



**HAL**  
open science

## Perceptual–Neural–Physical Sound Matching

Han Han, Vincent Lostanlen, Mathieu Lagrange

► **To cite this version:**

Han Han, Vincent Lostanlen, Mathieu Lagrange. Perceptual–Neural–Physical Sound Matching. Proceedings of the International Conference on Acoustics, Speech, and Signal Processing, Jun 2023, Rhodes Island, Greece, France. pp.1-5, 10.1109/ICASSP49357.2023.10095391 . hal-04027307

**HAL Id: hal-04027307**

**<https://hal.science/hal-04027307v1>**

Submitted on 13 Mar 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

# PERCEPTUAL-NEURAL-PHYSICAL SOUND MATCHING

Han Han, Vincent Lostanlen, and Mathieu Lagrange

Nantes Université, École Centrale Nantes, CNRS, LS2N, UMR 6004, F-44000 Nantes, France

## ABSTRACT

Sound matching algorithms seek to approximate a target waveform by parametric audio synthesis. Deep neural networks have achieved promising results in matching sustained harmonic tones. However, the task is more challenging when targets are nonstationary and inharmonic, e.g., percussion. We attribute this problem to the inadequacy of loss function. On one hand, mean square error in the parametric domain, known as “P-loss”, is simple and fast but fails to accommodate the differing perceptual significance of each parameter. On the other hand, mean square error in the spectrotemporal domain, known as “spectral loss”, is perceptually motivated and serves in differentiable digital signal processing (DDSP). Yet, spectral loss is a poor predictor of pitch intervals and its gradient may be computationally expensive; hence a slow convergence. Against this conundrum, we present Perceptual-Neural-Physical loss (PNP). PNP is the optimal quadratic approximation of spectral loss while being as fast as P-loss during training. We instantiate PNP with physical modeling synthesis as decoder and joint time–frequency scattering transform (JTFS) as spectral representation. We demonstrate its potential on matching synthetic drum sounds in comparison with other loss functions.

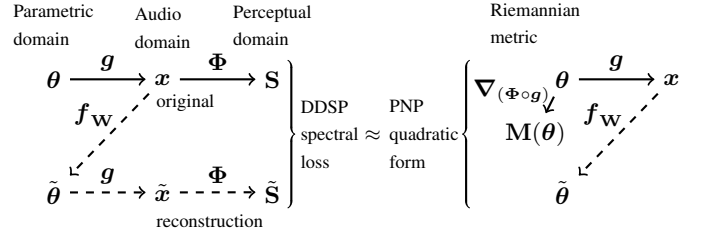
**Index Terms**— sound matching, auditory similarity, scattering transform, deep convolutional networks, physical modeling synthesis.

## 1. INTRODUCTION

Given an audio synthesizer  $g$ , the task of sound matching [1] consists in retrieving the parameter setting  $\theta$  that “matches” a target sound  $x$ ; i.e., such that a human ear judges the generated sound  $g(\theta)$  to resemble  $x$ . Sound matching has applications in automatic music transcription, virtual reality, and audio engineering [2, 3]. Of particular interest is the case where  $g(\theta)$  solves a known partial differential equation (PDE) whose coefficients are contained in the vector  $\theta$ . In this case,  $\theta$  reveals some key design choices in acoustical manufacturing, such as the shape and material properties of the resonator.

Over the past decade, the renewed interest for deep neural networks (DNN’s) in audio content analysis has led researchers to formulate sound matching as a supervised learning problem [4]. Intuitively, the goal is to optimize the synaptic weights  $\mathbf{W}$  of a DNN  $f_{\mathbf{W}}$  so that  $f_{\mathbf{W}}(x_n) = \tilde{\theta}_n$  approximates  $\theta_n$  over a training set of pairs  $(x_n, \theta_n)$ . Because  $g$  automates the mapping from parameter  $\theta_n$  to sound  $x_n$ , this training procedure incurs no real-world audio acquisition nor human annotation. However, prior publications have pointed out that the approximation formula  $\tilde{\theta}_n \approx \theta_n$  lacks a perceptual meaning: depending on the choice of target  $x_n$ , some deviations  $(\tilde{\theta}_n - \theta_n)$  may be judged to have a greater effect than others [5, 6, 7].

