienced in critical listening.

Five listening samples were used from each of the three test subsets. Each sample consisted of a 4-second segment with the exception of the tube amplifier emulation samples, where each note is 2 seconds long. Participants were asked to rate samples according to

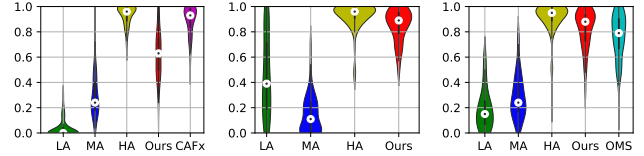


Fig. 3. Listening test violins plots. (Left) Tube amplifier emulation results show our approach does not outperform CAFx [5], but still yields good quality. (Center) Non-speech sounds removal results show near human quality. (Right) Music mastering results show our method performs equivalent to the state-of-the-art (OMS).

similarity relative to a reference sound, i.e. a sample processed by a sound engineer or a recording from an analog device.

The samples consisted of an unprocessed sample as low-anchor (LA), a mid-anchor (MA), a hidden reference as high-anchor (HA), and the output of our DeepAFx method with the Inception encoder (Ours). The mid-anchor corresponds to a *tube distortion* [44] with drive of +8 dB; a *gate* [45] with default settings; and peak normalization to -1dBFS for each task, respectively. All samples for the first two tasks were loudness normalised to -23dBFS. The samples for the music mastering task were monophonic and not loudness normalized so as to evaluate loudness as part of the task.

The results of the listening test can be seen in Figure 3 as a violin plot. For tube amplifier emulation, the median scores in figure order are (0.0, 0.24, 0.96, 0.63, 0.93), respectively. The CAFx method ranks above our model, which is reasonable considering it corresponds to a custom, state-of-the-art network. Our result for this task is still remarkable, however, since we are only using a multiband compressor in a field where complex models or deep neural proxies have predominated. For the non-speech sounds removal task, the median scores in figure order are (0.39, 0.11, 0.89, 0.96), respectively. Our model is rated almost as high as the hidden-anchor, which is again remarkable given the typical expertise and time required for the task. Also note, 1) the mid-anchor is rated lower than the low-anchor because of agitating gating artifacts from Fx processing and 2) the high variance of the low anchor, which contains different speech content (breaths and voice pops) and appears to confuse listeners when rating similarity. For the automatic music mastering task, the median scores in figure order are (0.15, 0.24, 0.95, 0.88, 0.79). Our model is rated above OMS, which we consider state of the art. We further conducted multiple post-hoc paired t-tests with Bonferroni correction [46] for each task and condition vs. our method. We found the rank ordering of the results are significant for $p < 0.01$, with the exception of our approach compared to OMS ($p \approx 0.0378$).

5. CONCLUSION

We present a new approach to differentiable signal processing that enables the use of stateful, third-party black-box audio effects (plugins) within any deep neural network. The main advantage of our approach is that it can be used with general, unknown plugins, and can be retrained for a wide variety of tasks by substituting in new audio training data and/or customizing with specific Fx. We believe that this framework has the potential for broader use, including optimizing legacy signal processing pipelines, style matching, and generative modeling for audio synthesis. For future work, we hope to address limitations including optimization of discrete parameters, gradient approx. accuracy, removing the need for paired training data, and robustness to black-box software brittleness.

