

## DIFFERENTIABLE GREY-BOX MODELLING OF PHASER EFFECTS USING FRAME-BASED SPECTRAL PROCESSING

*Alistair Carson*\*

Acoustics and Audio Group  
University of Edinburgh  
Edinburgh, UK  
alistair.carson@ed.ac.uk

*Simon King*

Centre for Speech Technology Research  
University of Edinburgh  
Edinburgh, UK  
simon.king@ed.ac.uk

*Cassia Valentini-Botinhao*

Centre for Speech Technology Research  
University of Edinburgh  
Edinburgh, UK  
cvbotinh@inf.ed.ac.uk

*Stefan Bilbao*

Acoustics and Audio Group  
University of Edinburgh  
Edinburgh, UK  
sbilbao@ed.ac.uk

### ABSTRACT

Machine learning approaches to modelling analog audio effects have seen intensive investigation in recent years, particularly in the context of non-linear time-invariant effects such as guitar amplifiers. For modulation effects such as phasers, however, new challenges emerge due to the presence of the low-frequency oscillator which controls the slowly time-varying nature of the effect. Existing approaches have either required foreknowledge of this control signal, or have been non-causal in implementation. This work presents a differentiable digital signal processing approach to modelling phaser effects in which the underlying control signal and time-varying spectral response of the effect are jointly learned. The proposed model processes audio in short frames to implement a time-varying filter in the frequency domain, with a transfer function based on typical analog phaser circuit topology. We show that the model can be trained to emulate an analog reference device, while retaining interpretable and adjustable parameters. The frame duration is an important hyper-parameter of the proposed model, so an investigation was carried out into its effect on model accuracy. The optimal frame length depends on both the rate and transient decay-time of the target effect, but the frame length can be altered at inference time without a significant change in accuracy.

### 1. INTRODUCTION

A broad class of audio effects found in almost all genres of popular music is that of time-varying modulation effects, and includes phasing, flanging, chorus, and tremolo. Model-based digital implementations of these effects are straightforward [1], but many musicians prefer the timbre and character of the original analog or electro-mechanical devices used to create these effects, and these may be considerably more difficult to model. Circuit-based simulations [2, 3] can produce physically accurate results but are highly

optimized to a specific device, and require complete knowledge of the circuit and its component values.

In general, modelling of audio effects using machine learning has become an active area of research in recent years, with a particular focus on modelling non-linear time-invariant effects such as guitar amplifiers and distortion pedals. Approaches in modelling these systems include fully black-box methods using recurrent or convolutional neural networks (RNNs / CNNs) [4, 5], as well as grey-box models which use some prior knowledge of the reference system, such as differentiable state-space models [6, 7] and differentiable DSP-based models [8, 9]. Modelling of effects with time-varying input-dependent behaviour such as dynamic range compression has also been explored through the use of CNNs with long receptive fields [10] and temporal feature-wise linear modulation [11]; as well as a differentiable DSP model proposed in [12]. Modelling of time-varying modulation effects presents a unique challenge due to the modulation of system behaviour by a low frequency oscillator (LFO). The deep-learning approach proposed in [13] can emulate a wide range of effects, however the use of bi-directional long-term short-memory networks (LSTMs) makes the model non-causal and therefore not suited for real-time use. Wright et al. [14] proposed a real-time model for phasing and flanging effects using RNNs, but required manual measurement and estimation of the LFO prior to training as this was an input to the model. Earlier work by Kiiski et al. [15] proposed a grey-box model of phasing, but in which no machine learning was used.

This work presents a differentiable DSP model of a phaser effect that can be trained through gradient descent to jointly learn the underlying LFO signal and the time-varying spectral response of an analog reference pedal. The model utilises frequency domain approximations of IIR filters to accelerate training times, as has been employed in [9, 12]. Through experiment, we investigate the conflicting demands of time and frequency resolution when using this method in the context of time-varying effects.

The paper is structured as follows: Section 2 provides background on analog and digital phasing effects; Section 3 outlines the proposed model; Section 4 describes the target systems and data; Section 5 describes the experimental procedure and results; and Section 6 concludes the paper with an outlook on areas of future work. Source code and audio examples are provided at the accompanying web-page [16].

\* A Carson is funded by the Scottish Graduate School of Arts and Humanities (SGSAH).

Copyright: © 2023 Alistair Carson et al. This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International License, which permits unrestricted use, distribution, adaptation, and reproduction in any medium, provided the original author and source are credited.

## 2. ANALOG AND DIGITAL PHASE SHIFTERS

The phaser is a time-varying filter effect in which the phase of an input signal is subject to a periodic modulation and combined together with the ‘dry’ (unaffected) signal to create audible notches in the frequency spectrum [17]. The movement of these notches with time gives the characteristic perceived sweeping effect. The phasing effect is often confused with flanging, and indeed both effects are caused by modulating notches in the spectra. The key difference is that a flanger generates an infinite series of harmonically spaced notches, whereas the phaser response has finite non-uniformly spaced notches [18]. Furthermore, implementations of phasers and flangers differ both historically and in present day commercial products. Originally, the flanging effect was created by playing two tape machines in near unison with a small variable time delay between the two reels [19]. Nowadays, flangers are typically implemented using time-varying digital comb filters, or ‘bucket-brigade’ delay lines in analog pedals. In contrast, phasers create phase shifts through cascaded all-pass filters, in which break frequencies are modulated by a low-frequency oscillator (LFO).

Digital phaser effects are usually implemented as linear time-varying signal-based models, in which the all-pass filters are discrete-time approximations of idealised all-pass filters [1, 18, 20]. The motivation of this work is to explore whether such a model can be embedded in a machine learning framework to emulate the response of an analog phaser pedal. The continuous and discrete time signal processing concepts of phasing, which underpin the proposed model, are outlined in the remainder of this section.

### 2.1. Continuous-time phasing

We will examine phasers of the topology shown in Figure 1 composed of a cascade of  $K$  identical first-order all-pass filters. (In many practical implementations of the phaser,  $K = 4$  [18], but other choices of  $K$  have been employed [15].) The continuous-time transfer function of each such first-order section is [18]:

$$A(s) = \frac{s - \omega_b}{s + \omega_b} \quad (1)$$

where  $s = \sigma + j\omega$  is the complex frequency and  $\omega_b \geq 0$  is known as the *break-frequency* of the all-pass filter. Because  $A$  is all-pass, for real frequencies  $s = j\omega$  it may be written as  $A = e^{j\Theta}$ , where:

$$\Theta(\omega) = \pi - 2 \arctan(\omega/\omega_b). \quad (2)$$

From (2) it can be observed that  $\omega_b$  is the frequency at which  $\Theta(\omega_b) = \pi/2$ . Furthermore, the all-pass filter inverts DC, and the phase response tends to zero as  $\omega \rightarrow \infty$ .