The paradigm of differentiable digital signal processing (DDSP) has brought a principled methodology to address this issue [8]. The key idea behind DDSP is to chain the learnable encoder  $f_{\mathbf{W}}$  with the known decoder  $g$  and a non-learnable but differentiable feature map  $\Phi$ . In DDSP,  $f_{\mathbf{W}}$  is trained to minimize the perceptual distance



**Fig. 1.** Graphical outline of the proposed method. Given a known synthesizer  $g$  and feature map  $\Phi$ , we train a neural network  $f_{\mathbf{W}}$  to estimate  $\tilde{\theta}$  and minimize the “perceptual–neural–physical” (PNP) quadratic form  $\langle \tilde{\theta} - \theta | \mathbf{M}(\theta) | \tilde{\theta} - \theta \rangle$  where  $\mathbf{M}$  is the Riemannian metric associated to  $(\Phi \circ g)$ . Hence, PNP approximates DDSP spectral loss yet does not need to backpropagate  $\nabla_{(\Phi \circ g)}(\tilde{\theta})$  at each epoch. Transformations in solid (resp. dashed) lines can (resp. cannot) be cached during training.

between vectors  $\Phi(\tilde{x}_n) = (\Phi \circ g \circ f_{\mathbf{W}})(x_n)$  and  $\Phi(x_n)$  on average over samples  $x_n$ . Yet, a practical shortcoming of DDSP is that it requires to backpropagate the “spectral loss”  $\|\Phi(\tilde{x}_n) - \Phi(x_n)\|_2$  over each DNN prediction  $\tilde{\theta}_n$ ; and so at every training step, since  $\mathbf{W}$  is iteratively updated by stochastic gradient descent (SGD).

In this article, we propose a new learning objective for sound matching, named perceptual–neural–physical (PNP). Our main contribution is to compute the Riemannian metric  $\mathbf{M}$  associated to the Jacobian  $\nabla_{(\Phi \circ g)}$  over each sample  $\theta_n$  (see Section 2.1). With  $\mathbf{M}(\theta_n)$ , we train  $f_{\mathbf{W}}$  to minimize a locally linear approximation of spectral loss, making PNP comparable to DDSP. Yet, unlike in DDSP, the computation of  $\nabla_{(\Phi \circ g)}$  is independent from the encoder  $f_{\mathbf{W}}$ : thus, it may be parallelized and cached during DNN training. A second novelty of our paper resides in its choice of application: namely, differentiable sound matching for percussion instruments. This requires not only a fine characterization of the spectral envelope, as in the DDSP of sustained tones; but also of attack and release transients. For this purpose, we need  $g$  and  $\Phi$  to accommodate sharp spectrotemporal modulations. Specifically, we rely on original differentiable implementations of the functional transformation method (FTM) for  $g$  and the joint time–frequency scattering transform (JTFS) for  $\Phi$ .<sup>1, 2</sup>

## 2. METHODS

### 2.1. Accelerating spectral loss with Riemannian geometry

We assume the synthesizer  $g$  and the feature map  $\Phi$  to be continuously differentiable. Let us denote by  $\mathcal{L}^{\text{DDSP}}$  the “spectral loss”

<sup>1</sup>Companion website: [https://github.com/lylyhan/perceptual\\_neural\\_physical](https://github.com/lylyhan/perceptual_neural_physical)

<sup>2</sup>Audio examples: <https://pnp.cargo.site/>

associated to the triplet  $(\Phi, \mathbf{f}_W, \mathbf{g})$ . Its value at a parameter set  $\theta$  is:

$$\begin{aligned} \mathcal{L}_\theta^{\text{DDSP}}(\mathbf{W}) &= \frac{1}{2} \|\Phi(\tilde{\mathbf{x}}) - \Phi(\mathbf{x})\|_2^2 \\ &= \frac{1}{2} \|(\Phi \circ \mathbf{g} \circ \mathbf{f}_W \circ \mathbf{g})(\theta) - (\Phi \circ \mathbf{g})(\theta)\|_2^2 \end{aligned} \quad (1)$$

by definition of  $\tilde{\mathbf{x}}$  and  $\mathbf{x}$ . Using  $\tilde{\theta}$  as shorthand for  $(\mathbf{f}_W \circ \mathbf{g})(\theta)$ , we conduct a first-order Taylor expansion of  $(\Phi \circ \mathbf{g})$  near  $\theta$ . We obtain:

$$\Phi(\tilde{\mathbf{x}}) = \Phi(\mathbf{x}) + \nabla_{(\Phi \circ \mathbf{g})}(\theta) \cdot (\tilde{\theta} - \theta) + O(\|\tilde{\theta} - \theta\|_2^2), \quad (2)$$

where the Jacobian matrix  $\nabla_{(\Phi \circ \mathbf{g})}(\theta)$  contains  $P = \dim \Phi(\mathbf{x})$  rows and  $J = \dim \theta$  columns. The manifold formed by differentiable map  $(\Phi \circ \mathbf{g})$  and the open set  $\theta \subset \mathbb{R}^J$  induces a Riemannian metric  $\mathbf{M}$ , i.e., an inner product on the tangent space at each point  $\theta$ :

$$\mathbf{M}(\theta)_{j,j'} = \sum_{p=1}^P (\nabla_{(\Phi \circ \mathbf{g})}(\theta)_{p,j}) (\nabla_{(\Phi \circ \mathbf{g})}(\theta)_{p,j'}). \quad (3)$$

