

LIGHTWEIGHT AND INTERPRETABLE NEURAL MODELING OF AN AUDIO DISTORTION EFFECT USING HYPERCONDITIONED DIFFERENTIABLE BIQUADS

Shahan Nercessian, Andy Sarroff, and Kurt James Werner

iZotope, Inc., Cambridge, MA

ABSTRACT

In this work, we propose using differentiable cascaded biquads to model an audio distortion effect. We extend trainable infinite impulse response (IIR) filters to the hyperconditioned case, in which a transformation is learned to directly map external parameters of the distortion effect to its internal filter and gain parameters, along with activations necessary to ensure filter stability. We propose a novel, efficient training scheme of IIR filters by means of a Fourier transform. Our models have significantly fewer parameters and reduced complexity relative to more traditional black-box neural audio effect modeling methodologies using finite impulse response filters. Our smallest, best-performing model adequately models a BOSS MT-2 pedal at 44.1 kHz, using a total of 40 biquads and only 210 parameters. Its model parameters are interpretable, can be related back to the original analog audio circuit, and can even be intuitively altered by machine learning non-specialists after model training. Quantitative and qualitative results illustrate the effectiveness of the proposed method.

Index Terms— digital audio effects, IIR filters, hyperconditioning, differentiable DSP

1. INTRODUCTION

Digital emulation of analog audio circuits is a topic of extensive study within the audio signal processing community. Such emulations allow for the widespread use of audio effects without the need for dedicated hardware. Virtual analog techniques including wave digital filters [1, 2, 3] are often used to model said circuits. While these approaches can provide high fidelity emulations, they require significant domain knowledge, and often employ expensive online solvers [4].

Deep learning has emerged as an alternative, data-driven technique for modeling audio effects. In [5, 6], generic black-box architectures were proposed to model a variety of effects, but use bidirectional layers having poor real-time implications. A feedforward variant of WaveNet was proposed for emulating a tube amp [4] and various distortion pedals [7]. The architecture boasts real-time performance and can model external parameters of the underlying effects, but still imposes a significant amount of complexity. Moreover, these

black-box models lack interpretability, making them difficult to understand or troubleshoot.

Recently, there has been a push towards implementing differentiable audio processors in deep learning frameworks, enabling end-to-end training [8]. [9] considered learning fixed infinite impulse response (IIR) filters to approximate a single configuration of an audio effect, and thus, cannot model external parameters. We have also found (and verified through talks with the original authors) that their implementation of IIR filters as recurrent layers restricts the number of filters that can be practically learned, limiting their applicability [10].

In this paper, we propose modeling an audio distortion effect using multiple differentiable biquads. We continue the trend towards lightweight, interpretable audio deep learning models, using several IIR filters as building blocks. Unlike previous attempts to learn trainable IIR filters using deep learning [9, 11], we use hyperconditioning (HC) [12, 13] in order to capture changes in filtering stages based on the user-facing parameters controlling the effect. We consider many cascaded biquad (CB) representations and introduce activation functions necessary for their stable training. We propose an efficient training scheme using frequency domain sampling, which is more tractable than a naïve, prohibitively slow training using recurrent layers. The resulting models are lightweight and interpretable, yet perform competitively with a black-box baseline model.

The remainder of this paper reads as follows: We review biquads in Sec. 2. We introduce the proposed method in Sec. 3. We illustrate its effectiveness, comparing it to a black-box baseline in Sec. 4. We draw conclusions in Sec. 5.

2. BACKGROUND INFORMATION

Biquads (2nd order IIR filters) have the difference equation

$$y[n] = b_0x[n] + b_1x[n-1] + b_2x[n-2] - a_1y[n-1] - a_2y[n-2] \quad (1)$$

where coefficients $\mathbf{b} = [b_0, b_1, b_2]$ and $\mathbf{a} = [1, a_1, a_2]$ are the feedforward and feedback gains of the filter, respectively [14]. For real \mathbf{b} and \mathbf{a} , the transfer function of a biquad is given by

$$H(z) = \frac{b_0 + b_1z^{-1} + b_2z^{-2}}{1 + a_1z^{-1} + a_2z^{-2}} = b_0 \frac{(z - q_0)(z - q_1)}{(z - p_0)(z - p_1)} \quad (2)$$

