

# DIFFERENTIABLE ALL-POLE FILTERS FOR TIME-VARYING AUDIO SYSTEMS

Chin-Yun Yu<sup>b\*</sup>, Christopher Mitcheltree<sup>b\*</sup>, Alistair Carson<sup>‡</sup>, Stefan Bilbao<sup>‡</sup>, Joshua D. Reiss<sup>b</sup>, and György Fazekas<sup>b</sup>

<sup>b</sup>Centre for Digital Music, Queen Mary University of London, London, UK

<sup>‡</sup>Acoustics and Audio Group, University of Edinburgh, Edinburgh, UK

{ chin-yun.yu, c.mitcheltree }@qmul.ac.uk, alistair.carson@ed.ac.uk

## ABSTRACT

Infinite impulse response filters are an essential building block of many time-varying audio systems, such as audio effects and synthesisers. However, their recursive structure impedes end-to-end training of these systems using automatic differentiation. Although non-recursive filter approximations like frequency sampling and frame-based processing have been proposed and widely used in previous works, they cannot accurately reflect the gradient of the original system. We alleviate this difficulty by re-expressing a time-varying all-pole filter to backpropagate the gradients through itself, so the filter implementation is not bound to the technical limitations of automatic differentiation frameworks. This implementation can be employed within audio systems containing filters with poles for efficient gradient evaluation. We demonstrate its training efficiency and expressive capabilities for modelling real-world dynamic audio systems on a phaser, time-varying subtractive synthesiser, and feed-forward compressor. We make our code and audio samples available and provide the trained audio effect and synth models in a VST plugin<sup>1</sup>.

## 1. INTRODUCTION

Infinite impulse response (IIR) filters are commonly used in many time-varying audio processing units, such as subtractive synthesisers, phaser effects, and dynamic range compression. Their recursive computation, using results from previous time steps, allows modelling a wide range of responses with low computational costs. Since differentiable DSP (DDSP) [1] emerged as an effective solution in attaining controllable audio systems, there have been attempts to incorporate recursive filters into automatic differentiation frameworks such as PyTorch. However, naive implementations result in deep computational graphs due to recursion that cannot be parallelised [2] and slow down training speed [3, 4, 5].

A common acceleration approach is to evaluate the filters in the frequency domain [3] and approximate time-varying behaviour by filtering on overlapping frames in parallel [6, 7]. Despite their popularity, these approximations have some potential drawbacks. Frame windowing for overlap-add smears the spectral peaks in the frequency domain, resulting in filters with artificially high resonance to counter this effect. Sampling the filters in the frequency domain sometimes truncates the IR length, and the circular convolution with the truncated IR caused by frequency sampling can lead to artefacts at frame boundaries. Most importantly, systems

trained in this way are not guaranteed the same results when operating sample-by-sample at audio rate to achieve low latency.

In this paper, we propose a solution to these problems by deriving and implementing an efficient backpropagation algorithm for a time-varying all-pole filter. Our filter can be used in various audio systems by separating the system's poles and zeros and explicitly handling the poles (where the recursion is) with our proposed method fully end-to-end. Our contributions are threefold:

1. We significantly increase the forward and backpropagation speed of time-varying recursive all-pole filters without introducing any approximation to the filter.
2. The systems trained with our implementation can be converted to real-time without generalisation issues besides the order of the zeros and poles.
3. We show that our filter efficiently and accurately models time-varying analog audio circuits with recursive structures.

## 2. RELATED WORKS

Differentiable training of IIR filters has been explored in [8, 9, 10, 3, 11, 4, 7]. To sidestep the problems inherent in training over a recursion, some authors approximate IIR filters in the frequency domain using the fast Fourier transform (FFT) [9, 10, 3, 11, 4], or the short-time Fourier transform (STFT) [7] for time-varying effects. This is known as the frequency sampling (FS) method. It approximates the filter as time-invariant and with a finite impulse response (FIR) over the duration of a short frame, thus the accuracy of FS depends heavily on the choice of STFT parameters: the hop-size, the frame length, the FFT length, and the windowing function [12]. In machine learning applications, these choices add extra hyper-parameters to models, which may require prior knowledge of the target system to set appropriately.

Time-varying all-pole filters have been used for decades in linear prediction (LP) voice synthesis [13]. Training them jointly with neural networks was first proposed in LPCNet [14]. The authors achieve training efficiency by using inverse-filtered speech as the target because the inverse filter has no recursion. Other works seek to parallelise LP with frame-based processing, either filtering in the time domain [15, 2] or via FS [6, 16] for each frame.

Dynamic processing effects like compressors and limiters also employ time-varying recursive filters. The filter is usually first order, and the coefficients are time-varying and dependent on the *attack* or *release* phases of the gain reduction signal [17]. In the differentiable learning context, frequency sampling can be used if the compressor's attack and release time are configured to be the same [4, 5], which simplifies the filter to be time-invariant. Colonel et al. [18] propose dividing the gain reduction signal into attack and release passages and filtering them separately with different filters, and Guo et al. [19] downsample the signal to reduce the number of recursions.

\* Equal contribution.

<sup>1</sup><https://diffapf.github.io/web/>

A distinct approach that provides significant acceleration is deriving the closed-form solution of the gradients for the filter parameters and implementing it in a highly optimised way. Bhattacharya et al. [8] derive the instantaneous backpropagation algorithm for peak and shelving filters and train them jointly to match the response of head-related transfer functions. Forgione et al. and Yu et al. [2, 20] decompose a time-invariant IIR filter into zeros and poles, and they show that backpropagation through an all-pole filter can be expressed using the same all-pole filter. This method is fast because the underlying filters do not have to be implemented in an automatic differentiation framework. Nevertheless, these solutions are made for specific filters or are only applied to time-invariant systems.

### 3. PROPOSED METHODOLOGY

Consider an  $M^{\text{th}}$ -order time-varying all-pole filter:

$$\begin{aligned} y(n) &= f_{\mathbf{a}(n)}(x(n)) \\ &= x(n) - \sum_{i=1}^M a_i(n)y(n-i) \end{aligned} \quad (1)$$

where  $\mathbf{a}(n) = [a_1(n), \dots, a_M(n)]$  are filter coefficients at time  $n$  and  $M \in \mathbb{Z}^+$ . In some applications,  $\mathbf{a}(n)$  varies at a control rate  $F_c$  much lower than the audio sampling rate  $F_s$  as  $\mathbf{a}(m)$ , and can then be up-sampled to the audio rate before the filter is applied.

The following sections describe the proposed method, in which the exact gradients for each parameter of filter  $f_{\mathbf{a}(n)}$  are derived and expressed in a form that can be computed efficiently. Our method aligns most with the instantaneous backpropagation algorithms proposed in [8, 20, 2], generalising their contributions to a time-varying all-pole filter that can be used in various recursive filters and time-varying audio systems. We refer to this method as the time domain (TD) method.