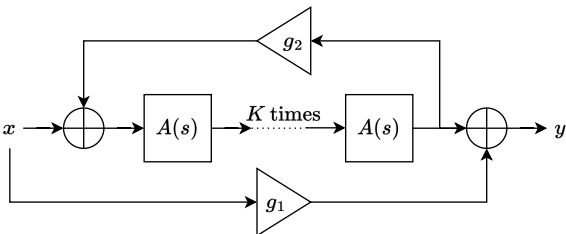


Figure 1: A typical phaser structure in continuous time.

The cascade arrangement in Figure 1 includes a through path with gain  $g_1$ , and a feedback path with gain  $g_2$ , with  $0 \leq g_2 < 1$ .

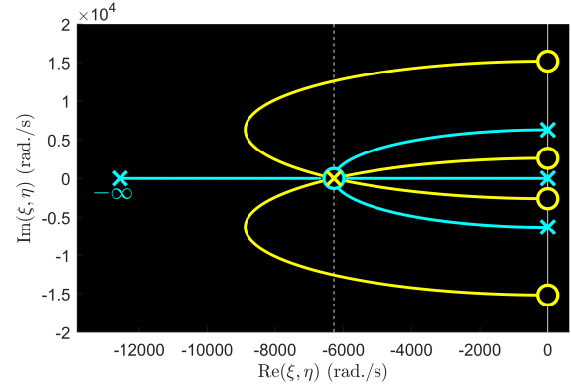


Figure 2: Root locus plot, for  $g_1 = 1$ , and for for  $g_2$  between 0 (open loop, indicated by circles) and 1 (indicated by crosses). Pole locations are indicated by blue lines, and zero locations by yellow lines. The real frequency  $-\omega_b$  (where here,  $\omega_b = 2\pi \cdot 1000$  rad/s) is indicated by a dashed line.

The following transfer function results:

$$H(s) = g_1 + \frac{A^K}{1 - g_2 A^K}. \quad (3)$$

The poles  $s = \xi_k$ , and zeros  $s = \eta_k$ ,  $k = 0 \dots, K - 1$ , are:

$$\xi_k = \omega_b \frac{1 + \lambda_k}{1 - \lambda_k} \quad \text{where } \lambda_k = \frac{e^{j2\pi k/K}}{\sqrt[K]{g_2}} \quad (4a)$$

$$\eta_k = \omega_b \frac{1 + \beta_k}{1 - \beta_k} \quad \text{where } \beta_k = \sqrt[K]{\frac{g_1}{1 - g_2}} e^{j\pi(2k+1)/K} \quad (4b)$$

Figure 2 shows a root locus plot of the pole and zero trajectories for  $0 \leq g_2 \leq 1$ , and when  $g_1 = 1$ . Under open loop ( $g_2 = 0$ ) conditions, two notches are located on the non-negative imaginary axis. As  $g_1 \rightarrow 0$  the zeros move away from the imaginary axis and as such, parameter  $g_1$  controls the perceived depth of the phaser effect. Under fully closed loop conditions, two poles are located on the non-negative imaginary axis (one at DC). Variations in the transfer function magnitude  $|H|$  with both  $g_1$  and  $g_2$  are shown in Figure 3.

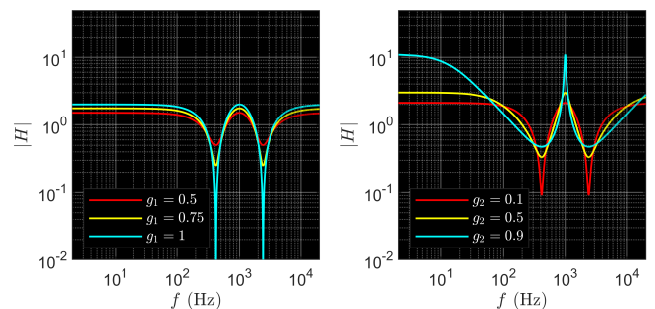


Figure 3: Magnitude response for  $H$ , under open-loop conditions at left, for different values of  $g_1$ , as indicated, and under closed-loop conditions at right, for  $g_1 = 1$ , and different values of  $g_2$ , as indicated. As before,  $\omega_b = 2\pi \cdot 1000$  rad/s.

## 2.2. Discrete-time phasing

The continuous-time transfer function (1) can be approximated in discrete-time via the bilinear transform to give the discrete-time all-pass section,  $A_d(z)$ , defined as:

$$A_d(z) = \frac{p - z^{-1}}{1 - pz^{-1}} \quad \text{where} \quad p = \frac{1 - \tan(\omega_b T/2)}{1 + \tan(\omega_b T/2)}. \quad (5)$$

Here,  $T = 1/F_s$  is the sampling period, for sample rate  $F_s$  in Hz. The corresponding discrete-time phaser is shown in Figure 4. The discrete-time transfer function  $H_d(z)$  of  $K$  cascaded all-pass sections, including the residual connection and feedback loop is:

$$H_d(z) = g_1 + \frac{A_d^K}{1 - g_2 z^{-\phi} A_d^K}. \quad (6)$$

Note that a digital delay  $z^{-\phi}$ ,  $\phi \in \mathbb{Z}^+$  has been included in the feedback loop. Without this, a delay-free loop would be present in the resulting difference equation. Efficient modelling of delay-free loops found in phaser pedals has proven challenging: the state-space model of a phaser pedal presented in [3] required a Newton-Raphson solver at run-time; and Kiiski et al. reported that the fictitious delay line in their DSP model ( $\phi = 1$ ) resulted in perceptual differences in the feedback effect when compared with the analog reference device [15]. The effect of the delay line on the magnitude response of the system can be seen in Figure 5. It is clear that the fictitious delay line in the feedback loop drastically changes the magnitude response for high-frequencies, and should ideally be avoided in virtual analog models of phasers.