The real-valued square matrix  $\mathbf{M}(\theta)$  defines a positive semidefinite kernel which, once plugged into Equation 2, serves to approximate  $\mathcal{L}_\theta^{\text{DDSP}}(\mathbf{W})$  in terms of a quadratic form over  $(\tilde{\theta} - \theta)$ :

$$\|\Phi(\tilde{\mathbf{x}}) - \Phi(\mathbf{x})\|_2^2 = \langle \tilde{\theta} - \theta | \mathbf{M}(\theta) | \tilde{\theta} - \theta \rangle + O(\|\tilde{\theta} - \theta\|_2^3). \quad (4)$$

The advantage of the approximation above is that the metric  $\mathbf{M}$  may be computed over the training set once and for all. This is because Equation 3 is independent of the encoder  $\mathbf{f}_W$ . Furthermore, since  $\theta$  is low-dimensional, we may store  $\mathbf{M}(\theta)$  on RAM. From this perspective, we define the perceptual–neural–physical loss (PNP) associated to  $(\Phi, \mathbf{f}_W, \mathbf{g})$  as the linearization of spectral loss at  $\theta$ :

$$\begin{aligned} \mathcal{L}_\theta^{\text{PNP}}(\mathbf{W}) &= \frac{1}{2} \langle (\mathbf{f}_W \circ \mathbf{g})(\theta) - \theta | \mathbf{M}(\theta) | (\mathbf{f}_W \circ \mathbf{g})(\theta) - \theta \rangle \\ &= \mathcal{L}_\theta^{\text{DDSP}}(\mathbf{W}) + O(\|(\mathbf{f}_W \circ \mathbf{g})(\theta) - \theta\|_2^3). \end{aligned} \quad (5)$$

According to the chain rule, the gradient of PNP loss at a given training pair  $(\mathbf{x}_n, \theta_n)$  with respect to some scalar weight  $\mathbf{W}_i$  is:

$$\frac{\partial \mathcal{L}_\theta^{\text{PNP}}}{\partial \mathbf{W}_i}(\theta_n) = \left\langle \mathbf{f}_W(\mathbf{x}_n) - \theta_n \middle| \mathbf{M}(\theta_n) \middle| \frac{\partial \mathbf{f}_W}{\partial \mathbf{W}_i}(\mathbf{x}_n) \right\rangle. \quad (6)$$

Observe that replacing  $\mathbf{M}(\theta_n)$  by the identity matrix in the equation above would give the gradient of *parameter loss* (P-loss); that is, the mean squared error between the predicted parameter  $\tilde{\theta}$  and the true parameter  $\theta$ . Hence, we may regard PNP as a perceptually motivated extension of P-loss, in which parameter deviations are locally recombined and rescaled so as to linearly approximate a DDSP objective.

The matrix  $\mathbf{M}(\theta)$  is constant in  $\mathbf{W}$ . Hence, its value may be cached across training epochs, and even across hyperparameter settings of the encoder. In comparison with P-loss, the only computational overhead of PNP is the bilinear form in Equation 6. However, this computation is performed in the parametric domain, i.e., in low dimension ( $J = \dim \theta$ ). Hence, its cost is negligible in front of the forward ( $\mathbf{f}_W$ ) and backward pass ( $\partial \mathbf{f}_W / \partial \mathbf{W}_i$ ) of DNN training.

## 2.2. Damped least squares

The principal components of the Jacobian  $\nabla_{(\Phi \circ \mathbf{g})}(\theta)$  are the eigenvectors of  $\mathbf{M}(\theta)$ . We denote them by  $\mathbf{v}_j$  and the corresponding eigenvalues by  $\sigma_j^2$ : for each of them, we have  $\mathbf{M}(\theta)\mathbf{v}_j = \sigma_j^2\mathbf{v}_j$ . The  $\mathbf{v}_j$ 's form an orthonormal basis of  $\mathbb{R}^J$ , in which we can decompose

the parameter deviation  $(\tilde{\theta} - \theta)$ . Recalling Equation 5, we obtain an alternative formula for PNP loss:

$$\mathcal{L}_\theta^{\text{PNP}}(\mathbf{W}) = \frac{1}{2} \sum_{j=1}^J \sigma_j^2 | \langle (\mathbf{f}_W \circ \mathbf{g})(\theta) - \theta | \mathbf{v}_j \rangle |^2 \quad (7)$$

The eigenvalues  $\sigma_j^2$  stretch and compress the error vector along their associated direction  $\mathbf{v}_j$ , analogous to the magnification and suppression of perceptually relevant and irrelevant parameter deviations. In practice however, when  $\sigma_j^2$  cover drastic ranges or contain zeros, as presented below in Section 4.3, the error vector is subject to extreme distortion and potential instability due to numerical precision errors. These scenarios, commonly referred to as  $\mathbf{M}$  being ill-conditioned, can lead to intractable learning objective  $\mathcal{L}_\theta^{\text{PNP}}$ .