#### 3.1. Unwinding the Recursion

We first rewrite the recursive Eq. (1) so there is no  $y$  variable on the right-hand side:

$$y(n) = x(n) + \sum_{d=1}^{\infty} b_d(n)x(n-d). \quad (2)$$

$\mathbf{b}(n) = [b_1(n), b_2(n), \dots]$  is the IIR of filter  $f_{\mathbf{a}(n)}$  at time  $n$ . We will use Eq. (2) and  $\mathbf{b}(n)$  to help us derive the gradient of  $\mathbf{a}(n)$  in the next two sections. To get the exact value of  $\mathbf{b}(n)$  in terms of  $\mathbf{a}(n)$ , let us think of (1) as the recursive definition of connecting any time step before  $n$  to  $n$ . In other words,  $y(n)$  is a summation of  $M$  sub-problems  $y(n-i)$ , each weighted by a unique coefficient  $-a_i(n)$ . Let us think of  $i$  as the step size going backwards from  $n$  with  $M$  being the maximum step size we can take. All possible combinations of steps that span a distance  $d \in \mathbb{Z}^+$  can be defined recursively as

$$\mathcal{G}_d = \bigcup_{i=1}^{\min(d, M)} \{[i; \mathbf{q}] : \mathbf{q} \in \mathcal{G}_{d-i}\} \quad (3)$$

with boundary condition  $\mathcal{G}_0 = \{[]\}$  and  $[;]$  is array concatenation.

If we compare Eq. (1) and (2), we can see that  $b_d(n)$  is in fact the sum of the coefficient combinations of  $a_{\cdot}(\cdot)$  going back from

$n$  with step combinations  $\mathcal{G}_d$ . Thus, the value of it is

$$b_d(n) = \sum_{\mathbf{q} \in \mathcal{G}_d} (-1)^{|\mathbf{q}|} \prod_{j=1}^{|\mathbf{q}|} a_{q_j} \left( n - \sum_{k=1}^j Q(\mathbf{q})_k \right) \quad (4)$$

where  $|\mathbf{q}|$  is the number of elements in array  $\mathbf{q}$  and  $Q(\mathbf{q}) = [0; \mathbf{q}]$ .

#### 3.2. Gradients of $x(n)$

Assuming we have computed  $f_{\mathbf{a}(n)}$  up to step  $N$ , evaluated  $y(n \leq N)$  with a differentiable function  $\mathcal{L}(y(n))$ , and have its instantaneous gradients  $\frac{\partial \mathcal{L}}{\partial y(n)}$ , we can backpropagate through  $f_{\mathbf{a}(n)}$  as

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial x(n)} &= \sum_{d=-\infty}^N \frac{\partial \mathcal{L}}{\partial y(d)} \frac{\partial y(d)}{\partial x(n)} \\ &= \frac{\partial \mathcal{L}}{\partial y(n)} \frac{\partial y(n)}{\partial x(n)} + \sum_{d=1}^{N-n} \frac{\partial \mathcal{L}}{\partial y(n+d)} \frac{\partial y(n+d)}{\partial x(n)} \quad (5) \\ &= \frac{\partial \mathcal{L}}{\partial y(n)} + \sum_{d=1}^{N-n} b_d(n+d) \frac{\partial \mathcal{L}}{\partial y(n+d)}. \end{aligned}$$

We use the fact that  $\frac{\partial y(n)}{\partial x(n-d)} = b_d(n)$  and  $\frac{\partial y(<n)}{\partial x(n)} = 0$  from (2). Unfortunately, Eq. (4) and therefore (5) are expensive to compute. We aim to express backpropagation using  $f_{\mathbf{a}(n)}$ , which is much more efficient to compute due to its recursion. If we can reparameterise Eq. (5) to look like (1), then  $f_{\mathbf{a}(n)}$  can be reused to compute  $\frac{\partial \mathcal{L}}{\partial x(n)}$ . We do this by writing  $b_d(n+d)$  in terms of  $a_{\cdot}(\cdot)$  which requires us to evaluate (4) at  $n+d$  to get

$$\begin{aligned} b_d(n+d) &= \sum_{\mathbf{q} \in \mathcal{G}_d} (-1)^{|\mathbf{q}|} \prod_{j=1}^{|\mathbf{q}|} a_{q_j} \left( n+d - \sum_{k=1}^j Q(\mathbf{q})_k \right) \\ &= \sum_{\mathbf{q} \in \mathcal{G}_d} (-1)^{|\mathbf{q}|} \prod_{j=1}^{|\mathbf{q}|} a_{q_j} \left( n + \sum_{k=j}^{|\mathbf{q}|} q_k \right) \quad (6) \\ &= \sum_{\mathbf{q} \in \mathcal{G}_d} (-1)^{|\mathbf{q}|} \hat{a}_{q_{|\mathbf{q}|}}(n) \prod_{j=1}^{|\mathbf{q}|-1} \hat{a}_{q_j} \left( n + \sum_{k=j+1}^{|\mathbf{q}|} q_k \right) \\ &= \sum_{\tilde{\mathbf{q}} \in \mathcal{G}_d} (-1)^{|\tilde{\mathbf{q}}|} \prod_{j=1}^{|\tilde{\mathbf{q}}|} \hat{a}_{\tilde{q}_j} \left( n + \sum_{k=1}^j Q(\tilde{\mathbf{q}})_k \right) \end{aligned}$$

where  $\tilde{\mathbf{q}} = [q_{|\mathbf{q}|}, \dots, q_1] \in \mathcal{G}_d$  and  $\hat{a}_i(n) = a_i(n+i)$ . Eq. (6) is now the same as (4) but uses  $\hat{\mathbf{a}}(n) = [\hat{a}_1(n), \dots, \hat{a}_M(n)]$  as coefficients and the plus sign inside the product changes to a minus sign, which means the filter should be applied in the reverse direction ( $n = N \rightarrow n = -\infty$ ) and is supported by the non-causal indexing in (5) ( $n+d$  instead of  $n-d$ ). We can now express (5) in terms of  $f_{\mathbf{a}(n)}$  using (6) and the equivalence between  $\mathbf{a}(n)$  and  $\mathbf{b}(n)$  from (1) and (2) to get

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial x(n)} &= \frac{\partial \mathcal{L}}{\partial y(n)} - \sum_{i=1}^M \hat{a}_i(n) \frac{\partial \mathcal{L}}{\partial x(n+i)} \\ &= \text{FLIP} \circ f_{\text{FLIP} \circ \hat{\mathbf{a}}(n)} \circ \text{FLIP} \circ \frac{\partial \mathcal{L}}{\partial y(n)} \quad (7) \end{aligned}$$

where  $\text{FLIP}(x(n)) = x(-n)$  and  $f_1 \circ f_2(x) = f_1(f_2(x))$ .  $\text{FLIP}$  and  $\hat{\mathbf{a}}(n)$  are trivial to compute using memory indexing. The backpropagation algorithm and how to arrange  $\mathbf{a}(n)$  into  $\hat{\mathbf{a}}(n)$  is shown in Figure 1.