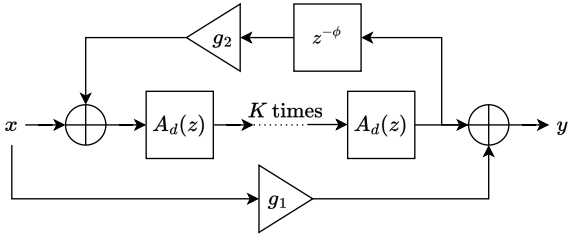


Figure 4: A typical DSP phaser structure.

## 2.3. Time-varying behaviour

In analog phase shifters, the break frequency of the all-pass sections,  $\omega_b$ , is periodically modulated by a sub-audio rate LFO, typically with frequency in the range 0.05 Hz to 5 Hz. The mapping between LFO voltage and break frequency is often nonlinear and asymmetrical: for example, in the MXR Phase 90 the LFO varies the voltage across JFETs which in turn alters the break frequency of the all-pass [3]. Certain devices, such as the Uni-Vibe and its replicas, are optical: the LFO varies the voltage across a light source (an incandescent bulb in original units [21]) surrounded by light dependent resistors (LDRs). The resulting fluctuation in current through the LDRs controls the all-pass break frequencies. Accurate prediction of the LFO characteristics (including fundamental frequency and waveshape) are therefore critical to the quality of discrete-time emulations of such effects, including phasers.

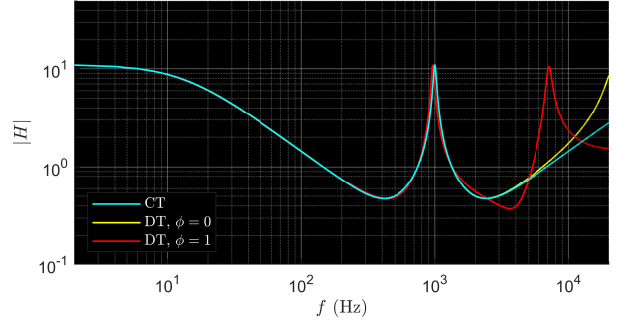


Figure 5: Comparison of magnitude responses of continuous-time (CT) and discrete-time (DT) phaser models, with constant parameters  $\omega_b = 2\pi \cdot 1000$  rad/s,  $g_1 = 1.0$ ,  $g_2 = 0.9$ , and  $T = 1/44100$  s. Distinct cases of the discrete-time model without ( $\phi = 0$ ) and with ( $\phi = 1$ ) a fictitious delay are illustrated.

## 3. MODELLING METHOD

This section outlines the proposed model of a differentiable phase shifting algorithm which, given input-output audio recordings of a reference device, can be trained to emulate the time-varying behaviour. Consider an arbitrary phase-shifting device that has been sampled in discrete time as:

$$y[n] = f(x[n], \theta[n]) \quad (7)$$

where  $x$  is the input signal,  $y$  is the output signal,  $\theta$  are the time-varying parameters of the system,  $n$  is the sample index and  $f$  is a linear function of  $x$ . We seek to develop a model of the form:

$$\hat{y}[n] = g(x[n], \hat{\theta}[n]) \quad (8)$$

whose output  $\hat{y}[n]$  is perceptually indistinguishable from  $y[n]$ . The function  $g$  is assumed differentiable so that the model can be trained using gradient descent to find model parameters  $\hat{\theta}$  that minimise an objective loss function  $\mathcal{L}(y, \hat{y})$ . The proposed model architecture is shown in Figure 6.

### 3.1. Frame-based processing

Utilising frame-based spectral processing [22, 23], the proposed model assumes that the target phaser can be treated as a linear time-invariant (LTI) system over the duration of a short frame of length  $W$  seconds. Suppose that an input audio signal  $x[l]$  of length  $L$  samples is segmented into  $N_f = \lceil L/H \rceil$  frames of length  $N = \lfloor WF_s \rfloor$  with hop-size  $H$  samples. (The resulting frame rate is  $F_f = F_s/H$  in Hz.) The  $m$ th frame,  $m = 0, \dots, N_f - 1$ , is defined as the  $N \times 1$  column vector  $\mathbf{x}_m = [x[mH], \dots, x[mH + N - 1]]^T$ . The short-time Fourier transform vector  $\mathbf{X}_m$ ,  $m = 0 \dots, N_f - 1$  is derived from  $\mathbf{x}_m$  through windowing and Fourier transformation as follows:

$$\mathbf{X}_m = \mathbf{U} \mathbf{Q} \mathbf{x}_m. \quad (9)$$

Here,  $\mathbf{Q} = [\mathbf{W} \mathbf{Z}]^T$  is an  $N' \times N$  windowing matrix, where where  $N'$  is the DFT length. It includes the diagonal  $N \times N$  matrix  $\mathbf{W}$ , containing samples of the Hann window on the diagonal, and an  $N \times (N' - N)$  all-zero matrix  $\mathbf{Z}$  implementing zero-padding.  $\mathbf{U}$  is an  $N' \times N'$  DFT matrix. At each frame, spectral

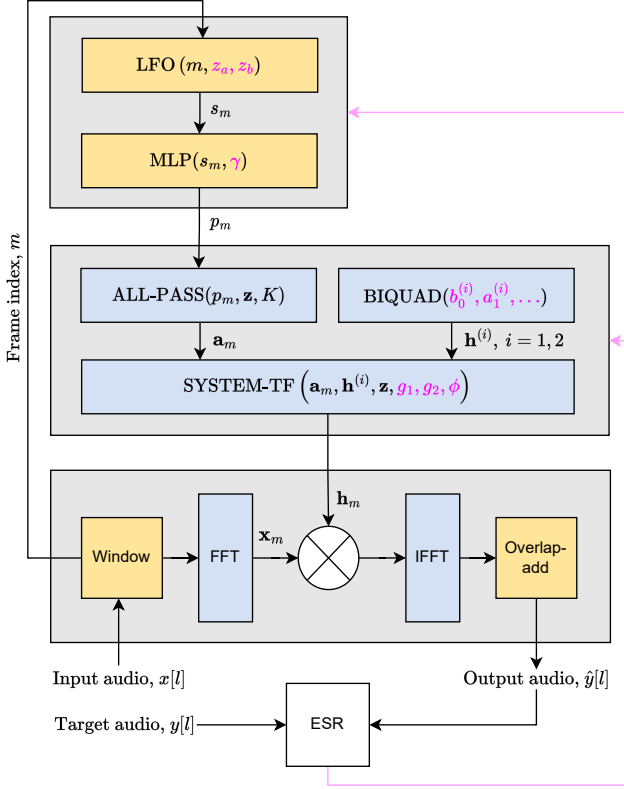


Figure 6: Structure of the proposed model. Black arrows indicate signal flow, with magenta indicating the flow of gradients to the learnable parameters. The blue and yellow boxes indicate operations in the frequency domain and time domain respectively. The subscript  $m$  denotes parameters varying at the frame rate, with all other parameters held constant for the duration of a training epoch. Parameters in magenta are updated once per epoch by the optimizer.