Reminiscent of the damping mechanism introduced in Levenberg-Marquardt algorithm when solving nonlinear optimization problems, we update Equation 5 as

$$\mathcal{L}_\theta^{\text{PNP}}(\mathbf{W}) = \frac{1}{2} \langle \tilde{\theta} - \theta | \mathbf{M}(\theta) + \lambda I | \tilde{\theta} - \theta \rangle \quad (8)$$

The damping term  $\lambda I$  up-shifts all eigenvalues of  $\mathbf{M}$  by a constant positive amount  $\lambda$ , thereby changing its condition number. At the limit of  $\lambda \rightarrow 0$ ,  $\mathcal{L}_\theta^{\text{PNP}}$  reduces to a quadratic form which is asymptotically equivalent to spectral loss as P-loss approaches zero. At the limit of  $\lambda \rightarrow \infty$ ,  $\mathbf{M}$  is negligible in front of  $\lambda I$  thus  $\mathcal{L}_\theta^{\text{PNP}}$  boils down to P-loss. Alternatively, Equation 8 may also be viewed as a L2 regularization with coefficient  $\lambda$ .

To further address potential convergence issues,  $\lambda$  may be scheduled or adaptively changed according to epoch validation loss. We adopt delayed gratification mechanism to decrease  $\lambda$  by a factor of 5 when epoch validation loss is going down, and fix  $\lambda$  otherwise.

## 3. APPLICATION TO DRUM SOUND MATCHING

### 3.1. Perceptual: Joint time–frequency scattering (JTFS)

The joint time–frequency scattering transform (JTFS) is a nonlinear convolutional operator which extracts spectrotemporal modulations in the constant- $Q$  scalogram [9, 10]. Its kernels proceed from a separable product between two complex-valued wavelet filterbanks, defined over the time axis and over the log-frequency axis respectively. After convolution, we apply pointwise complex modulus and temporal averaging to each JTFS coefficient. These coefficients are known as scattering “paths”  $p$ . We apply a logarithmic transformation to the feature vector  $\text{JTFS}(\mathbf{x}_n)$  corresponding to each sound  $\mathbf{x}_n$ , yielding

$$\mathbf{S}_{n,p} = \Phi(\mathbf{x}_n)_p = (\Phi \circ \mathbf{g})(\theta_n)_p = \log \left( 1 + \frac{\text{JTFS}(\mathbf{x}_n)_p}{\varepsilon} \right), \quad (9)$$

We set  $\varepsilon = 10^{-3}$ , which is the order of magnitude of the median value of JTFS across all examples  $\mathbf{x}_n$  and paths  $p$ .

The multiresolution structure of JTFS is reminiscent of spectrotemporal receptive fields (STRF), and thus may serve as a biologically plausible predictor of neurophysiological responses in the primary auditory cortex [11]. At a higher level of music cognition, a recent study has shown that Euclidean distances in  $\Phi$  space predict auditory judgments of timbre similarity within a large vocabulary of instrumental playing techniques, as collected from a group of professional composers and non-expert music listeners [12].

We use the GPU implementation of [13] to compute JTFS with the same parameters as [12]:  $Q_1 = 12$ ,  $Q_2 = 1$ , and  $Q^{\text{tr}} = 1$  filters per octave respectively. We set the temporal averaging to  $T = 3$

seconds and the frequential averaging to  $F = 2$  octaves; hence a total of  $P = 20762$  paths. We refer to [14] for further details on the ability of  $\nabla_{(\Phi \circ g)}$  to extract “mesostructures” in nonstationary audio signals.

### 3.2. Neural: Deep convolutional network (convnet)

EfficientNet is a convolutional neural network architecture that balances the scaling of the depth, width and input resolution of consecutive convolutional blocks [15]. Achieving state-of-the-art performance on image classification with significantly less trainable parameters, its most light-weight version EfficientNet-B0 also succeeded in benchmarking audio classification tasks [16]. We adopt EfficientNet-B0 as our encoder  $\mathbf{f}_W$ , resulting in 4M learnable parameters. We append a linear dense layer of  $J = \dim \theta$  neurons and a 1D batch normalization before tanh activation. The goal of batch normalization is to gaussianize the input, such that the activated output is capable of uniformly cover the normalized prediction range. The input to  $\mathbf{f}_W$  is the log-scaled CQT coefficients of each example, spanning 10 octaves with 12 filters per octave.