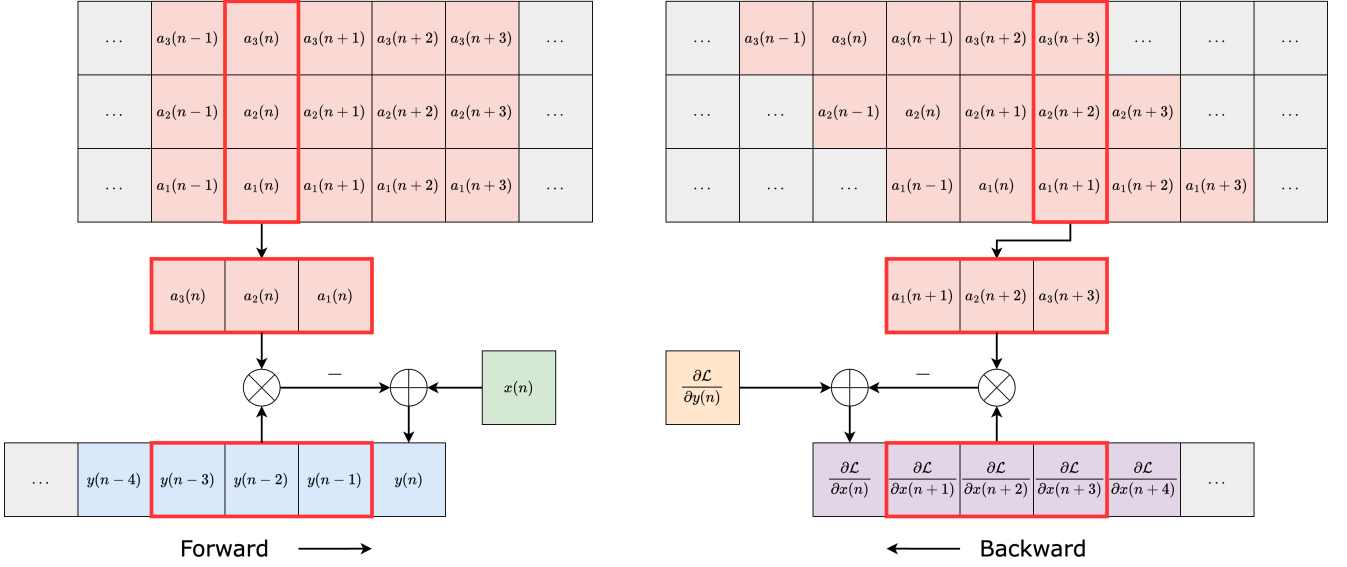


Figure 1: The forward (left) and backpropagation (right) flow chart of a third-order time-varying all-pole filter.

### 3.3. Gradients of $\mathbf{a}(n)$

Let  $u_i(n) = -a_i(n)y(n-i)$  so  $y(n) = x(n) + u_1(n) + \dots + u_M(n)$ . Because of the chain rule and  $\frac{\partial y(n)}{\partial x(n)} = \frac{\partial y(n)}{\partial u_i(n)} = 1$ ,  $x(n)$  and  $u_i(n)$  should have the same derivatives ( $\frac{\partial \mathcal{L}}{\partial x(n)} = \frac{\partial \mathcal{L}}{\partial u_i(n)}$ ). Since we can compute  $\frac{\partial \mathcal{L}}{\partial x(n)}$  from (7), the gradients of the coefficients are simply

$$\frac{\partial \mathcal{L}}{\partial a_i(n)} = \frac{\partial \mathcal{L}}{\partial u_i(n)} \frac{\partial u_i(n)}{\partial a_i(n)} = -\frac{\partial \mathcal{L}}{\partial x(n)} y(n-i). \quad (8)$$

In summary, we can calculate all gradients with one pass of  $f_{\mathbf{a}(n)}$  and the multiplications in (8), which are fast to compute. We implement an efficient  $f_{\mathbf{a}(n)}$  for both forward and backward computation using Numba and register it as a custom operator in PyTorch. The implementation is available on GitHub<sup>2</sup>.

## 4. APPLICATIONS

We demonstrate our all-pole filter implementation on three commonly used dynamic audio systems: a phaser, a subtractive synthesiser, and a compressor. All three systems have time-varying recursive structures that are not easy to train in a differentiable way and would typically be modelled using FS approaches.

Although the filtering order of poles and zeros matters in time-varying systems, in this work, we rearranged the poles of each system into one all-pole filter for maximum training efficiency. Due to the relatively slowly varying filter coefficients, we found this approach to be sufficient. We direct interested readers to our parallel work in speech synthesis [21], an exact all-pole system.

### 4.1. Phaser

We test our filter implementation on a virtual analog phaser model, based on [22, 7]. At the core of the model is a differentiable LFO that operates at the control rate  $F_c$ . The oscillator is implemented

as a damped oscillator with learnable frequency  $f_0$ , decay rate  $\sigma$ , and phase  $\phi$ :

$$s(m) = e^{-\sigma^2 m / F_c} \cos(2\pi f_0 m / F_c + \phi) \quad (9)$$

where  $m$  is the control-rate sample index. The inclusion of parameter  $\sigma$  alleviates some non-convexity issues when learning the frequency  $f_0$ , as shown in [23]. Note that the oscillator is unconditionally stable for all  $\sigma$ . The oscillator is passed through a multi-layer perceptron (MLP) network to obtain the control signal  $p(m)$ . The MLP, with parameters  $\Theta$ , contains 3x8 hidden layers, with tanh activation functions on all layers including the final. The control signal is then up-sampled with linear interpolation to obtain  $p(n)$  and modulates the coefficients of four cascaded first-order all-pass filters (APF), each with the difference equation:

$$y_k(n) = p(n) \cdot [x_k(n) + y_k(n-1)] - x_k(n-1) \quad (10)$$

where  $x_k$  and  $y_k$  are the input and state of the  $k^{\text{th}}$  APF, respectively,  $0 \leq k < 4$ . Note that the MLP tanh output activation ensures the all-pass poles remain within the unit circle. The APFs are arranged in series, with a through path of gain  $g_1$  and feedback loop  $g_2$  as shown in Figure 2. It is common to include a unit delay in the feedback path for ease of implementation [22, 24], however here we use instantaneous feedback for a more realistic virtual analog model. In Figure 2, BQ represents a biquad filter with coefficients  $\mathbf{b}^{(\text{bq})} = [b_0^{(\text{bq})}, b_1^{(\text{bq})}, b_2^{(\text{bq})}]$  and  $\mathbf{a}^{(\text{bq})} = [a_1^{(\text{bq})}, a_2^{(\text{bq})}]$ . The entire model is a sixth-order time-varying IIR filter. Here we approximate the system as having the difference equation:

$$y(n) = f_{\mathbf{a}(n)} \circ f_{\mathbf{b}(n)} \circ x(n) \quad (11)$$

where  $f_{\mathbf{b}(n)}(\cdot)$  is a time-varying FIR filter:

$$f_{\mathbf{b}(n)}(x(n)) = \sum_{i=0}^M b_i(n)x(n-i) \quad (12)$$

and  $f_{\mathbf{a}(n)}(\cdot)$  is a time-varying all-pole filter (see Eq. (1)). Here  $M = 6$  for  $\mathbf{b}(n)$  and  $\mathbf{a}(n)$ , which are functions of  $p(n)$ ,  $g_1$ ,  $g_2$ ,

<sup>2</sup><https://github.com/yoyololicon/torchlpc>

$\mathbf{b}^{(bq)}$ , and  $\mathbf{a}^{(bq)}$ . In previous work [7], the control parameters of a similar time-varying filter were learned through gradient descent using the FS method. This frequency sampling approach had some limitations, however. Firstly, the optimal frame size for the best training accuracy depended on the rate of the target LFO, which we ideally should not assume as prior knowledge. Secondly, it was not fully investigated whether the trained model could then be implemented in the time domain at inference to avoid latency.