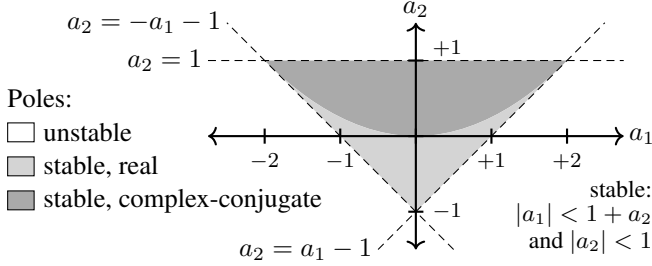


Fig. 1. Biquad triangle ensuring filter stability.

where pairs of zeros $\{q_0, q_1\}$ and poles $\{p_0, p_1\}$ dictate the nature of the filtering, and b_0 controls its gain. Stability is ensured so long as both $|p_0| < 1$ and $|p_1| < 1$, achieved when a_1 and a_2 are constrained to be inside the so-called biquad triangle (Fig. 1).

The frequency response of a biquad can be evaluated over a vector of linearly-spaced digital frequencies ω by evaluating Eq. 2 with $e^{j\omega}$ as its argument. Higher order IIR filters can be created by cascading K biquads in series. The composite complex frequency response of such a cascade is given by

$$H_{\text{cascade}}(e^{j\omega}) = \prod_{k=0}^{K-1} H_k(e^{j\omega}) \quad (3)$$

which we will refer to as the cascaded `freqz` operation. While biquads model arbitrary 2nd order filters, the RBJ cookbook [15] defines explicit formulae for parametric peaking and shelving filters with specified center/cutoff frequency f in Hz, gain g in decibels (dB), and Q/slope Q at a given sampling rate f_{SR} . The parameterization of these filters is interpretable for end-users and stable by design. The design of cascaded peaking and shelving filters has been explored in several works [16, 17, 18], which form the basis for the ubiquitous parametric equalizer (EQ).

3. MODELING USING DIFFERENTIABLE BIQUADS

A high-level block diagram of the proposed model is illustrated in Fig. 2. The input audio \mathbf{x} is subjected to a parametric delay layer, yielding \mathbf{y}_{-1} , followed by S filtering stages. Each filtering stage s consists of a cascade of K biquads, a multiplicative gain α_s (specified in dB), and a \tanh nonlinearity (NL), except for the final stage where the NL is omitted. Each biquad is described by P filter parameters, whose value depends on the specific representation that is used. In addition to the audio signal from the previous stage \mathbf{y}_{s-1} , each filtering stage also takes as input C_s external parameters from a user, forming the conditioning vector \mathbf{c}_s , which acts as a signal capable of altering the underlying parameters of the filtering stage. The vector \mathbf{c}_s can be made to vary over time to accommodate parameter automation. The output of the final stage $\mathbf{y} = \mathbf{y}_{S-1}$ is an estimate of the ground truth output from the effect that is being modeled.

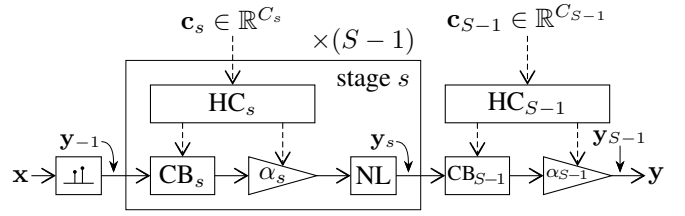


Fig. 2. Proposed neural network model architecture.

3.1. Parametric delay layer

The input is passed through a delay layer parameterized by 2 trainable parameters: a delay (in non-integer samples) and a linear gain. The delay layer is critical for accommodating any potential timing offsets between inputs and outputs in the dataset. The delay parameter maps to a linear fractional delay filter consisting of 2 taps which sum to 1 [19]. The delay filter computation could be extended to a more generic sinc interpolation kernel or be replaced with a non-sparse convolution, but neither are considered here in order to minimize complexity. The filtered result is multiplied by a linear gain term, which in particular, gives the network the ability to potentially invert the signal. As the \tanh NL used in each block is symmetric, we found this to be sufficient for modeling a circuit containing an odd number of inverting amplifier stages.