### 3.3. Physical: Functional transformation method (FTM)

We are interested in the perpendicular displacement  $\mathbf{X}(t, \mathbf{u})$  on a rectangular drum face, which can be solved from the following partial differential equation defined in the Cartesian coordinate system  $\mathbf{u} = (u_1, u_2)$ .

$$\left( \frac{\partial^2 \mathbf{X}}{\partial t^2}(t, \mathbf{u}) - c^2 \nabla^2 \mathbf{X}(t, \mathbf{u}) \right) + S^4 (\nabla^4 \mathbf{X}(t, \mathbf{u})) + \frac{\partial}{\partial t} (d_1 \mathbf{X}(t, \mathbf{u}) + d_3 \nabla^2 \mathbf{X}(t, \mathbf{u})) = 0 \quad (10)$$

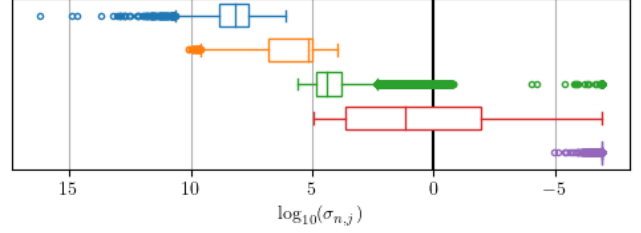
In addition to the standard traveling wave equation in the first above parenthesis, the fourth-order spatial and first-order time derivatives incorporate damping factors induced by stiffness, internal friction in the drum material and air friction in the external environment, rendering the solution a closer simulation to reality. Specifically,  $\alpha, S, c, d_1, d_3$  designate respectively the side length ratio, stiffness, traveling wave speed, frequency-independent damping and frequency-dependent damping of the drum. Even though real world drums are mostly circular, a rectangular drum model is equally capable of eliciting representative percussive sounds in real world scenarios. The circular drum model simply requires a conversion of Equation 10 into the Polar coordinate system. We bound the four sides of this  $l$  by  $l\alpha$  rectangular drum at zero at all time. For simplicity, we simulate the excitation by setting the initial condition at  $t_0 = 0$  to be  $\mathbf{X}(t_0, \mathbf{u} = (0.4l, 0.4l\alpha)) = 0.03$  meters, and 0 otherwise.

We implement generator  $\mathbf{g}$  as a PDE solver to this high-order damped wave equation, namely the functional transformation method (FTM) [17, 18]. FTM solves the PDE by transforming the equation into its Laplace and functional space domain, where an algebraic solution can be obtained. It then finds the time-space domain solution via inverse functional transforms, expressed in an infinite modal summation form

$$\mathbf{x}(t) = \mathbf{X}(t, \mathbf{u}) = \sum_{m \in \mathbb{N}^2} K_m(\mathbf{u}, t) \exp(\sigma_m t) \sin(\omega_m t) \quad (11)$$

The coefficients  $K_m(\mathbf{u}, t)$ ,  $\sigma_m$ ,  $\omega_m$  are derived from the original PDE parameters in the following ways.

$$\omega_m^2 = (S^4 - \frac{d_3^2}{4}) \Gamma_{m_1, m_2}^2 + (c^2 + \frac{d_1 d_3}{2}) \Gamma_{m_1, m_2} - \frac{d_1^2}{4} \quad (12)$$



**Fig. 2.** Distributions of the five sorted eigenvalues of  $\mathbf{M}(\theta_n)$ . For the sake of comparison between PNP and P-loss, the bold line indicates the eigenvalues of the identity matrix (see Equation 6).

$$\sigma_m = \frac{d_3}{2} \Gamma_{m_1, m_2} - \frac{d_1}{2} \quad (13)$$

$$K_m(\mathbf{u}, t) = y_u^m \delta(t) \sin\left(\frac{\pi m_1 u_1}{l}\right) \sin\left(\frac{\pi m_2 u_2}{l\alpha}\right) \quad (14)$$

where  $\Gamma_{m_1, m_2} = \pi^2 m_1^2 / l^2 + \pi^2 m_2^2 / (l\alpha)^2$ , and  $y_u^m$  is the  $m^{th}$  coefficient associated to the eigenfunction  $\sin(\pi m \mathbf{u} / l)$  that decomposes  $y_u(\mathbf{u})$ .

Without losing connections to the acoustical manufacturing of the drum yet better relating  $\mathbf{g}$ 's input with perceptual dimensions, we reparametrize the PDE parameters  $\{S, c, d_1, d_3, \alpha\}$  into  $\theta = \{\log \omega_1, \tau_1, \log p, \log D, \alpha\}$ , detailed in Section 3.4 of [19]. We prescribe sonically-plausible ranges for each parameter in  $\theta$ , normalize them between  $-1$  and  $1$ , uniformly sample in the hyper-dimensional cube, and obtain a dataset of 100k percussive sounds sampled at 22050 HZ. The train/test/validation split is  $8 : 1 : 1$ .