Here, we instead implement Eq. (11) directly in the time domain during training using the method proposed in Section 3.

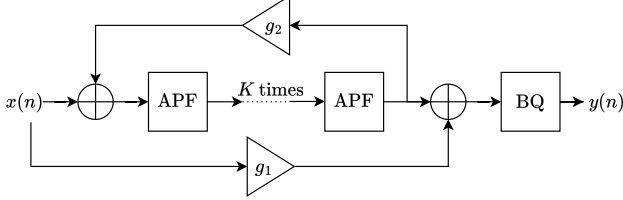


Figure 2: Discrete-time phaser model considered in this work, where  $K = 4$ . APF represents a time-varying all-pass filter with difference equation (10) and BQ is a biquad filter.

#### 4.2. Time-varying Subtractive synthesiser

We test our filter implementation on a subtractive synthesiser roughly modelled after the *Roland TB-303 Bass Line* synth<sup>3</sup> which defined the acid house electronic music movement of the late 1980s. The TB-303 is an ideal synth for our use case because its defining feature is a resonant low-pass filter where the cutoff frequency is modulated quickly using an envelope to create its signature squelchy, “liquid” sound. Although the original TB-303’s circuit contains a 4-pole diode ladder filter, for simplicity and demonstration purposes, we implement our synthesiser using a biquad filter. Our synth is differentiable and consists of three main components: a monophonic oscillator, a time-varying biquad filter, and a waveshaper for adding distortion to the output.

The oscillator is the same as in the one in TorchSynth [25] and uses hyperbolic tangent waveshaping to generate sawtooth or square waves, and can sweep continuously between them. It is defined by the following equations:

$$\psi(n) = 2\pi n f_0 / F_s + \phi \pmod{2\pi} \quad (13)$$

$$o(n) = \rho_{osc} s_{saw}(\psi(n)) + (1 - \rho_{osc}) s_{sq}(\psi(n)) \quad (14)$$

$$e(n) = \begin{cases} \left( \frac{N_{on} - n}{N_{on}} \right)^{\rho_{env}} & 0 \leq n \leq N_{on} \\ 0 & \text{otherwise} \end{cases} \quad (15)$$

$$s(n) = g_{osc} e(n) o(n) \quad (16)$$

where  $F_s$  and  $f_0$  are the sampling rate and fundamental frequency in hertz, and  $\phi$  is the phase in radians.  $\rho_{osc}$  is a continuous control parameter to sweep between the wave shapes where 0 makes a square wave ( $s_{sq}(\cdot)$ ), and 1 makes a saw wave ( $s_{saw}(\cdot)$ ). The output audio of the oscillator is multiplied by gain  $g_{osc}$  and is then shaped using a decaying envelope  $e(n)$  of length  $N_{on}$  note on samples with control parameter  $\rho_{env}$ .

A time-varying biquad filter  $h_{bq}(\cdot)$  (same as Equations 11 and 12 for the phaser, but  $M = 2$ ) is then applied to the oscillator output audio. This filter takes as input 5 time-varying filter coefficients  $\{a_1(n), a_2(n), b_0(n), b_1(n), b_2(n)\}$  at sample rate which

can be passed in directly or generated from filter cutoff and resonance modulation signals  $m_{fc}(n)$  and  $m_q(n)$ . These modulation signals are then used to calculate the coefficients for a biquad lowpass filter using the corresponding equations in the Audio EQ Cookbook<sup>4</sup>.

Finally, the output of the filter is fed through a hyperbolic tangent waveshaper which adds distortion:

$$f_{dist}(x(n)) = \tanh(g_{dist} x(n)). \quad (17)$$

The amount of distortion is controlled by parameter  $g_{dist}$  which modifies the gain of the input of the waveshaper  $x(n)$ .

The entire synth is therefore controllable using 8 global parameters  $\{f_0, F_s, \phi, N_{on}, \rho_{osc}, \rho_{env}, g_{osc}, g_{dist}\}$  and 2 or 5 time-varying parameters  $\{m_{fc}(n), m_q(n)\}$  or  $\{a_1(n), a_2(n), b_0(n), b_1(n), b_2(n)\}$ . It is defined by composing Equation 16,  $h_{bq}(\cdot)$ , and Equation 17 together as follows:

$$y_{synth}(n) = f_{dist} \circ h_{bq} \circ s(n). \quad (18)$$

#### 4.3. Feed-forward Compressor

The compressor we consider here is a simple feed-forward compressor from [17], which is defined as:

$$x_{rms}(n) = \alpha_{rms} x^2(n) + (1 - \alpha_{rms}) x_{rms}(n-1) \quad (19)$$

$$g(n) = \min \left( 1, \left( \frac{\sqrt{x_{rms}(n)}}{10^{\frac{CT}{20}}} \right)^{\frac{1-R}{R}} \right) \quad (20)$$

$$\hat{g}(n) = \begin{cases} \alpha_{at} g(n) + (1 - \alpha_{at}) \hat{g}(n-1) & g(n) < \hat{g}(n-1) \\ \alpha_{rt} g(n) + (1 - \alpha_{rt}) \hat{g}(n-1) & \text{otherwise} \end{cases} \quad (21)$$

$$y(n) = x(n) \hat{g}(n) \gamma. \quad (22)$$

$R$ ,  $CT$ , and  $\gamma$  are the ratio, threshold, and make-up gain, respectively.  $\alpha_{rms/at/rt}$  are the average smoothing coefficients.  $\alpha_{at/rt}$  are chosen based on whether the compressor is operated in the *attack* phase or *release* phase. This compressor’s training efficiency bottleneck is described in (21), as the coefficient of a recursive filter is computed *on the fly*, so we cannot use  $f_{a(n)}$  directly. Moreover, it operates at the audio rate, so it is unsuitable for frame-based approximation.

##### 4.3.1. Custom backward function

We backpropagate gradients through (21) using the proposed method in Sec. 3. To handle the if-else statement, we write the gain reduction filter (21) in Numba and record each if-else decision inside the recursion into a binary mask  $\zeta(n)$ :

$$\zeta(n) = \begin{cases} 1 & g(n) < \hat{g}(n-1) \\ 0 & \text{otherwise} \end{cases}. \quad (23)$$

Using this, Eq. (21) equals the following time-varying IIR:

$$\begin{aligned} \beta(n) &= \alpha_{at}^{\zeta(n)} \alpha_{rt}^{1-\zeta(n)} \\ \hat{g}(n) &= \beta(n) g(n) + (1 - \beta(n)) \hat{g}(n-1). \end{aligned} \quad (24)$$

<sup>3</sup><https://www.roland.com/uk/promos/303day/>

<sup>4</sup><https://www.w3.org/TR/audio-eq-cookbook/>

Table 1: Parameter limits in our differentiable DSP models.