3.2. Hyperconditioned blocks

Previous attempts to learn IIR filters in audio deep learning [9, 11] have only learned fixed filters which do not change with external user input. Conventionally, conditioning information is injected into networks via addition [20], concatenation [21], or adaptive instance normalization [22]. Conversely, HC learns transformations of conditioning signals which ultimately control the underlying model weights (for instance, the coefficients of a convolutional kernel) [12]. This type of adaptation is more direct and expressive, validating its use in recent audio style transfer applications [13]. We inject conditioning information (i.e. the values of user-facing effect parameters scaled to be in the range $[0, 1]$) into the model by applying HC to the cascade of biquads and gain in a filtering stage. Specifically, we transform conditioning vectors into underlying parameters at filtering stage s using a single affine transformation, implemented as a 1×1 convolution.

Inspection of the schematic of the circuit being modeled can help appropriately set the value of S and K , and to determine which filtering stages have external parameters as input. Injecting the knowledge of the conditioning sites into the model can further reduce parameter counts and impose more inductive bias into the model. For blocks with $C_s > 0$, we use HC networks to transform conditioning signals into their corresponding parameters. Otherwise, we simply use trainable weights which learn a fixed gain and filtering for the stage.

3.3. Cascaded biquad representations

We consider 3 parameterizations of CBs in this work. While [9] suggests that filter stability can be largely ensured with proper initialization when learning fixed IIR filter configurations (as network optimization will naturally promote learning the stable solution), this is not sufficient for hyperconditioned filters that change with user input, as in this work. We ensure stable training and inference by subjecting parameters to representation-specific activations imposing their stability constraints. We generically refer to these activations as `p2c`, characterizing all of the operations that are performed to map parameters to their ultimate filter coefficients and gains. Small “epsilon” terms avoid marginal stability and/or undefined behavior, but we omit them here for readability.

3.3.1. Coefficient representation

In the 1st parameterization, we characterize the k th biquad at stage s by their coefficient vectors $\mathbf{b}_{s,k}$ and $\mathbf{a}_{s,k}$. Since each CB is followed by a gain, we set $b_{0,s,k} = 1$ for all filters with no remarkable loss to model expression. The vector of filter parameters is $\mathbf{v}_{s,k} = [b_{1,s,k}, b_{2,s,k}, a_{1,s,k}, a_{2,s,k}]$ and $P = 4$. In order to maintain stability, we must enforce that $a_{1,s,k}$ and $a_{2,s,k}$ are inside the biquad triangle. This is accomplished by performing the following activations in order:

$$a_{1,s,k} \leftarrow 2 \cdot \tanh a_{1,s,k} \quad (4)$$

$$a_{2,s,k} \leftarrow [(2 - |a_{1,s,k}|) \tanh a_{2,s,k} + |a_{1,s,k}|] / 2 \quad (5)$$

3.3.2. Pole/zero representation

The 2nd parameterization characterizes each biquad by the real and imaginary part of a single pole and zero, constraining the remaining pole and zero to be their complex conjugates. In this case, $\mathbf{v}_{s,k} = [\Re\{q_{s,k}\}, \Im\{q_{s,k}\}, \Re\{p_{s,k}\}, \Im\{p_{s,k}\}]$, and again, $P = 4$. Stability is ensured by constraining $|p_{s,k}| < 1$ using the activation function

$$p_{s,k} \leftarrow p_{s,k} \tanh (|p_{s,k}|) / |p_{s,k}| \quad (6)$$

and we compute $\mathbf{b}_{s,k} = [1, -2\Re\{q_{s,k}\}, \Re\{q_{s,k}\}^2 + \Im\{q_{s,k}\}^2]$, $\mathbf{a}_{s,k} = [1, -2\Re\{p_{s,k}\}, \Re\{p_{s,k}\}^2 + \Im\{p_{s,k}\}^2]$.

3.3.3. Parametric EQ representation

In our final parameterization, we model each CB with a parametric EQ. Each parametric EQ consists of a low shelf, a high shelf, and $K - 2$ peaking filters. In this case, we have $\mathbf{v}_{s,k} = [f_{s,k}, g_{s,k}, Q_{s,k}]$ and $P = 3$. The parametric EQ parameters are converted to biquad filter coefficients using the RBJ cookbook formulae [15].