processing is applied via element-wise complex multiplication in the frequency domain, followed by an inverse Fourier transform, truncation and windowing to yield the  $N \times 1$  output vectors  $\mathbf{y}_m$ ,  $m = 0, \dots, N_f - 1$ :

$$\mathbf{y}_m = \frac{1}{N'} \mathbf{Q}^T \mathbf{U}^* \mathbf{H}_m \mathbf{X}_m. \quad (10)$$

Here,  $\mathbf{H}_m$  is an  $N' \times N'$  diagonal matrix containing values of a transfer function (incorporating Hermitian symmetry) at frame  $m$  on its diagonal.  $\mathbf{U}^*$  is the conjugate transpose of the DFT matrix  $\mathbf{U}$ . Finally, the output time series  $\hat{y}[l]$  is obtained from the frames  $\mathbf{y}_m$  through an overlap-add procedure, with hop size  $H$ . It can be shown that for  $\mathbf{H}_m = \mathbf{I}$ , exact reconstruction of the input signal can be obtained if  $N/H$  is an integer greater than two—this property is known as constant-overlap-add and is enforced in the proposed model. Considering the transfer function in (10) to represent a digital filter with  $M \in \mathbb{Z}$  non-zero taps, then the DFT length must be  $N' \geq N + M - 1$  to avoid temporal aliasing in the output frames. For IIR filtering (where  $M \rightarrow \infty$ ), some degree of temporal aliasing is inevitable but can be practically suppressed by choosing  $M$  as the 60dB decay time (in samples) of the filter.

In this work the decay time of the target system was not known, but we found a DFT length of  $N' = 2^{\lceil \log_2(N) \rceil}$  to be sufficient to train the models.

### 3.2. LFO generator

The proposed model has an LFO module which governs the time-varying behaviour of the spectral processing. The LFO produces samples at the frame rate,  $F_f$ , and is defined as the real part of a damped complex exponential:

$$s_m(z_a, z_b) = \text{Re}(z_b z_a^m) = |z_b| |z_a|^m \cos(m \angle z_a + \angle z_b) \quad (11)$$

where  $z_a, z_b \in \mathbb{C}$  are the complex frequency and complex amplitude respectively. Hayes et al. showed that Wirtinger's calculus can be used to compute the partial derivatives of real-valued, complex-variable functions such as (11) to enable sinusoidal frequency estimation by gradient descent [24]. In this work we extend this method to include a learnable starting phase and amplitude (defined by  $z_b$ ). The Wirtinger derivatives of (11) are:

$$\frac{\partial s_m}{\partial z_a} \triangleq \frac{1}{2} \left( \frac{\partial s_m}{\partial \text{Re}(z_a)} - j \frac{\partial s_m}{\partial \text{Im}(z_a)} \right) = \frac{m z_b z_a^{m-1}}{2} \quad (12a)$$

$$\frac{\partial s_m}{\partial z_b} \triangleq \frac{1}{2} \left( \frac{\partial s_m}{\partial \text{Re}(z_b)} - j \frac{\partial s_m}{\partial \text{Im}(z_b)} \right) = \frac{z_a^m}{2}. \quad (12b)$$

In the model implementation, the LFO parameters were initialised to:

$$z_a = 0.7 \exp(j\zeta/F_f), \quad z_b = 1.0 \quad (13)$$

where  $\zeta$  is a random number sampled from a standard normal distribution. The damped amplitude envelope was only applied during training; during inference  $z_a$  was normalised to the unit circle to give a lossless LFO.

### 3.3. Multi-layer perceptron waveshaper

The output of the LFO generator is passed through a multi-layer perceptron (MLP) to allow the model to learn non-sinusoidal control signals like those described in Section 2.3. The conceptual motivation for this module was to emulate the linear and non-linear mapping of the LFO signal (in volts or amperes) to the break-frequencies of the all-pass filters (in  $\text{rad s}^{-1}$ ) in an analog phaser device. Therefore we treat the output of the MLP as the model's prediction of the normalised time-varying break-frequency signal, such that the all-pass parameter is given by:

$$p_m = \frac{1 - \tan(d_m)}{1 + \tan(d_m)} \quad \text{where} \quad d_m = \text{MLP}(s_m, \gamma), \quad (14)$$

with  $\gamma$  being the MLP parameters. The MLP consisted of three hidden layers; with 8 neurons per layer; hyperbolic tangent activation functions in the hidden layers; and linear activation in the output layer. The input and output features were scalar, based on the assumption that all  $K$  all-pass filters are identical in the reference device.

### 3.4. Model transfer function

The frame-dependent transfer function of the model has the form:

$$\mathbf{h}_m \triangleq \text{diag}(\mathbf{H}_m) = \mathbf{h}^{(1)} \cdot \left( g_1 + \frac{\mathbf{h}^{(2)} \cdot \mathbf{a}_m}{1 - |g_2| \mathbf{z}^{-I(\phi)\phi} \cdot \mathbf{h}^{(2)} \cdot \mathbf{a}_m} \right) \quad (15)$$

where  $\mathbf{a}_m$  is the frame-dependent all-pass kernel:

$$\mathbf{a}_m = \left( \frac{p_m - \mathbf{z}^{-1}}{1 - p_m \mathbf{z}^{-1}} \right)^K \quad (16)$$

and  $\mathbf{z}$  is the  $N'$ -element vector  $\mathbf{z} = [e^{\frac{2\pi j(0)}{N'}}, \dots, e^{\frac{2\pi j(N'-1)}{N'}}]^T$ . (Here and elsewhere in this section, operations on vectors are assumed to be applied element-wise.) The parameters  $g_1, g_2, \phi \in \mathbb{R}$  retain their physical meanings from Section 2 but are now initialised as learnable parameters of the model.  $I(\cdot)$  is the Heaviside step function and was included to prevent the model learning a non-casual transfer function. The number of all-pass filters is a hyper-parameter and was fixed for all experiments to  $K = 4$ . The model includes two frequency domain representations of bi-quad filters, given by:

$$\mathbf{h}^{(i)} = \frac{b_0^{(i)} + b_1^{(i)} \mathbf{z}^{-1} + b_2^{(i)} \mathbf{z}^{-2}}{1 + a_1^{(i)} \mathbf{z}^{-1} + a_2^{(i)} \mathbf{z}^{-2}} \quad (17)$$

where  $b_0^{(i)}, b_1^{(i)}, b_2^{(i)}, a_1^{(i)}, a_2^{(i)} \in \mathbb{R}$  are the learnable filter parameters for the  $i$ th biquad,  $i = 1, 2$ . These kernels were included to account for any further LTI filtering that an analog phaser might impart in addition to the core phasing effect described in Section 2. For example, low-pass filtering, DC blocking or gain adjustment.

### 3.5. Loss function

The proposed model uses the error-to-signal ratio (ESR) as the objective loss function during training:

$$\mathcal{L}(y, \hat{y}) = \frac{\sum_{l=0}^{L-1} (y[l] - \hat{y}[l])^2}{\sum_{l=0}^{L-1} y[l]^2} \quad (18)$$

where  $L$  is the length of the training data in samples. This loss function has been widely used in black-box and grey-box modelling of other audio effects [5, 12, 14], but the strict time-alignment required by (18) introduced some interesting challenges when it came to training the proposed model. Because the frequency and phase of the LFO are unknown parameters, the training data cannot be arbitrarily split into short segments (as is common when using long audio sequences as training data [5, 11]). For example, even if the initial random frequency guess was precisely correct, the model would ‘see’ a different starting phase for each segment (unless the segment length happened to be an integer multiple of the LFO period). In initial experiments, this phase discrepancy was accounted for using the current estimate of LFO frequency, but this caused noisy optimizer updates and convergence issues when the segment length was shorter than one LFO period in the target data.

This presents a dilemma in training this model: we need a sufficient duration of training data to capture the slowly time-varying features in the target system, but are constrained to training with single-batch gradient descent, meaning the time taken for one optimizer step increases linearly with the length of training data. The proposed solution to this problem was to use a short, spectrally-flat training signal. Details of this signal are outlined in Section 4.

### 3.6. Training details

All models were trained on audio with a sample rate of 44.1 kHz using an Adam optimizer [25] with an initial learning rate of  $10^{-3}$ . Models were trained for a maximum of 5000 epochs on a NVIDIA

Titan-X GPU. The training times varied depending on the length of training data and window size. For example, for an audio sequence of 10 s the training times were  $\sim 3$  hours and  $\sim 16$  hours for window lengths of 160 ms and 10 ms respectively. It is important to note that early designs of the proposed model were implemented in the time-domain and trained via back-propagation through time, but training times were deemed too slow to pursue this approach further ( $\sim 24$  hours for 1000 epochs on 1s of training audio). The prohibitive training times of IIR filters has been reported in previous work [9] [12].

### 3.7. Model inference

In this work, at model inference we use the same algorithm as in training (i.e. using frame-based spectral processing). This has the limitation of introducing a minimum latency of  $W$  seconds into the system, which may be unsuitable for real-time use. A time-domain implementation using IIR filters could conceivably be derived via inverse  $z$ -transform of the system transfer function (15), but this is left as a task for future work. The handling of the possibly non-integer delay-line length  $\phi$  would require some consideration, but could be implemented with an all-pass filter in the feedback loop shown in Figure 4.

## 4. TARGET SYSTEMS AND DATASETS

A custom dataset was collected consisting of 60s of a synthetic chirp-train signal followed by 60s of direct-input (DI) guitar recordings. The chirp-train signal was used as the input signal for model training; whereas the guitar recordings were reserved for testing.

### 4.1. Synthetic training signal

The chirp-train signal was synthesised as an impulse train with period 30 ms passed through a cascade of 64 all-pass filters (5) with  $p = 0.9$ . This type of signal has been used previously for estimating the LFO frequency and shape of LTV audio effects [14, 15]. Due to its spectral flatness, it was hypothesised that even a few seconds of this signal would be sufficient to train the model.

### 4.2. Digital phaser

As a simplified test problem, the proposed model was initially trained on data generated through a digital phaser with transfer function (6) ( $K = 4$ ) implemented through a time-domain recursion in MATLAB. This can be viewed as a specific instance of the model itself but without frame-based spectral processing and under known parameters, shown in Table 2. The LFO was set to a triangular wave, sweeping through break-frequencies from 4000  $\text{rad s}^{-1}$  to a maximum of 16 000  $\text{rad s}^{-1}$ . The maximum 60 dB decay time of this system (occurring at the minimum of the LFO cycle) was measured to be  $t_{60} = 38$  ms.

### 4.3. EHX Small Stone

The Electro-Harmonix (EHX) Small Stone is a commonly encountered analog phaser pedal. The pedal has a single knob to control the rate of the effect, and a binary ‘colour’ switch. High-level analysis of a circuit schematic [26] showed that it consisted of an LFO module, a series of four first order all-pass filter sections and a feedback loop (engaged when the ‘colour’ switch is on). This circuit therefore shares the same topology as the proposed model.

The Small Stone data was collected by processing both the training and testing audio through the pedal in one continuous recording. The audio was sent to the pedal via the output of a PreSonus AudioBox i2 interface, and the output of the pedal re-connected to the input of the audio interface. A calibration recording was obtained with the pedal in bypass-mode and used as the input to the models to negate the effect of the recording equipment on model training. Six unique parameter configurations were captured: three different positions of the ‘rate’ knob with colour switch ON (circuit with feedback) and colour switch OFF (no feedback). The LFO rates were estimated through manual inspection of the spectrogram, providing a pseudo-ground-truth  $f_0$  which could later be compared to the learned LFO signals — see Table 1.