In particular, fundamental frequency  $\omega_1$ , duration  $\tau_1$  falls into ranges  $[40, 1000]$  Hz and  $[0.4, 3]$  seconds respectively. Inhomogeneous damping rate  $p$ , frequential dispersion  $D$  and aspect ratio  $\alpha$  ranges are  $[10^{-5}, 0.2]$ ,  $[10^{-5}, 0.3]$ , and  $[10^{-5}, 1]$ .

## 4. RESULTS

### 4.1. Baselines

We train  $f_W$  with 3 different losses - multi-scale spectral loss [20], parameter loss, and PNP loss. We use a batch size of 64 samples for spectral loss, and 256 samples for P-loss and PNP loss. The training proceeds for 70 epochs, where around 20% of the training set is seen at each epoch. We use Adam optimizer with learning rate  $10^{-3}$ . Table 1 reports the training time per epoch on a single Tesla V100 16GB GPU.

### 4.2. Evaluation with JTFS-based spectral loss

We propose to use the L2 norm of JTFS coefficients error averaged over test set for evaluation. As a point of reference, we also include the average multi-scale spectral error, implemented as in Section 4.1. One of the key distinctions between Euclidean JTFS distance and multi-scale spectral error is the former's inclusion of spectro-temporal modulations information. Meanwhile unlike mean squared parameter error, both metrics reflect the perceptual closeness instead of parametric retrieval accuracy for each proposed model.

### 4.3. Discussion

Despite being the optimal quadratic approximation of spectral loss, it is nontrivial to apply the bare PNP loss form as Equation 5 in

Loss	$\Phi$	Pitch	JTFS distance (avg. on test set)	MSS (avg. on test set)	Training time per epoch
P-loss	—	Known	<b>22.23</b> $\pm$ 2.17	<b>0.31</b> $\pm$ 0.013	49 minutes
$\mathcal{L}_\theta^{\text{DDSP}}$	$\Phi_{\text{MSS}}$	Known	31.86 $\pm$ 0.332	0.335 $\pm$ 0.005	54 minutes
$\mathcal{L}_\theta^{\text{PNP}}$	$\Phi_{\text{JTFS}}$	Known	23.58 $\pm$ 0.877	0.335 $\pm$ 0.005	49 minutes
$\mathcal{L}_\theta^{\text{DDSP}}$	$\Phi_{\text{JTFS}}$	—	—	—	est., > 1 day
P-loss	—	Unknown	61.91 $\pm$ 6.26	1.02 $\pm$ 0.094	53 minutes
$\mathcal{L}_\theta^{\text{DDSP}}$	$\Phi_{\text{MSS}}$	Unknown	138.95 $\pm$ 37.12	1.59 $\pm$ 0.307	59 minutes
$\mathcal{L}_\theta^{\text{PNP}}$	$\Phi_{\text{JTFS}}$	Unknown	<b>61.21</b> $\pm$ 1.207	<b>0.97</b> $\pm$ 0.019	49 minutes

**Table 1.** Report of average JTFS distance and MSS metrics evaluated on test set. Six models are trained with two modalities: 1. the inclusion of pitch retrieval i.e. regressing  $\theta = \{\tau, \log p, \log D, \alpha\}$  vs.  $\theta = \{\log \omega_1, \tau, \log p, \log D, \alpha\}$ , and 2. the choice of loss function: P-loss, MSS loss, or PNP loss with adaptive damping mechanism. The best performing models with known and unknown pitch are P-loss and PNP loss respectively. Training with MSS loss is more time consuming than training with P-loss or PNP loss. Training with differentiable JTFS loss is unrealistic in the interest of time.

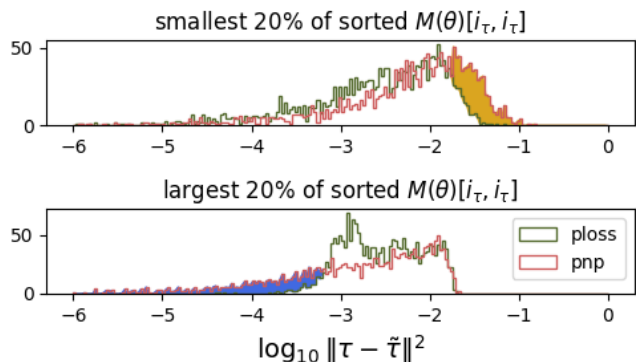
experimental settings. On one hand,  $\Phi \circ g$  potentially has undesirable property that exposes the Riemannian metric calculations to numerical precision errors. On the other hand, extreme deformation of the optimization landscape may lead to the same numerical instability facing stochastic gradient descent with spectral loss. We report on a few remedies that helped stabilize learning with PNP loss, and offer insights on future directions to take.