We apply an activation on $f_{s,k}$ such that it is in $[0, f_{\text{SR}}/2]$ and $f_{s,0} \leq f_{s,1} \leq \dots \leq f_{s,K-1}$, defined as

$$f_{s,k} \leftarrow (f_{\text{SR}}/K) \sum_{k=0}^{K-1} \lfloor |f_{s,k}| - f_{s,k} \rfloor \quad (7)$$

where $\lfloor \cdot \rfloor$ denotes the round operator. No activation is applied to $g_{s,k}$, though gains specified in dB are inherently converted back to a linear scale in the RBJ cookbook formulae. Values for $Q_{s,k}$ must be non-negative, and for shelving filters, must also be ≤ 1 . Accordingly, $Q_{s,k}$ is subjected to the activation

$$Q_{s,k} \leftarrow Q_{\max,k} \cdot \sigma(Q_{s,k}) \quad (8)$$

where $\sigma(\xi) = 1/(1 + e^{-\xi})$ and

$$Q_{\max,k} = \begin{cases} 3, & \text{if } 0 < k < K - 1 \\ 1, & \text{otherwise} \end{cases} \quad (9)$$

This parameterization has one fewer degree of freedom than the other representations but is more interpretable.

3.4. Implementation and training

It is becoming increasingly well-known that differentiable IIR filters can be constructed exactly using recurrent layers [9, 11, 16]. If we consider the number of filters which we may want to cascade, and that such layers would operate on audio samples at $f_{\text{SR}} \geq 44.1$ kHz, it is apparent that this would be prohibitively slow to train in practice. Alternatively, we propose training differentiable biquads by implementing the cascaded `freqz` computation in Eq. 3 and `p2c` procedures for each representation with differentiable operators. Setting ω to match the frequency axis of an underlying N -point fast Fourier transform (FFT) allows us to perform filtering in the frequency domain. If N is made large enough, the amount of time-aliasing incurred via frequency domain sampling of the system response is small enough that it sufficiently approximates the true IIR response (at least for training purposes). This methodology was initially considered in [16] for the task of parametric EQ matching, and we have now extended it to the more general problem of (time-domain) audio effect modeling. A short-time Fourier transform, to this end, allows us to train on arbitrary length signals. During inference, filtering stage parameters are inferred from input parameter settings, and the resulting filters are implemented recursively.

3.5. Model complexity and interpretability

The number of parameters for the models proposed here is $2 + \sum_{s=0}^{S-1} (1 + KP)(1 + C_s)$. While neural networks are often criticized for their lack of interpretability, our models are transparent and intuitive. Particularly with the parametric EQ representation, each stage of the model can be hand-tuned by a machine learning non-specialist after training. Note that the biases of HC affine transformations correspond to the “quiescent” state of a stage with user inputs set to 0. Therefore, by altering the biases of hyperconditioned filtering stages or the parameters of fixed ones, we can globally scale gains and resonances, and/or shift frequency responses to the left/right. The “action” of external controls can be altered by adjusting the proper elements of the learned HC matrices.

s	C_s	Parameter name	Model	Params.	MSE
2	1	DIST	Coefficient	274	0.1708
6	2	LOW, HIGH	Pole/zero	274	0.0885
7	2	MID, MID FREQ	Param. EQ	210	0.0629
8	1	LEVEL	WaveNet	22960	0.0088

Table 1. Conditioning sites.

Table 2. Model comparisons.

Reference	98.6	H
Param. EQ	71.6	H
WaveNet	83.2	H
Anchor	17.7	H

Fig. 3. MUSHRA scores with 95% confidence intervals.

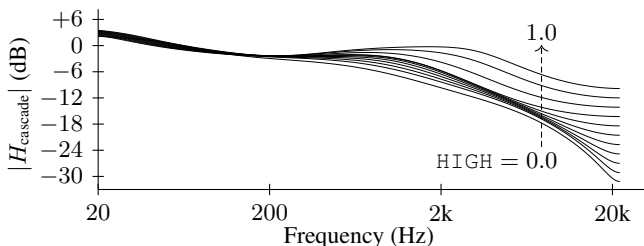


Fig. 4. MT-2 magnitude responses ($s = 6$), sweeping HIGH.

4. EXPERIMENTAL RESULTS

We seek to emulate the BOSS MT-2 distortion pedal, which has $S = 10$ cascaded filtering stages and $C_{\max} = 6$ external parameters. The parameters and their conditioning sites in the circuit are summarized in Tab. 1. We use an internally collected dataset of direct input guitar signals operating at $f_{\text{SR}} = 44.1$ kHz, segmenting them into 1 s clips. Distorted audio is generated using a SPICE model of the analog circuit. A large dataset is created by applying the effect with several random parameter settings to each clip.