6. REFERENCES

- [1] D. T. Yeh, J. S. Abel, A. Vladimirescu, and J. O. Smith III, "Numerical methods for simulation of guitar distortion circuits," *Computer Music Journal*, 2008.
- [2] F. Eichas and U. Zölzer, "Black-box modeling of distortion circuits with block-oriented models," in *DAFx*, 2016.
- [3] B. De Man, R. Stables, and J. D. Reiss, *Intelligent Music Production*, Focal Press, 2020.
- [4] D. Moffat and Mark B Sandler, "Approaches in intelligent music production," in *Arts*. MDPI, 2019.
- [5] M. A. Martínez Ramírez, E. Benetos, and J. D. Reiss, "Deep learning for black-box modeling of audio effects," in *Applied Sciences*. MDPI, 2020.
- [6] A. Wright, E.-P. Damskägg, L. Juvela, and V. Välimäki, "Real-time guitar amplifier emulation with deep learning," in *Applied Sciences*. MDPI, 2020.
- [7] S. H. Hawley, B. Colburn, and S. I. Mimitakis, "SIGNALTRAIN: Profiling audio compressors with deep neural networks," in *AES Convention*, 2019.
- [8] J. Rämö and V. Välimäki, "Neural third-octave graphic equalizer," in *DAFx*, 2019.
- [9] D. Sheng and G. Fazekas, "A feature learning siamese model for intelligent control of the dynamic range compressor," in *IJCNN*, 2019.
- [10] S. I. Mimitakis, N. J. Bryan, and P. Smaragdis, "One-shot parametric audio production style transfer with application to frequency equalization," in *ICASSP*, 2020.
- [11] J. Engel, L. Hantrakul, C. Gu, and A. Roberts, "DDSP: Differentiable digital signal processing," in *ICLR*, 2020.
- [12] B. Kuznetsov, J. D. Parker, and F. Esqueda, "Differentiable IIR filters for machine learning applications," in *DAFx*, 2020.
- [13] L. M. Milne-Thomson, *The calculus of finite differences*, American Mathematical Soc., 2000.
- [14] T. Salimans, J. Ho, X. Chen, S. Sidor, and I. Sutskever, "Evolution strategies as a scalable alternative to reinforcement learning," in *arXiv:1703.03864*, 2017.
- [15] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," *ML*, 1992.
- [16] A. Jacovi and et al., "Neural network gradient-based learning of black-box function interfaces," in *ICLR*, 2019.
- [17] J. C. Spall, "An overview of the simultaneous perturbation method for efficient optimization," *The Johns Hopkins APL Technical Digest*, 1998.
- [18] J. C. Spall, "Multivariate stochastic approximation using a simultaneous perturbation gradient approximation," *TAC*, 1992.
- [19] U. Zölzer, *DAFX: digital audio effects*, J. Wiley & Sons, 2011.
- [20] Y. You, Z. Zhang, C.-J. Hsieh, J. Demmel, and K. Keutzer, "Imagenet training in minutes," in *ICPP*, 2018.
- [21] J. O. Smith III, *Introduction to Digital Filters: with Audio Applications*, vol. 2, W3K Publishing, 2007.
- [22] J. Lee, N. J. Bryan, J. Salamon, Z. Jin, and J. Nam, "Metric learning vs classification for disentangled music representation learning," in *ISMIR*, 2020.
- [23] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "MobileNetV2: Inverted residuals and linear bottlenecks," in *CVPR*, 2018.
- [24] D. Robillard, "LV2 plugins," 2015, <http://lv2plug.in/>.
- [25] V. Sadovnikov, "LSP parametric equalizer x32 mono," 2015, <https://lsp-plugin.in/>.
- [26] M. A. Martínez Ramírez and J. D. Reiss, "Modeling of non-linear audio effects with end-to-end deep neural networks," in *ICASSP*, 2019.
- [27] E.-P. Damskägg, L. Juvela, E. Thuillier, and V. Välimäki, "Deep learning for tube amplifier emulation," in *ICASSP*, 2019.
- [28] K. Foltman, M. Schmidt, C. Holschuh, T. H. Johanssen, and T. Szilagy, "Calf studio gear," 2020, <http://calf-studio-gear.org/doc>.
- [29] M. Stein, J. Abeßer, C. Dittmar, and G. Schuller, "Automatic detection of audio effects in guitar and bass recordings," in *AES Convention*, 2010.
- [30] B. Owsinski, *The Mixing Engineer's Handbook*, Nelson Education, 2013.
- [31] M. Terrell, J. D. Reiss, and M. Sandler, "Automatic noise gate settings for drum recordings containing bleed from secondary sources," *EURASIP JASP*, 2011.
- [32] S. H. Dumpala and K.N.R.K. R. Alluri, "An algorithm for detection of breath sounds in spontaneous speech with application to speaker recognition," in *SPECOM*, 2017.
- [33] É. Székely, G. E. Henter, and J. Gustafson, "Casting to corpus: Segmenting and selecting spontaneous dialogue for TTS with a CNN-LSTM speaker-dependent breath detector," in *ICASSP*, 2019.
- [34] G. J. Mysore, "Can we automatically transform speech recorded on common consumer devices in real-world environments into professional production quality speech?—a dataset, insights, and challenges," *SPL*, 2014.
- [35] S. I. Mimitakis, K. Drossos, A. Floros, and D. Katerelos, "Automated tonal balance enhancement for audio mastering applications," in *AES Convention*, 2013.
- [36] M. Hilsamer and S. Herzog, "A statistical approach to automated offline dynamic processing in the audio mastering process," in *DAFx*, 2014.
- [37] S. I. Mimitakis, K. Drossos, T. Virtanen, and G. Schuller, "Deep neural networks for dynamic range compression in mastering applications," in *AES Convention*, 2016.
- [38] "LANDR," 2020, <https://www.landr.com/>.
- [39] V. Sadovnikov, "LSP limiter mono," 2015, <https://lsp-plugin.in/>.
- [40] J. D. Reiss and A. McPherson, *Audio effects: Theory, Implementation and Application*, CRC Press, 2014.
- [41] M. Senior, "Mixing secrets for the small studio," 2011, Taylor & Francis, <https://www.cambridge-mt.com/ms/mtk/>.
- [42] B. Series, "Recommendation ITU BS.1534-3," 2014.
- [43] N. Jillings, B. De Man, D. Moffat, and J. D. Reiss, "Web Audio Evaluation Tool: A browser-based listening test environment," in *SMC*, 2015.
- [44] "Invada tube distortion," 2020, <https://launchpad.net/invada-studio/>.
- [45] V. Sadovnikov, "LSP gate mono manual," 2015, https://lsp-plugin.in/?page=manuals§ion=gate_mono.
- [46] S. Holm, "A simple sequentially rejective multiple test procedure," *Scandinavian Journal of Statistics*, 1979.