Table 1: *Manually estimated LFO rates of the Small Stone under different parameter configurations. \* denotes approximate values.*

Label	Colour	Rate knob position	$T_0^*$ [s]	$f_0^*$ [Hz]
SS-A	OFF	3 o'clock	0.44	2.28
SS-B		12 o'clock	1.60	0.625
SS-C		9 o'clock	11.6	0.086
SS-D	ON	3 o'clock	0.70	1.4
SS-E		12 o'clock	2.56	0.38
SS-F		9 o'clock	18	0.056

## 5. EXPERIMENTS AND RESULTS

The focus of the experiments presented here is on the effect of the window length  $W$  (in seconds) on model accuracy in the context of both model training and inference. In all experiments, the accuracy metric was the resultant ESR (18) on the test dataset. It was noted that the training convergence was sensitive to the initialisation of the MLP parameters,  $\gamma$ . To address this issue, each training procedure was re-initialised and repeated three times. The iteration with the lowest ESR was retained. A frame overlap of 75% was used across all experiments, with the frame length in samples  $N$  truncated to a multiple of four to ensure constant-overlap-add [18].

### 5.1. Experiment 1: training frame size sweep

As an initial experiment, instances of the model were trained with frame lengths ranging from 10 ms to 160 ms on the following data:

- (a) Digital phaser with LFO rate  $T_0 = 2$  s (DP-2).
- (b) Small Stone with parameter configuration A (SS-A).
- (c) Small Stone with parameter configuration D (SS-D).

The training data was truncated to 2.67 s in duration to accelerate training, and was deemed sufficient given that it contained at least one LFO cycle for all case studies.

The ESR obtained in experiment 1 can be seen in Figure 7. In the case of the digital phaser, frame lengths of 40 ms to 160 ms all resulted in a error-to-signal of less than 1% on the testing data – implying an accurate match between the model output and target waveform with minimal artefacts introduced by windowing. The frame lengths of 10 ms and 20 ms produced worse results, suggesting an insufficient number of bins in the transfer function shape the spectrum and/or severe time aliasing in the output. The learned parameters during this experiment can be found in Table 2

and show a good estimation of the parameters in the target model. Figure 8 shows the synthesised LFO signals, compared to the target triangular wave.

In the case of the Small Stone, the resulting model accuracy depended on the presence of feedback in the circuit. The minimum test loss was approximately 1.5% without feedback and 10% with feedback. This result is expected due to the increase in circuit complexity and longer decay time associated with the feedback case. In both cases, window lengths of 40 ms and 80 ms provided the best performance—suggesting a good trade-off between time and frequency resolution. Despite the discrepancies in numerical results, the perceptual differences are difficult to distinguish, informally, from the target system—however, some differences in the low-frequencies are noted for short window lengths. The reader is referred to the accompanying web-page for audio examples. It is interesting to observe the learned LFO signals of the Small Stone, as shown in Figure 9. In both cases, the MLP module has consistently predicted a similar wave-shape across the frame-rates. When engaged, the colour switch appears to increase both the depth and the period of break-frequency modulation.

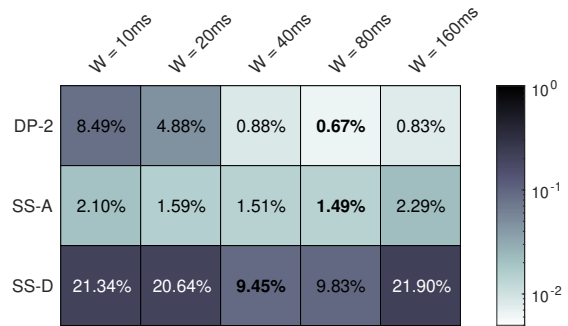


Figure 7: *ESR for different training window lengths  $W$  on the test audio of three case studies: digital phaser with  $T_0 = 2$  s (DP-2), Small Stone with parameter configurations A and D (SS-A, SS-D).*

### 5.2. Experiment 2: training frame size vs LFO rate

Experiment 2 investigated the effect of frame-length on model accuracy in more detail, considering as a case study the digital phaser with LFO periods  $T_0 = 0.5$  s, 2 s and 8 s. The length of training data was held constant at 10 s, and the prior knowledge of target LFO periods informed the choice of frame-lengths:

$$W_b = T_0 2^{b/2} / 100 \quad \text{where } b = 0, \dots, 10 \quad (19)$$

Figure 10 (top) shows the results of the experiment, with minimum test loss against training frame size for different rates of phaser effect. Firstly, we see that the model accuracy increases for longer LFO periods. This is intuitive, as in the limit  $T_0 \rightarrow \infty$  the target system becomes linear and time-invariant (LTI) so we expect the artefacts of frame-based processing to diminish. Also intuitively, the results suggest the optimum window length depends on the target LFO period. In the bottom figure, the window length has been normalised to target LFO period. In this case, the optimum  $W/T_0$  ratio shows consistency across the phaser rates, with  $W/T_0 \approx 5\%$  giving, on average, the lowest loss.

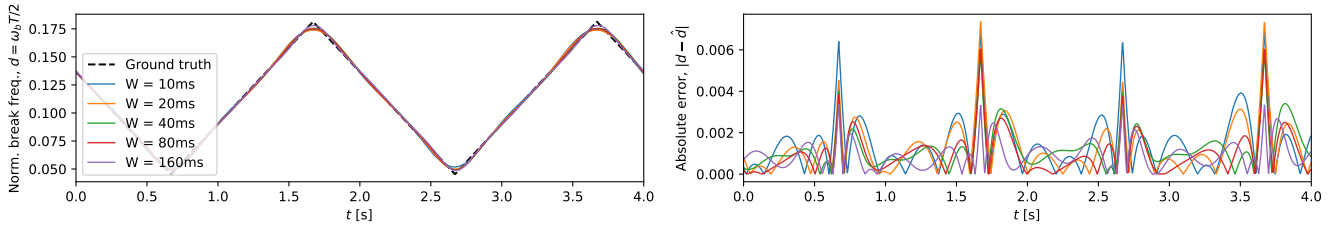


Figure 8: Outputs of the MLP module (left) and the absolute error (right) compared to the triangular LFO in the target digital phaser, with  $T_0 = 2\text{s}$  (plotted left as ground truth).