	Phaser					Time-varying subtractive synth						Feed-forward compressor					
	$f_0$	$\sigma$	$\phi$	$g_1$	$g_2$	$\rho_{\text{osc}}$	$\rho_{\text{env}}$	$g_{\text{osc}}$	$g_{\text{dist}}$	$m_{\text{fc}}(n)$	$m_{\text{q}}(n)$	$R$	$CT$	$\alpha_{\text{at}}$	$\alpha_{\text{rt}}$	$\alpha_{\text{rms}}$	$\gamma$
min.	$-F_c/2$	$-\infty$	$-\pi$	$-\infty$	0	0.0	0.1	0.01	0.01	0.1 kHz	0.7071	1.0	$-\infty$	0.0	0.0	0.0	0.0
max.	$F_c/2$	$\infty$	$\pi$	$\infty$	1	1.0	6.0	1.00	4.00	8.0 kHz	8.0	$\infty$	$\infty$	1.0	1.0	1.0	$\infty$

We can now use  $\beta(n)$  and  $f_{\text{a}(n)}$  to backpropagate the gradients through the compressor. The pseudo-code of the differentiable gain reduction filter is summarised in Algorithm 1, and the implementation can be found on GitHub<sup>5</sup>.

**Algorithm 1:** Differentiable gain reduction filter (21).

```

def forward( $g, \alpha_{\text{at}}, \alpha_{\text{rt}}$ ):
     $\hat{g}(-1) \leftarrow 1$ 
     $\hat{g}(n) \leftarrow \text{Eq. (21)}$ 
     $\zeta(n) \leftarrow \text{Eq. (23)}$ 
    return  $\hat{g}, \zeta$ 

def backward( $\frac{\partial \mathcal{L}}{\partial \hat{g}}, \hat{g}, g, m, \alpha_{\text{at}}, \alpha_{\text{rt}}$ ):
     $\beta(n) \leftarrow \text{Eq. (24)}$ 
     $\frac{\partial \mathcal{L}}{\partial \beta(n)g(n)} \leftarrow \text{filtering } \frac{\partial \mathcal{L}}{\partial \hat{g}(n)}$  with Eq. (7) and
         $a_1(n) \equiv \beta(n) - 1$ 
     $\frac{\partial \mathcal{L}}{\partial g(n)} \leftarrow \frac{\partial \mathcal{L}}{\partial \beta(n)g(n)} \beta(n)$ 
     $\frac{\partial \mathcal{L}}{\partial \beta(n)} \leftarrow \frac{\partial \mathcal{L}}{\partial \beta(n)g(n)} (g(n) - \hat{g}(n-1))$ 
     $\frac{\partial \mathcal{L}}{\partial \alpha_{\text{at}}} \leftarrow \sum_n \frac{\partial \mathcal{L}}{\partial \beta(n)} \zeta(n)$ 
     $\frac{\partial \mathcal{L}}{\partial \alpha_{\text{rt}}} \leftarrow \sum_n \frac{\partial \mathcal{L}}{\partial \beta(n)} (1 - \zeta(n))$ 
    return  $\frac{\partial \mathcal{L}}{\partial g}, \frac{\partial \mathcal{L}}{\partial \alpha_{\text{at}}}, \frac{\partial \mathcal{L}}{\partial \alpha_{\text{rt}}}$ 
    
```

## 5. EXPERIMENTS

All three systems are trained to model some target analog audio in an end-to-end fashion using gradient descent. We use the Hanning window for all the frame-based approaches. The ranges of all the interpretable parameters are summarised in Table 1. We provide audio samples for all experiments on the accompanying website.

### 5.1. Modelling the EHX Small Stone analog phaser

Here we explore using the DSP phaser model outlined in Sec. 4.1 to model an analog phaser pedal: the Electro-Harmonix Small-Stone. This is the same system modelled in [7] using the FS method. The circuit consists of four cascaded analog all-pass filters, a through-path for the input signal, and a feedback path [26] – so topologically the circuit is similar to the discrete-time phaser model considered in this paper. The pedal consists of one knob which controls the LFO rate, and a switch that engages the feedback loop. Six different parameter configurations are considered:

- SS-A: feedback off, rate knob 3 o’clock ( $f_0 \approx 2.3$  Hz)
- SS-B: feedback off, rate knob 12 o’clock ( $f_0 \approx 0.6$  Hz)
- SS-C: feedback off, rate knob 9 o’clock ( $f_0 \approx 0.09$  Hz)
- SS-D: feedback on, rate knob 3 o’clock ( $f_0 \approx 1.4$  Hz)
- SS-E: feedback on, rate knob 12 o’clock ( $f_0 \approx 0.4$  Hz)
- SS-F: feedback on, rate knob 9 o’clock ( $f_0 \approx 0.06$  Hz)

Table 2: Phaser evaluation results. “Method” refers to the training method; whereas the ESR was computed over the test dataset using both FS and TD at inference.

Dataset	$L/F_s$	Method	ESR (%)	
			FS	TD
SS-A	10 ms	FS	1.46	1.53
		TD	<b>1.34</b>	<b>1.36</b>
SS-B	40 ms	FS	1.37	1.49
		TD	<b>1.35</b>	<b>1.34</b>
SS-C	160 ms	FS	<b>1.62</b>	<b>1.80</b>
		TD	2.56	2.23
SS-D	10 ms	FS	22.47	23.47
		TD	<b>21.64</b>	<b>23.33</b>
SS-E	40 ms	FS	15.43	16.69
		TD	<b>13.63</b>	<b>13.87</b>
SS-F	160 ms	FS	8.79	9.83
		TD	<b>7.83</b>	<b>8.79</b>

The training data consists of a 30 s chirp-train both dry (input) and processed through the pedal (target). At each training iteration, the input signal is processed through the model in a single batch, and the loss function is computed as the error-to-signal ratio (ESR) between the model output and target. The learnable model parameters are  $\{g_1, g_2, f_0, \sigma, \phi, \Theta, \mathbf{b}^{(\text{bq})}, \mathbf{a}^{(\text{bq})}\}$ , as defined in Sec. 4.1, giving a total of 182 model parameters. An Adam optimiser with a learning rate  $5 \times 10^{-4}$  is employed to carry out parameter updates every iteration for a maximum of 10k iterations. The test data includes the training data plus the next 10 seconds of the same audio signal, which contains guitar playing. This ensures the learned LFO phase is always aligned to the same point in time.

As reported in [7], the accuracy and convergence of model training depends on the choice of hop-size and window-length. Furthermore, even for a fixed choice of hyper-parameters, the training convergence depends on the initial parameter values, which are pseudo-randomly initialised [7]. We observed that for some random seeds, the LFO would not converge to the correct frequency  $f_0$  and/or decay rate  $\sigma$ . In a successful run, the learned  $f_0$  is approximately equal to that of the target, and  $\sigma$  converges to zero.

As a baseline, we train the model using the FS method with a single hop-size  $L$  for each parameter configuration. The window-length  $N_{\text{WIN}}$  is set to four times the hop-size, and the FFT length to  $2^{\lceil \log_2(N_{\text{WIN}}) \rceil + 1}$ . The training process is repeated up to five times with different seeds until a model converges. The proposed time domain (TD) implementation is then trained with the same hyper-parameters (where relevant) and initial seed as the FS method. Here, the hop-size determines the control rate.

To evaluate the two methods, we train the phaser model with the respective method and then test using both FS and TD at inference time. The test ESR can be seen in Table 2. For all datasets, it can be seen that both methods result in a very similar test loss, with the TD method doing slightly better in five out of the six datasets.