In addition to the proposed models, we train a feedforward variant of WaveNet [7], serving as a baseline black-box model. The 1st layer of the model is a 1×1 convolutional layer outputting 16 channels, followed by 10 dilated convolutional layers, with dilation rate increasing by a factor of 2 at each layer. The number of residual, skip, and dilation channels per layer is 16. Dilation channels are generated using causal convolutions with a kernel size of 3 and standard gated-filter activation [20]. External parameter settings comprise a conditioning vector of length C_{\max} , which are injected into the gate and filter of each layer by means of 1×1 convolutions and addition. Residual and skip channels are generated using 1×1 convolutions. Skip channels are mixed using a 1×1 convolution with linear activation, yielding the final output. We train models for 100 epochs with a batch size of 16 and ADAM optimizer. We consider all representations discussed here with $K = 4$, $N = 4096$, and HC as in Tab. 1. As in [9], we use a time-domain mean squared error (MSE) loss, noting that pre-emphasis could improve perceptual relevance [23].

We report the loss evaluated on a held-out test set alongside parameter counts for all models in Tab. 2. While the

baseline model does outperform the proposed models in terms of MSE, our models perform competitively given that they have on the order of 100 times fewer parameters and added interpretability. Interestingly, we found that the parametric EQ representation actually outperformed the other representations mentioned here despite its lower parameter count. It is also the most interpretable and user-tunable representation introduced here.

We performed a Multiple Stimuli with Hidden Reference and Anchor (MUSHRA) [24] listening test to evaluate models subjectively. For each of 20 trials, participants were presented with test items and asked to rank how accurately each item sounded like the SPICE reference (0–100 scale). We used the following test items for each trial: the SPICE reference itself, the input guitar signal, peak-normalized to 6.0, and processed through a tanh NL (serving as an anchor), the output of the WaveNet baseline, and the output of our proposed method with parametric EQ representation. We had a total of 28 participants involved in the music technology field. Due to COVID-19 restrictions, participants conducted tests in their home environments with their own choice of headphones.

The results of our listening tests are shown in Fig. 3. We observe that all pairwise comparisons are statistically significant with 95% confidence, and that our approach performs competitively, albeit just slightly worse than the baseline. Informally, we observe that the emulation of the proposed method is still quite good, especially given its lower complexity, and that the WaveNet model often introduces unnatural resonances around 8 and 10 kHz. Lastly, Fig. 4 illustrates that HC can alter filter frequency responses in intuitive ways. For example, increasing the HIGH parameter increases the response of stage 6 in the high end. For audio examples and additional figures, please visit <https://sites.google.com/izotope.com/icassp2021-audio-demo>.

5. CONCLUSIONS

We proposed modeling an audio effect using hyperconditioned differentiable biquads. Several biquad representations were considered, and an FFT-based approximation enabled their efficient training. Our models are lightweight and interpretable, and perform comparably to a black-box approach. Future work will extend the use of our trainable IIR filtering methodologies to different audio modeling tasks.