First and foremost, our preliminary experiments show that training PNP loss without damping  $\lambda = 0$  subjects to convergence issues due to the high condition numbers in empirical Ms as illustrated in Section 2.2. Fig. 2 shows the sorted eigenvalue distribution of all Ms in test set, where Ms are rank-2,3 or 4 matrices with eigenvalues ranging from 0 to  $10^{20}$ . This could be an implication that entries of  $\theta$  contain implicit linear dependencies in generator  $g$ , or that local variations of certain  $\theta$  fail to linearize differences in the output of  $g$  or  $\Phi \circ g$ . As an example, the aspect ratio  $\alpha$  influences the modal frequencies and decay rates via [19, Equations 12–13], where in fact  $(1/\alpha + 1/\alpha^2)$  could be a better choice of variable that linearizes  $g$ .

To address Ms’ ill conditions we attempted at numerous damping mechanisms to update  $\lambda$ : namely, constant  $\lambda$ , scheduled  $\lambda$  decay, and adaptive  $\lambda$  decay. The intuition is to have  $\mathcal{L}_\theta^{\text{PNP}}$  start in the parameter loss regime and move towards the spectral loss regime while training. The best performing model is achieved with adaptive  $\lambda$  decay (see Section 2.2). We propose to divide  $\lambda$  by a factor of 5 if the model breaks the best epoch validation loss record and keep it the same otherwise. In practice, we initialize  $\lambda$  to match the largest empirical  $\sigma_j^2 \approx 10^{20}$ , and then adaptively decay it to  $3 \times 10^{14}$  in 20 epochs. This indicates that  $f_{\mathbf{w}}$  is able to learn with damped PNP loss if  $\lambda$  is large enough to compensate for rank deficiency in M.

The diagonal elements of  $M(\theta)$  can be regarded as both the applied weights’ magnitudes and proxies for the perceptual importance of  $\theta$ ’s accuracy. Inspecting the results of  $\tau = \theta[i_\tau]$  regression, we observe in Fig. 3 that in comparison with P-loss model, PNP model improves retrieval accuracy for sounds inducing larger perceptual difference to changes in  $\tau$  (in blue), at the expense of lowered accuracy for the opposite (in yellow). This suggests a trade-off behavior aligned with PNP loss’ weighting scheme.

We believe that more of PNP loss’ mathematical potential can be exploited in the future, notably in cases where parameterization without domain-specific knowledge renders the failure of P-loss, and its use in hybrid optimization schemes. We plan to investigate the scalability of each loss function under reparameterizations, as well as other damping schemes and optimizers. The current update mechanism, originated from the Levenberg-Marquardt algorithm, aims to improve the conditioning of a matrix inversion problem in the



**Fig. 3.** Histogram of log squared  $\tau$  estimation error for perceptually significant (largest 20%  $M(\theta)[i_\tau, i_\tau]$ ) and less significant (smallest 20%  $M(\theta)[i_\tau, i_\tau]$ ) sounds. The yellow and blue regions indicate the weighting-induced difference in retrieval accuracy.

Gauss-Newton algorithm. However when used jointly with stochastic gradient descent, each  $\lambda$  update may change the optimization landscape drastically. The resulting optimization behavior is thus not fully understood. We consider interfacing nonlinear least squares solver with SGD and forming a hybrid learning scheme in future work.

## 5. CONCLUSION

Knowledge on human auditory perception aiding data-driven approaches to machine listening tasks have been exemplified in a multitude of applications [21]. In this article we presented another case of this synergy named Perceptual-Neural-Physical (PNP) autoencoding, a bilinear form learning objective for sound matching task. In our application, PNP optimizes the retrieval of physical parameters from sounds in a perceptually-motivated metric space, enabled by differentiable implementations of physical model and computational proxy of neurophysiological construct of human auditory system.

We demonstrated PNP’s mathematical relationship to spectral loss and parameter loss. Using this formulation, we motivated and established one way of interpolating between optimizing in parameter and spectral loss regimes. We presented damping mechanisms to facilitate its learning under ill-conditioned empirical settings and provided future plans for further exploiting its mathematical potential.