<sup>5</sup><https://github.com/yoyololicon/torchcomp>

Of course, only one hop-size has been considered here, so a more detailed analysis across a range of hop-sizes would be an interesting area of further work. However, these specific hop-sizes were chosen based on the recommendations in our previous work [7], in which they were heuristically found to give the best results (for the respective datasets) using the FS training method.

In early experiments we observed that the TD implementation could become unstable during training, causing the output signal to explode. This occurred even for a simplified problem of  $g_2 = 0$  and with the exclusion of the BQ filter. It was verified that the APF poles were within the unit circle for all  $n$ , albeit close in some cases ( $p(n) > 0.98$ ). We, therefore, suspect this instability was due to numerical inaccuracies associated with the transient response of the filters (e.g. as described in [12]) because changing from single-precision to double-precision resolved this problem. Therefore, we recommend operating at double precision when using the proposed filter implementation if instability arises.

## 5.2. Modelling the Roland TB-303 Acid Synth

We model analog TB-303 audio with our time-varying subtractive synth. The dataset is made from Sample Science’s royalty free *Abstract 303* sample pack<sup>6</sup> consisting of 100 synth loops at 120 BPM recorded dry from a hardware TB-303 clone. All loops are concatenated together, resampled to 48 kHz, and then pitch and note on durations are extracted using Ableton Live 11’s melody-to-midi conversion algorithm which we found to be highly accurate for these monophonic melody loops. Since the TB-303 is a 16 note sequencer, the resulting annotated notes are truncated or zero-padded to 6000 samples (one 16th note at 120 BPM and 48 kHz) with any notes shorter than 4000 samples in duration thrown out. This is then split into 60%, 20%, and 20% train, validation, and test sets, respectively, resulting in a total of 42.5 seconds of audio.

We use the same modulation extraction approach as [27] and DDSP [1] to model the synth. First, frame-by-frame features are extracted from the target audio and are processed by a neural network which predicts the temporal and global control parameters for our differentiable synth as defined in Section 4.2. Temporal parameters are linearly interpolated from frame-rate to sample-rate when required. The synth then generates some reconstructed audio from the control parameters which can be compared against the target audio using a differentiable loss function, thus enabling the system to be trained end-to-end using gradient descent. A diagram of the entire setup is shown in Figure 3.

Since the filter modulations of the TB-303 are very fast (around 125 ms for 16th notes at 120 BPM), we use a Mel spectrogram with 1024 FFT size, 128 Mel bins, and a short hop length of 32 samples which results in 188 frames for 6000 samples of audio. The neural network is the same architecture as LFO-net [27] except with 4 or 5 additional 2-layer MLP networks applied to the latent embedding averaged across the temporal axis to predict the global parameters of the synth. It contains 730 K parameters. We conduct experiments with five different synth filter configurations:

1. Time domain biquad coefficients (Coeff TD)
2. Frequency sampling biquad coefficients (Coeff FS)
3. Time domain low-pass biquad (LP TD)
4. Frequency sampling low-pass biquad (LP FS)
5. Time domain recurrent neural network (LSTM)

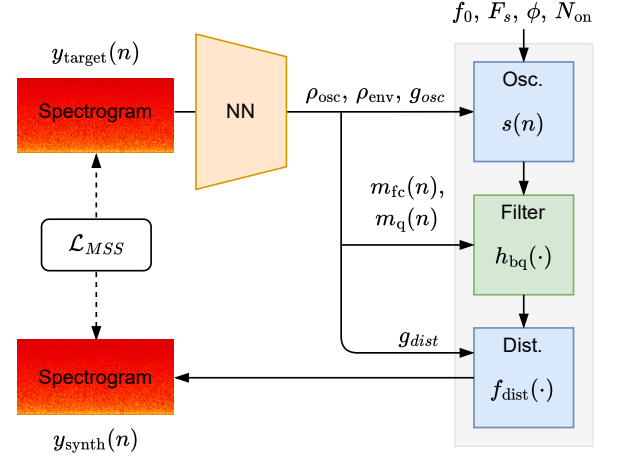


Figure 3: Diagram of the differentiable synth modelling process. Our time domain filter component is shown in green.

As discussed in Section 4.2, for configurations 1 and 2 the neural network outputs a 5-dimensional temporal control parameter of changing biquad coefficients. Before filtering, these coefficients are post-processed using the biquad triangle parameterisation of [3] to ensure stability. For configurations 3 and 4, a 2-dimensional modulation signal is returned, representing a changing filter cutoff and its Q factor (which is constant over time). These are then converted to five biquad coefficients. The raw coefficient filter configuration gives the synth as many degrees of freedom as possible whereas the lowpass filter configuration is based on the TB-303’s analog design consisting of an envelope modulated lowpass filter with a global resonance control knob. Finally, the recurrent neural network filter (configuration 5) is based on the architecture in [5, 27] and enables us to compare against learning a time-varying IIR filter directly from scratch. It is conditioned with the same 2-dimensional modulation signal as configurations 3 and 4.

We train all models for 200 epochs on batches of 34 notes using the AdamW optimiser and single precision – the numerical issues noted in Sec. 5.1 were not observed here. For the FS synth filter configurations we train separate models for  $N_{WIN} \in [128, 256, 512, 1024, 2048, 4096]$  while keeping the hop-size  $L$  fixed at 32 samples to match the frame-rate temporal outputs of the modulation extraction neural network. As in the phaser experiments, the FFT length is set to  $2^{\lceil \log_2(N_{WIN}) \rceil + 1}$ . The LSTM configuration is trained with 64 hidden units. Note pitch and durations are provided to the synth for reconstruction and the phase for the oscillator is random to improve robustness and reflect how synths behave in reality. As a result, the target and reconstructed audio may be misaligned which is why we use multi-resolution STFT loss (MSS) [28] for training which is phase agnostic.

We evaluate the different filter configurations by comparing their MSS loss values on the 20% test split. The FS configurations that operate at frame-rate are also evaluated at sample-rate by linearly interpolating their filter coefficients during inference. We also calculate the Fréchet Audio Distance (FAD) [29] for each model which has been shown to correlate with human perception. Since the individual audio files are very short, we first concatenate them into one audio file before calculating the FAD. To avoid harsh discontinuities in the concatenated file, we apply a 32-sample fade (one hop-size  $L$ ) to both ends of individual clips. The evaluation results are summarised in Table 3. Finally, in Table 4 we show the speed benchmarks of the different synth configurations.

<sup>6</sup><https://www.samplescience.info/2022/05/abstract-303.html>

Table 3: Synth evaluation results (FS = frequency sampling, TD = time domain) with 95% confidence intervals for FAD scores.