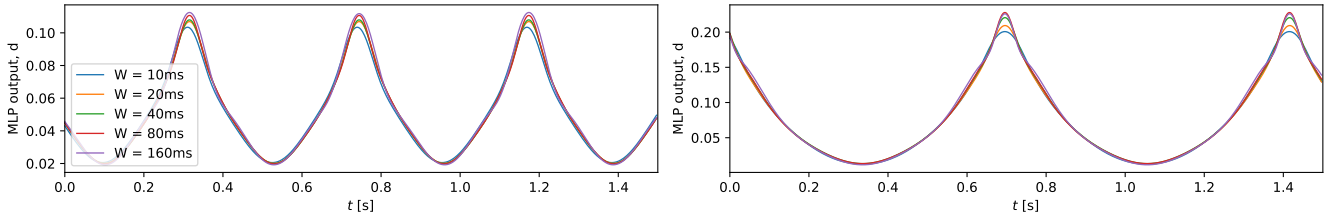


Figure 9: Outputs of the MLP module in the Small Stone modelling task for different training window lengths with feedback off (SS-A, left) and feedback on (SS-D, right). In both cases, the training audio was recorded with the pedal’s rate knob at 3 o’clock. NB the ground-truth signal is unknown so not plotted.

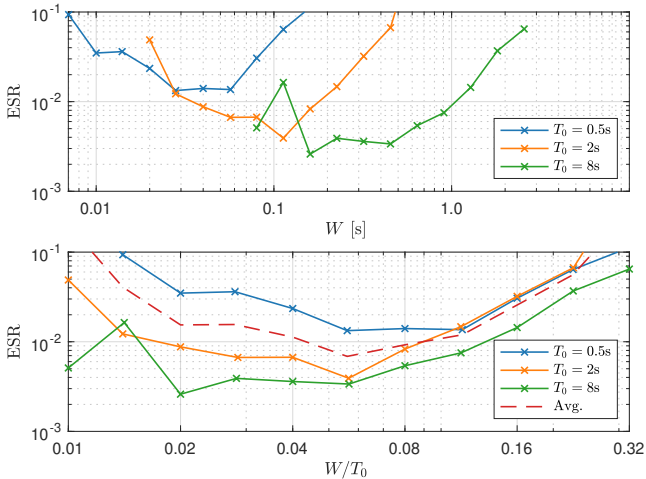


Figure 10: ESR against training window length  $W$  (top) and the ratio of training window length to target LFO period  $W/T_0$  (bottom) for the digital phaser.

### 5.3. Experiment 3: inference frame size

The aim of the final experiment was two-fold: to train instances of the model on all six parameter configurations of the Small Stone, and to investigate the effect of window length on model accuracy during inference. For each configuration, the training window length was informed by the results of Experiment 2 and set within 5-10% of the estimated LFO period (see Table 1). The training data was truncated to contain approximately three cycles of the LFO. After training, the models were tested using various window lengths at inference, with the results shown in Figure 11

In both feedback configurations, the models trained on higher

LFO rates (SS-A, SS-D) were most sensitive to changes in window size at inference time. This implies a fine balance between the window size being long enough to simulate the transient behaviour of the device, but short enough to not smear the LFO behaviour. In contrast, the accuracy of models trained on longer LFO periods (SS-C, SS-F) was mostly unchanged across the range of window sizes tested. This is a promising result, as it implies one can use a long window size for accelerated training; but a short window size for lower-latency playback at inference. However, the results suggest that there will always be a lower bound on the frame size that is determined by the decay time of the system modelled.

	$W=5\text{ms}$	$W=10\text{ms}$	$W=20\text{ms}$	$W=40\text{ms}$	$W=80\text{ms}$	$W=160\text{ms}$	$W=320\text{ms}$	$W=640\text{ms}$	$W=1.28\text{s}$	$W=2.56\text{s}$
SS-A	7.2%	2%	1.6%	1.7%	1.7%	3.7%	23%	74%	93%	>100%
SS-B	6.4%	2.2%	1.7%	1.7%	1.7%	1.7%	2.3%	10%	45%	82%
SS-C	4.3%	2.2%	1.8%	1.8%	1.8%	1.8%	1.8%	1.8%	2%	5%
SS-D	52%	21%	12%	12%	15%	34%	79%	>100%	>100%	>100%
SS-E	31%	18%	14%	14%	14%	12%	19%	40%	72%	96%
SS-F	7.9%	7.3%	7.7%	8.1%	8.3%	8.3%	8.2%	8%	8.8%	12%

$10^{-2}$   $10^{-1}$   $10^0$

Figure 11: ESR using various window sizes at inference time across the Small Stone parameter configurations. The green squares indicate the training window lengths.

Table 2: Example learned parameters of digital phaser model with  $T_0 = 2$  s (DP-2). The parameters are from the best performing model in Experiment 1, obtained using a window size of 80ms during training. Note that the biquad feedforward coefficients have been normalised.

Parameter	Description	Target value	Initial Value	Learned Value (to 3.s.f)
$f_0$	LFO rate [Hz]	0.5	0.0627	0.500
$g_1$	Wet mix	1.0	1.0	0.999
$g_2$	Feedback gain	0.7	0.01	0.700
$\phi$	Feedback delay-line length (samples)	1	0.5	0.995
$[b_{0_1}, b_{1_1}, b_{2_1}]$	Biquad 1 feedforward coeffs.	[1, 0, 0]	[1, 0, 0]	[1.00, -0.0641, 0.0336]
$[a_{1_1}, a_{2_1}]$	Biquad 1 feedback coeffs.	[0, 0]	[0, 0]	[-0.0629, 0.0336]
$[b_{0_2}, b_{1_2}, b_{2_2}]$	Biquad 2 feedforward coeffs.	[1, 0, 0]	[1, 0, 0]	[1.00, -0.0214, 0.0140]
$[a_{1_2}, a_{2_2}]$	Biquad 2 feedback coeffs.	[0, 0]	[0, 0]	[-0.0238, 0.0136]

## 6. CONCLUSIONS AND FURTHER WORK

This work has presented a differentiable DSP model of a phaser that uses frame-based spectral processing to implement a time-varying filter in the frequency domain. The model was based on a generalised continuous-time model of a phaser effect with several free parameters learnable via gradient descent. It was shown that the model can recover the parameters of a reference digital phaser and learn the correct frequency, starting phase and waveform of the underlying low-frequency oscillator (LFO), without seeing the ground-truth LFO during training. Furthermore, the model was trained to emulate an analog reference device. Informal listening found the model perceptually convincing in this task for a range of parameter configurations. Formal listening tests are important, but left for future work. It was found that the objective model accuracy depended on the training window length and required a manual estimation of the target LFO frequency for the best results. Future work will aim to remove the need for this initial estimation, perhaps through a multi-resolution training process. A key limitation of the proposed model is the inherent latency introduced by the frame-based approach, which could be problematic for real-time model inference. Future work will aim to remove this by implementing an equivalent audio-rate time-domain recursion. Finally, further work may involve extending the general approach proposed in this paper to grey-box modelling of other time-varying filters and delay-based audio effects such as auto-wah, flangers and chorus.