6. REFERENCES

- [1] A. Fettweis, “Wave digital filters: Theory and practice,” *Proc. IEEE*, vol. 74, no. 2, pp. 270–327, Feb. 1986.
- [2] K. J. Werner, *Virtual analog modeling of audio circuitry using wave digital filters*, Ph.D. diss., Stanford Univ., 2016.
- [3] K. J. Werner, A. Bernardini, A. Sarti, and J. O. Smith III, “Modeling circuits with arbitrary topologies and active linear multiports using wave digital filters,” *IEEE Trans. Circuits Syst. I—Reg. Papers*, vol. 65, no. 12, pp. 4233–4246, Dec. 2018.
- [4] E. Damskögg, L. Juvela, E. Thuillier, and V. Välimäki, “Deep learning for tube amplifier emulation,” in *Proc. IEEE Int. Conf. Acoustic. Speech, Signal Process. (ICASSP)*, Brighton, UK, May 2019, pp. 471–475.
- [5] M. A. Martínez Ramírez and J. D. Reiss, “Modeling nonlinear audio effects with end-to-end deep neural networks,” in *Proc. IEEE Int. Conf. Acoustic. Speech, Signal Process. (ICASSP)*, Brighton, UK, May 2019, pp. 171–175.
- [6] M. A. Martínez Ramírez, E. Benetos, and J. D. Reiss, “A general-purpose deep learning approach to model time-varying audio effects,” in *Proc. Int. Conf. Digital Audio Effects (DAFx-19)*, Birmingham, UK, Sept. 2019, pp. 189–196.
- [7] A. Wright, E. Damskögg, L. Juvela, E., and V. Välimäki, “Real-time guitar amplifier emulation with deep learning,” *Appl. Sci.*, vol. 10, no. 3, 2020, Article #766.
- [8] J. Engel, L. Hantrakul, C. Gu, and A. Roberts, “DDSP: Differentiable digital signal processing,” in *Proc. Int. Conf. Learning Representations (ICLR)*, Addis Ababa, Ethiopia, Apr. 2020, pp. 26–30.
- [9] B. Kuznetsov, J. Parker, and F. Esqueda, “Differentiable IIR filters for machine learning applications,” in *Proc. Int. Conf. Digital Audio Effects (eDAFx-20)*, Vienna, Austria, Sept. 2020, pp. 297–303.
- [10] B. Kuznetsov, “Personal communication,” Sept. 2020.
- [11] É. Bavu, A. Ramamonjy, H. Pujol, and A. Garcia, “Timescalenet : A multiresolution approach for raw audio recognition,” in *Proc. IEEE Int. Conf. Acoustic. Speech, Signal Process. (ICASSP)*, Brighton, UK, May 2019, pp. 5686–5690.
- [12] G. Oh and J. S. Valois, “HCNAF: Hyper-conditioned neural autoregressive flow and its application for probabilistic occupancy map forecasting,” in *Proc. IEEE/CVF Conf. Comput. Vision Pattern Recognition (CVPR)*, June 2020, pp. 14538–14547.
- [13] J. Serrà, S. Pascual, and C. Segura, “Blow: A single-scale hyperconditioned flow for non-parallel raw-audio voice conversion,” in *Proc. Conf. Neural Information Process. Syst. (NeurIPS)*, Vancouver, Canada, Dec. 2019.
- [14] S. K. Mitra and J. F. Kaiser, Eds., *Handbook for Digital Signal Processing*, J. Wiley & Sons, New York, 1993.
- [15] R. Bristow-Johnson, “RBJ audio EQ cookbook,” Available at <https://www.musicdsp.org/en/latest/Filters/197-rbj-audio-eq-cookbook.html>, accessed October 01, 2020.
- [16] S. Nercessian, “Neural parametric equalizer matching using differentiable biquads,” in *Proc. Int. Conf. Digital Audio Effects (eDAFx-20)*, Vienna, Austria, Sept. 2020, pp. 265–272.
- [17] P. Bhattacharya, P. Nowak, and U. Zölzer, “Optimization of cascaded parametric peak and shelving filters with backpropagation algorithm,” in *Proc. Int. Conf. Digital Audio Effects (eDAFx-20)*, Vienna, Austria, Sept. 2020, pp. 101–108.
- [18] J. S. Abel and D. P. Berners, “Filter design using second order peaking and shelving sections,” in *Proc. Int. Comput. Music Conf. (ICMC)*, Miami, FL, Nov. 2004.
- [19] J. O. Smith III, *Physical Audio Signal Processing*, W3K Publ., 2010.
- [20] A. van den Oord et al., “WaveNet: A generative model for raw audio,” *arXiv:1609.03499*, 2016.
- [21] K. Qian, Y. Zhang, S. Chang, X. Yang, and M. Hasegawa-Johnson, “AutoVC: Zero-shot voice style transfer with only autoencoder loss,” in *Proc. Int. Conf. on Machine Learning (ICML)*, Long Beach, CA, June 2019.
- [22] X. Huang and S. Belongie, “Arbitrary style transfer in real-time with adaptive instance normalization,” in *Proc. IEEE Int. Conf. Comput. Vision (ICCV)*, Venice, Italy, Oct. 2017, pp. 1510–1519.
- [23] A. Wright and V. Välimäki, “Perceptual loss function for neural modeling of audio systems,” in *Proc. IEEE Int. Conf. Acoustic. Speech, Signal Process. (ICASSP)*, Barcelona, Spain, May 2020, pp. 251–255.
- [24] *ITU-R recommendation BS. 1534-3: Method for the subjective assessment of intermediate quality level of audio*, International Telecommunication Union, Oct. 2015.