Filter	Method	$N_{\text{WIN}}$	MSS		FAD VGGish	
			FS	TD	FS	TD
Coeff.	FS	4096	1.66	1.78	$2.62 \pm 0.09$	$2.70 \pm 0.13$
		2048	1.64	1.65	<b><math>2.18 \pm 0.07</math></b>	$2.35 \pm 0.11$
		1024	1.53	1.58	$2.57 \pm 0.08$	$2.27 \pm 0.12$
		512	1.57	1.57	$2.87 \pm 0.10$	$2.46 \pm 0.10$
		256	<b>1.49</b>	1.48	$2.25 \pm 0.08$	<b><math>1.98 \pm 0.06</math></b>
		128	1.53	1.55	$3.37 \pm 0.14$	$2.73 \pm 0.12$
LP	FS	4096	1.96	1.98	$2.59 \pm 0.06$	<b><math>2.09 \pm 0.07</math></b>
		2048	1.95	2.04	$2.62 \pm 0.07$	$4.52 \pm 0.17$
		1024	1.89	2.15	$2.59 \pm 0.08$	$4.18 \pm 0.14$
		512	1.83	2.92	<b><math>2.13 \pm 0.06</math></b>	$3.38 \pm 0.08$
		256	<b>1.82</b>	2.89	$2.17 \pm 0.06$	$3.36 \pm 0.12$
		128	1.84	2.70	$2.34 \pm 0.09$	$3.93 \pm 0.12$
LSTM 64	TD	-	-	<b>1.38</b>	-	$2.49 \pm 0.21$
		-	-	-	-	$2.51 \pm 0.10$
LSTM 64	TD	-	-	1.76	-	$3.24 \pm 0.07$

Table 4: Synth CPU runtime benchmarks on an M1 Pro MacBook for one optimisation step (forward + backward, one thread, single precision, batch size of 34).

Synth	TD	FS $N_{\text{WIN}}$					
		128	256	512	1024	2048	4096
Coeff.	32 ms	57 ms	102 ms	201 ms	390 ms	833 ms	1795 ms
LP	29 ms	58 ms	98 ms	195 ms	376 ms	804 ms	1667 ms
LSTM 64	1322 ms	-	-	-	-	-	-

Looking at the evaluation and benchmarking results, we observe that both configurations of our TD filter perform well and generally match or outperform the corresponding FS implementations. Our method also provides roughly a 2x to 30x speedup over the FS and LSTM configurations. The low-pass FS methods perform significantly worse when applied in the time domain which we attribute to overfitting of the neural network to the window size of the filter and can be clearly heard in the resulting audio. This is less the case for the learned coefficient filters, especially according to the FAD metric. We hypothesise this could be due to the linear interpolation at inference of the learned coefficients which prevents harsh discontinuities from occurring at the frame boundaries, resulting in a smoother time-varying transfer function. We also found during training that the Coeff FS synths would sometimes not converge, even when using gradient clipping, whereas our TD implementations never experienced stability issues.

### 5.3. Modelling the LA-2A Leveling Amplifier

For the compressor experiments, the targets we model are 1) the feed-forward compressor (FF) in Sec. 4.3 and 2) a Universal Audio LA-2A analog compressor (LA). We optimise the proposed differentiable FF compressor ( $\nabla\text{FF}$ ) to match the target sounds, examining its capability to replicate and infer the parameters of dynamic range controllers. We test the following conditions:

- FF-A:  $R = 3$ , 1 ms attack and 100 ms release
- FF-B:  $R = 5$ , 30 ms attack and 30 ms release
- FF-C:  $R = 8$ , 0.1 ms attack and 200 ms release
- LA-D: compressor mode, 25 peak reduction
- LA-E: compressor mode, 50 peak reduction
- LA-F: compressor mode, 75 peak reduction

Table 5: Summary of compressor ESR (%) evaluation.

Method	FF-A	FF-B	FF-C	LA-D	LA-E	LA-F
FS	2.362	<b>0.00780</b>	4.649	11.29	9.485	7.783
$\nabla\text{FF}$	<b>0.015</b>	0.00785	<b>0.017</b>	<b>10.58</b>	<b>9.356</b>	<b>7.639</b>

Table 6: The learned parameters for matching a LA-2A.

Method	Data	$R$	$CT$ (dB)	Attack	Release	$\alpha_{\text{rms}}$	$\gamma$ (dB)
FS	D	9.1	-11.66	489.43 ms		0.008	0.69
	E	231.1	-19.08	44.62 ms		0.606	0.34
	F	2.9	-26.00	0.06 ms		0.002	-0.81
$\nabla\text{FF}$	D	39.0	-26.58	99.41 ms	0.06 ms	0.703	0.74
	E	13.1	-12.41	5.68 ms	420.56 ms	0.978	0.54
	F	5.4	-20.14	2.24 ms	229.15 ms	0.973	-0.13

Table 7: Compressor runtime benchmarks for different lengths of 44.1 kHz audio. Ran on an M1 Pro MacBook for one optimisation step (forward + backward, one thread, single precision).

Method	Sequence duration		
	30 s	60 s	120 s
FS	163.4 ms	320.8 ms	663.8 ms
$\nabla\text{FF}$	64.9 ms	117.9 ms	239.4 ms

The  $\alpha_{\text{rms}}$ ,  $CT$  and  $\gamma$  for  $\text{FF}_*$  are set to 0.03,  $-20$  and  $0$  dB, respectively. We train and evaluate our compressors on the Signal-Train dataset [30], which consists of paired data recorded in 44.1 kHz from the LA-2A compressor with different peak reduction values. Following [5], we select files with the sub-string 3c in the file name. Each LA-2A setting has 20 min of paired audio containing real-world musical sounds and synthetic test signals. We use the first 5 min for training and the rest for evaluation. The same input data is used for FF-A/B/C, and the target audio is generated by applying the FF compressor with the target parameters. We pick the simplified compressor [4] as our baseline, which uses the same FF compressor but has  $\alpha_{\text{at}} = \alpha_{\text{rt}}$ . We denote this baseline as FS, as it uses frequency sampling to compute (19) and (21).

The parameters we optimise are  $\{\hat{R}, CT, \hat{\alpha}_{\text{at/bt/rms}}, \gamma\}$ . We set  $\alpha_* = \text{sigmoid}(\hat{\alpha}_*)$ ,  $R = \exp(\hat{R}) + 1$ . The initial values are 50 ms attack/release,  $R = 2$ ,  $CT = -10$  dB,  $\alpha_{\text{rms}} = 0.3$ , and  $\gamma = 0$  dB. The conversion from time  $t$  in seconds to coefficients  $\alpha_*$  is  $1 - \exp(-\frac{2.2}{44100t})$ . We train each compressor without using mini-batching for at least 1000 epochs using stochastic gradient descent with a learning rate of 100 and 0.9 momentum, minimising the mean absolute error (MAE). For evaluation, we select the parameters with the lowest training loss. We apply the same pre-filter from [5] before calculating loss and evaluation metrics. We use double precision when computing gigantic FFTs for the FS method to avoid numerical overflow.

Table 5 shows that  $\nabla\text{FF}$  has a lower ESR than FS besides condition FF-B. The parameters learned for condition FF-A/B/C using  $\nabla\text{FF}$  are close to the ground truth. FS can only recover the parameters when attack and release are identical (FF-B). For LA-2A,  $\nabla\text{FF}$  reasonably captures the analog characteristics (fast attack and slow release, shown in Table 6) of condition E/F. We found FS tends to learn unrealistically large attack/release times and ratios, as seen in Table 6, likely due to its simplified design.

Table 7 shows our method is two to three times faster than FS. Training takes roughly 43 min for FS and 17 min for  $\nabla\text{FF}$  on an M1 Pro MacBook.