## 7. ACKNOWLEDGMENTS

The authors would like to thank Ben Hayes, Lauri Juvola and Alec Wright for helpful discussions on some underlying concepts of this work. Thank you to the anonymous reviewers for their comments.

## 8. REFERENCES

- [1] P. Dutilleul, M. Holters, S. Disch, and U. Zölzer, "Filters and delays," in *DAFX: Digital Audio Effects*, U. Zölzer, Ed., chapter 2, pp. 47–81. John Wiley & Sons, Ltd, 2011.
- [2] M. Holters and U. Zölzer, "Physical modelling of a wah-wah effect pedal as a case study for application of the Nodal DK Method to circuits with variable parts," in *Proc. 14th Int. Conf. Digital Audio Effects*, Paris, France, Sept. 2011.
- [3] F. Eichas, M. Fink, M. Holters, and U. Zölzer, "Physical modeling of the MXR Phase 90 guitar effect pedal," in *Proc. 17th Int. Conf. Digital Audio Effects*, Erlangen, Germany, Sept. 2014.
- [4] E.-P. Damskögg, L. Juvola, E. Thuillier, and V. Välimäki, "Deep learning for tube amplifier emulation," in *Proc. IEEE Int. Conf. on Acoust. Speech Sig. Proces.*, Brighton, UK, May 2019.
- [5] A. Wright, E.-P. Damskögg, L. Juvola, and V. Välimäki, "Real-time guitar amplifier emulation with deep learning," *Appl. Sci.*, vol. 10, no. 2, 2020.
- [6] J. Parker, F. Esqueda, and A. Bergner, "Modelling of nonlinear state-space systems using a deep neural network," in *Proc. 22nd Int. Conf. Digital Audio Effects*, Birmingham, UK, Sept. 2019.
- [7] F. Esqueda, B. Kuznetsov, and J. Parker, "Differentiable white-box virtual analog modeling," in *Proc. 24th Int. Conf. Digital Audio Effects*, Vienna, Austria, Sept. 2021.
- [8] B. Kuznetsov, J. Parker, and F. Esqueda, "Differentiable IIR filters for machine learning applications," in *Proc. 23rd Int. Conf. Digital Audio Effects*, Vienna, Austria, Sept. 2020.
- [9] S. Nercessian, A. Sarroff, and K. Werner, "Lightweight and interpretable neural modeling of an audio distortion effect using hyperconditioned differentiable biquads," in *Proc. IEEE Int. Conf. Acoust. Speech Sig. Proces.*, Toronto, Canada, June 2021.
- [10] C. Steinmetz and J. Reiss, "Efficient neural networks for real-time modeling of analog dynamic range compression," in *152nd Audio Eng. Soc. Conv.*, The Hague, Netherlands, May 2022.
- [11] M. Comunità, C. Steinmetz, H. Phan, and J. Reiss, "Modelling black-box audio effects with time-varying feature modulation," in *Proc. IEEE Int. Conf. on Acoust. Speech Sig. Proces.*, Rhodes, Greece, 2023.
- [12] A. Wright and V. Välimäki, "Grey-box modelling of dynamic range compression," in *Proc. 25th Int. Conf. Digital Audio Effects*, Vienna, Austria, Sept. 2022.
- [13] M. Martínez-Ramírez, E. Benetos, and J. Reiss, "Deep learning for black-box modeling of audio effects," *Appl. Sci.*, vol. 10, no. 2, 2020.
- [14] A. Wright and V. Välimäki, "Neural modeling of phaser and flanging effects," *J. Audio Eng. Soc.*, vol. 69, no. 7, pp. 517–529, 2021.
- [15] R. Kiiski, F. Esqueda, and V. Välimäki, "Time-variant gray-box modeling of a phaser pedal," in *Proc. 19th Int. Conf. Digital Audio Effects*, Brno, Czech Republic, Sept. 2016.
- [16] A. Carson, "Differentiable Grey-box Modelling of Phaser Effects — accompanying web-page [online]," <https://a-carson.github.io/ddsp-phaser/>, accessed 4/4/23.
- [17] B. Bartlett, "A scientific explanation of phasing (flanging)," *J. Audio Eng. Soc.*, vol. 18, no. 6, pp. 674–675, 1970.
- [18] J. O. Smith III, *Physical Audio Signal Processing*, <http://ccrma.stanford.edu/~jos/pasp/>, accessed 28/2/23, online book, 2010 edition.
- [19] H. Bode, "History of electronic sound modification," *J. Audio Eng. Soc.*, vol. 32, no. 10, pp. 730–739, 1984.
- [20] V. Välimäki, S. Bilbao, J. O. Smith, J. S. Abel, J. Pakarinen, and D. Berners, "Virtual analog effects," in *DAFX: Digital Audio Effects*, U. Zölzer, Ed., chapter 12, pp. 473–522. John Wiley & Sons, Ltd, 2011.
- [21] C. Darabundit, R. Wedelich, and P. Bischoff, "Digital grey box model of the Uni-Vibe effects pedal," in *Proc. 22nd Int. Conf. Digital Audio Effects (DAFX-19)*, Birmingham, UK, Sept. 2019.
- [22] T. Stockham, Jr., "High speed convolution and correlation," in *Proc. Spring Joint Comput. Conf.*, April 1966, p. pp. 229–233.
- [23] J. Bonada, X. Serra, X. Amatriain, and A. Loscos, "Spectral processing," in *DAFX: Digital Audio Effects*, U. Zölzer, Ed., chapter 10, pp. 393–445. John Wiley & Sons, Ltd, 2011.
- [24] B. Hayes, C. Saitis, and G. Fazekas, "Sinusoidal frequency estimation by gradient descent," in *Proc. IEEE Int. Conf. on Acoust. Speech Sig. Proces.*, Rhodes, Greece, 2023.
- [25] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proc. Int. Conf. Learning Representations*, San Diego, CA, 2015.
- [26] J. Sleep, "Small Stone Information [Online]," <http://generalguitargadgets.com/effects-projects/phase-shifters/small-stone-information/>, accessed 26/3/23.