## 6. REFERENCES

- [1] Andrew Horner, “Wavetable matching synthesis of dynamic instruments with genetic algorithms,” *Journal of the Audio Engineering Society*, vol. 43, no. 11, pp. 916–931, 1995.
- [2] Jordie Shier, Kirk McNally, George Tzanetakis, and Ky Grace Brooks, “Manifold learning methods for visualization and browsing of drum machine samples,” *Journal of the Audio Engineering Society*, vol. 69, no. 1/2, pp. 40–53, 2021.
- [3] Philippe Esling, Naotake Masuda, Adrien Bardet, Romeo Depres, Axel Chemla, et al., “Universal audio synthesizer control with normalizing flows,” in *Proceedings of the International Conference on Digital Audio Effects (DAFX)*, 2019.
- [4] Leonardo Gabrielli, Stefano Tomassetti, Carlo Zinato, and Francesco Piazza, “End-to-end learning for physics-based acoustic modeling,” *IEEE Transactions on Emerging Topics in Computational Intelligence*, vol. 2, no. 2, pp. 160–170, 2018.
- [5] Naotake Masuda and Daisuke Saito, “Synthesizer sound matching with differentiable DSP,” in *Proceedings of the International Society on Music Information Retrieval (ISMIR) Conference*, 2021, pp. 428–434.
- [6] Martin Roth and Matthew Yee-king, “A comparison of parametric optimization techniques for musical instrument tone matching,” *Journal of the Audio Engineering Society*, May 2011.
- [7] Matthew Yee-King, Leon Fedden, and Mark d’Inverno, “Automatic programming of vst sound synthesizers using deep networks and other techniques,” *IEEE Transactions on Emerging Topics in Computational Intelligence*, vol. 2, pp. 150–159, 2018.
- [8] Jesse Engel, Lamtharn (Hanoi) Hantrakul, Chenjie Gu, and Adam Roberts, “DDSP: Differentiable Digital Signal Processing,” in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020.
- [9] Joakim Andén, Vincent Lostanlen, and Stéphane Mallat, “Joint time–frequency scattering,” *IEEE Transactions on Signal Processing*, vol. 67, no. 14, pp. 3704–3718, 2019.
- [10] Mathieu Andreux, Tomás Angles, Georgios Exarchakis, Roberto Leonarduzzi, Gaspar Rochette, Louis Thiry, John Zarka, Stéphane Mallat, Joakim Andén, Eugene Belilovsky, Joan Bruna, Vincent Lostanlen, Muawiz Chaudhary, Matthew J. Hirn, Edouard Oyallon, Sixin Zhang, Carmine Cella, and Michael Eickenberg, “Kymatio: Scattering transforms in Python,” *Journal of Machine Learning Research*, vol. 21, no. 60, pp. 1–6, 2020.
- [11] Taishih Chi, Powen Ru, and Shihab A Shamma, “Multiresolution spectrotemporal analysis of complex sounds,” *The Journal of the Acoustical Society of America*, vol. 118, no. 2, pp. 887–906, 2005.
- [12] Vincent Lostanlen, Christian El-Hajj, Mathias Rossignol, Grégoire Lafay, Joakim Andén, and Mathieu Lagrange, “Time–frequency scattering accurately models auditory similarities between instrumental playing techniques,” *EURASIP Journal on Audio, Speech, and Music Processing*, vol. 2021, no. 1, pp. 1–21, 2021.
- [13] John Muradeli, Cyrus Vahidi, Changhong Wang, Han Han, Vincent Lostanlen, Mathieu Lagrange, and George Fazekas, “Differentiable time-frequency scattering in kymatio,” in *Proceedings of the International Conference on Digital Audio Effects (DAFX)*, 2022.
- [14] Cyrus Vahidi, Han Han, Changhong Wang, Mathieu Lagrange, György Fazekas, and Vincent Lostanlen, “Mesostructures: Beyond spectrogram loss in differentiable time-frequency analysis,” *arXiv preprint arXiv:2301.10183*, 2023.
- [15] Mingxing Tan and Quoc Le, “EfficientNet: Rethinking model scaling for convolutional neural networks,” in *Proceedings of the International conference on Machine Learning (ICML)*. PMLR, 2019, pp. 6105–6114.
- [16] Neil Zeghidour, Olivier Teboul, Félix de Chaumont Quitry, and Marco Tagliasacchi, “LEAF: A learnable frontend for audio classification,” *ICLR*, 2021.
- [17] L. Trautmann and Rudolf Rabenstein, *Digital Sound Synthesis by Physical Modeling Using the Functional Transformation Method*, Springer, 2003.
- [18] M. Schäfer, M. Werner, and R. Rabenstein, “Physical modeling in sound synthesis: Vibrating plates,” in *Proc. 26th International Congress on Sound and Vibration (ICSV26)*, Montreal, Canada, Jul. 2019, pp. 1–8.
- [19] Han Han and Vincent Lostanlen, “wav2shape: Hearing the Shape of a Drum Machine,” in *Proceedings of Forum Acusticum*, 2020, pp. 647–654.
- [20] Christian J. Steinmetz and Joshua D. Reiss, “auraloss: Audio focused loss functions in PyTorch,” in *Digital Music Research Network One-day Workshop (DMRN+15)*, 2020.
- [21] Laurie M. Heller, Benjamin Elizalde, Bhiksha Raj, and Soham Deshmukh, “Synergy between human and machine approaches to sound/scene recognition and processing: An overview of ICASSP special session,” in *Proceedings of the IEEE International Conference on Audio, Speech, and Signal Processing*. 2023, IEEE.