## 6. CONCLUSION AND FUTURE WORK

In this work, we propose an efficient backpropagation algorithm and implementation for an all-pole filter that can be used to model time-varying analog audio systems end-to-end using gradient descent. We demonstrate its advantages over previous frequency sampling approximations by using it to model a phaser, a time-varying subtractive synthesiser, and a feed-forward compressor. Our method outperforms frequency sampling in accuracy and training efficiency, especially when using the systems at the sample-rate level. We make our code and audio samples available and provide the trained audio effect and synth models in a VST plugin.

Our future work involves extending the backpropagation algorithm to work with differentiable initial conditions and applying it to relevant tasks. We also plan on benchmarking the forward-mode differentiation of our filter, investigating its numerical stability, and extending our gradient derivation to higher-order optimisation use cases.

## 7. ACKNOWLEDGEMENTS

Funded by UKRI and EPSRC as part of the “UKRI CDT in Artificial Intelligence and Music”, under grant EP/S022694/1, and by the Scottish Graduate School of Arts & Humanities (SGSAH).

## 8. REFERENCES

- [1] Jesse Engel, Lamtharn (Hanoi) Hantrakul, Chenjie Gu, and Adam Roberts, “DDSP: Differentiable digital signal processing,” in *International Conference on Learning Representations*, 2020.
- [2] Chin-Yun Yu and György Fazekas, “Singing voice synthesis using differentiable LPC and glottal-flow-inspired wavetables,” in *Proc. International Society for Music Information Retrieval*, 2023, pp. 667–675.
- [3] Shahan Nercessian, Andy Sarroff, and Kurt James Werner, “Lightweight and interpretable neural modeling of an audio distortion effect using hyperconditioned differentiable biquads,” in *ICASSP. IEEE*, 2021, pp. 890–894.
- [4] Christian J Steinmetz, Nicholas J Bryan, and Joshua D Reiss, “Style transfer of audio effects with differentiable signal processing,” *Journal of the Audio Engineering Society*, vol. 70, no. 9, pp. 708–721, 2022.
- [5] Alec Wright and Vesa Välimäki, “Grey-box modelling of dynamic range compression,” in *DAFx*, 2022, pp. 304–311.
- [6] Lauri Juvela, Bajibabu Bollepalli, Junichi Yamagishi, and Paavo Alku, “GELP: GAN-excited linear prediction for speech synthesis from mel-spectrogram,” in *Proc. INTERSPEECH*, 2019, pp. 694–698.
- [7] Alistair Carson, Simon King, Cassia Valentini Botinhao, and Stefan Bilbao, “Differentiable grey-box modelling of phaser effects using frame-based spectral processing,” in *DAFx*, 2023, pp. 219–226.
- [8] Purbaditya Bhattacharya, Patrick Nowak, and Udo Zölzer, “Optimization of cascaded parametric peak and shelving filters with backpropagation algorithm,” in *DAFx*, 2020, pp. 101–108.
- [9] Shahan Nercessian, “Neural parametric equalizer matching using differentiable biquads,” in *DAFx*, 2020, pp. 265–272.
- [10] Joseph T Colonel, Christian J Steinmetz, Marcus Michelen, and Joshua D Reiss, “Direct design of biquad filter cascades with deep learning by sampling random polynomials,” in *ICASSP. IEEE*, 2022, pp. 3104–3108.
- [11] Taejun Kim, Yi-Hsuan Yang, and Juhan Nam, “Joint estimation of fader and equalizer gains of DJ mixers using convex optimization,” in *DAFx*, 2022, pp. 312–319.
- [12] Julius O. Smith III, *Spectral Audio Signal Processing*, <https://ccrma.stanford.edu/~jos/sasp/>, accessed 2024-03-27, online book, 2011 edition.
- [13] John D. Markel and Augustine H. Gray, *Linear Prediction of Speech*, vol. 12 of *Communication and Cybernetics*, Springer, Berlin, Heidelberg, 1976.
- [14] Jean-Marc Valin and Jan Skoglund, “LPCNet: Improving neural speech synthesis through linear prediction,” in *ICASSP. IEEE*, 2019, pp. 5891–5895.
- [15] Achuth Rao MV and Prasanta Kumar Ghosh, “SFNet: A computationally efficient source filter model based neural speech synthesis,” *IEEE Signal Processing Letters*, vol. 27, pp. 1170–1174, 2020.
- [16] Suhyeon Oh, Hyungseob Lim, Kyunguen Byun, Min-Jae Hwang, Eunwoo Song, and Hong-Goo Kang, “ExcitGlow: Improving a WaveGlow-based neural vocoder with linear prediction analysis,” in *Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*. IEEE, 2020, pp. 831–836.
- [17] Udo Zölzer, *DAFx: Digital Audio Effects*, chapter Nonlinear Processing, pp. 110–112, John Wiley & Sons, 2011.
- [18] Joseph T Colonel and Joshua D Reiss, “Approximating ballistics in a differentiable dynamic range compressor,” in *Audio Engineering Society Convention 153*. Audio Engineering Society, 2022.
- [19] Zixun Guo, Chen Chen, and Eng Siong Chng, “DENT-DDSP: Data-efficient noisy speech generator using differentiable digital signal processors for explicit distortion modelling and noise-robust speech recognition,” in *Proc. INTERSPEECH*, 2022, pp. 3799–3803.
- [20] Marco Forgiione and Dario Piga, “dynoNet: A neural network architecture for learning dynamical systems,” *International Journal of Adaptive Control and Signal Processing*, vol. 35, no. 4, pp. 612–626, 2021.
- [21] Chin-Yun Yu and György Fazekas, “Differentiable time-varying linear prediction in the context of end-to-end analysis-by-synthesis,” *arXiv:2406.05128*, 2024.
- [22] Julius O. Smith III, *Physical Audio Signal Processing*, <https://ccrma.stanford.edu/~jos/pasp/>, accessed 2023-02-28, online book, 2010 edition.
- [23] Ben Hayes, Charalampos Saitis, and György Fazekas, “Sinusoidal frequency estimation by gradient descent,” in *ICASSP. IEEE*, 2023.
- [24] Roope Kiiski, Fabián Esqueda, and Vesa Välimäki, “Time-variant gray-box modeling of a phaser pedal,” in *DAFx*, 2016, pp. 31–38.
- [25] Joseph Turian, Jordie Shier, George Tzanetakis, Kirk McNally, and Max Henry, “One billion audio sounds from GPU-enabled modular synthesis,” in *DAFx*, 2021, pp. 222–229.
- [26] JD Sleep, “Small Stone Information [Online],” <https://generalguitargadgets.com/effects-projects/phase-shifters/small-stone-information/>, accessed 2023-03-26.
- [27] Christopher Mitcheltree, Christian J Steinmetz, Marco Comunità, and Joshua D Reiss, “Modulation extraction for LFO-driven audio effects,” in *DAFx*, 2023, pp. 94–101.
- [28] 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 *ICASSP. IEEE*, 2020, pp. 6199–6203.
- [29] Kevin Kilgour, Mauricio Zuluaga, Dominik Roblek, and Matthew Sharifi, “Fréchet audio distance: A reference-free metric for evaluating music enhancement algorithms,” in *Proc. INTERSPEECH*, 2019, pp. 2350–2354.
- [30] Scott Hawley, Benjamin Colburn, and Stylianos Ioannis Mimitakis, “Profiling audio compressors with deep neural networks,” in *Audio Engineering Society Convention 147*. Audio Engineering Society, 2019.