

# SCRAPL: SCATTERING TRANSFORM WITH RANDOM PATHS FOR MACHINE LEARNING

Anonymous authors

Paper under double-blind review

## ABSTRACT

The Euclidean distance between wavelet scattering transform coefficients (known as *paths*) provides informative gradients for perceptual quality assessment of deep inverse problems in computer vision, speech, and audio processing. However, these transforms are computationally expensive when employed as differentiable loss functions for stochastic gradient descent due to their numerous paths, which significantly limits their use in neural network training. Against this problem, we propose “Scattering transform with **R**andom **P**aths for machine **L**earning” (SCRAPL): a stochastic optimization scheme for efficient evaluation of multivariable scattering transforms. We implement SCRAPL for the joint time–frequency scattering transform (JTFS) which demodulates spectrotemporal patterns at multiple scales and rates, allowing a fine characterization of intermittent auditory textures. We apply SCRAPL to differentiable digital signal processing (DDSP), specifically, unsupervised sound matching of a granular synthesizer and the Roland TR-808 drum machine. We also propose an initialization heuristic based on importance sampling, which adapts SCRAPL to the perceptual content of the dataset, improving neural network convergence and evaluation performance. We make our audio samples available and provide SCRAPL as a Python package.

## 1 INTRODUCTION

A scattering transform (ST) is a wavelet-based nonlinear operator which decomposes a high-resolution input  $x$  into a collection  $\Phi x$  of low-resolution coefficients, known as *paths* (Mallat, 2012). Without loss of generality, let us consider a two-layer multivariable ST of a time-domain signal  $x(t)$ :

$$\Phi x(p, t, \lambda) = \rho \left( \left( |\mathbf{W}x| \circledast \Psi_p \right) (t, \lambda) \right). \quad (1)$$

In the equation above,  $\mathbf{W}$  is a wavelet transform; the vertical bars denote complex modulus; the circled asterisk  $\circledast$  denotes a multivariable convolution over time  $t$  and wavelet scale  $\lambda$ ;  $\Psi$  is a multivariable wavelet filterbank which is indexed by path  $p$ ;  $\Psi_0$ , i.e.,  $\Psi_p$  with  $p = 0$  is a multivariable low-pass filter; and  $\rho$  is a pointwise nonlinearity, e.g., path normalization and logarithmic transformation.

The design of the filterbank  $\Psi$  aims at a tradeoff between three properties: invariance to rigid motion, stability to small deformations, and separation of sparse patterns (Mallat, 2016). In speech and audio processing, examples of such  $\Psi$  include “plain” time ST (Andén & Mallat, 2014); joint time–frequency scattering (JTFS) (Andén & Mallat, 2014); and spiral ST (Lostanlen & Mallat, 2016). In computer vision, examples include “plain” 2-D ST (Bruna & Mallat, 2013); joint roto-translation ST Sifre & Mallat (2013); and scalo-roto-translation ST (Oyallon et al., 2014).

The squared Euclidean distance between scattering coefficients, or *ST distance* for short, is:

$$d_\Phi(x, \tilde{x}) = \sum_{p=0}^{P-1} \sum_{t=0}^{T-1} \sum_{\lambda=0}^{\Lambda-1} \left| \Phi x(p, t, \lambda) - \Phi \tilde{x}(p, t, \lambda) \right|^2, \quad (2)$$

where  $P$  is the number of paths;  $T$  is the number of time samples; and  $\Lambda$  is the number of scales. Behavioral studies suggest that ST distance is a good predictor of dissimilarity judgments between isolated sounds, for suitably chosen  $\Psi$  and  $\rho$  (Patil et al., 2012; Lostanlen et al., 2021; Tian et al., 2025). Relatedly, neurophysiology studies suggest that JTFS is a suitable idealized model of

spectrotemporal receptive fields in the auditory cortex of humans (Norman-Haignere & McDermott, 2018) and nonhuman mammals (Kowalski et al., 1996). These findings motivate the use of JTFS as part of a differentiable loss function for neural audio models (Vahidi et al., 2023).

As an illustration, let  $\mathbf{x}$  be a fixed reference and  $\tilde{\mathbf{x}} = F_{\mathbf{x}}(\mathbf{w})$  be its reconstruction by an autoencoder  $F$  with trainable weights  $\mathbf{w}$ . Denoting the set of path indices by  $\mathcal{P} = \{0, \dots, P-1\}$  and the vector of all time–frequency entries  $\Phi\mathbf{x}(p, t, \lambda)$  for each path  $p \in \mathcal{P}$  by  $\phi_p(\mathbf{x})$ , the ST loss function writes as:

$$\mathcal{L}_{\mathbf{x}}^{\Phi}(\tilde{\mathbf{x}}) = \frac{1}{P} \sum_{p=0}^{P-1} \mathcal{L}_{\mathbf{x}}^{\phi_p}(\tilde{\mathbf{x}}) \quad \text{where} \quad \forall p \in \mathcal{P}, \mathcal{L}_{\mathbf{x}}^{\phi_p}(\tilde{\mathbf{x}}) = P \|\phi_p(\mathbf{x}) - \phi_p(\tilde{\mathbf{x}})\|^2. \quad (3)$$

Unfortunately,  $\mathcal{L}_{\mathbf{x}}^{\Phi}(\tilde{\mathbf{x}})$  and its gradient  $\nabla \mathcal{L}_{\mathbf{x}}^{\Phi}(\tilde{\mathbf{x}})$  are expensive in memory and in operations. Certainly, algorithmic refinements such as FFT-based filtering, multirate processing, and depth-first search can reduce the cost of an ST path (Oyallon et al., 2018). Yet, the need to traverse all  $P$  paths remains an obstacle to the applicability of multivariable ST for gradient-based learning at scale.

In this article, we aim to accelerate the training of an autoencoder  $F$  whose loss is ST distance between reference and reconstruction, and so over a finite corpus  $\mathcal{X} = \{\mathbf{x}_0, \dots, \mathbf{x}_{N-1}\}$ . Formally:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{N} \sum_{n=0}^{N-1} (\mathcal{L}_{\mathbf{x}_n}^{\Phi} \circ F_{\mathbf{x}_n})(\mathbf{w}). \quad (4)$$

Given the decomposition in Equation 3, a naïve idea would be to replace each term  $\mathcal{L}_{\mathbf{x}_n}^{\Phi}$  in the equation above by some per-path loss  $\mathcal{L}_{\mathbf{x}_n}^{\phi_p}$ , where the  $p$ ’s would be drawn independently and uniformly at random in the path set  $\mathcal{P}$ . This is a crude form of stochastic approximation (Benveniste et al., 2012) which is motivated by the tree-like structure of ST: neglecting the overhead of the first layer ( $|\mathbf{W}\mathbf{x}|$ ), the computation of single-path gradient  $\nabla \mathcal{L}_{\mathbf{x}_n}^{\phi_p}$  is roughly  $P$  times more efficient than that of a full ST gradient  $\nabla \mathcal{L}_{\mathbf{x}_n}^{\Phi}$ . However, this speedup comes at the detriment of numerical precision: a deterministic quantity has been replaced by an estimator whose variance may be impractically large.

“Scattering transform with **R**andom **P**aths for machine **L**earning” (SCRAPL) is our proposed solution to this problem. Acknowledging that each single-path gradient makes for an inexpensive but noisy learning signal, we stabilize it via a combination of three stochastic optimization techniques and apply an architecture-informed importance sampling heuristic. Our contributions are:

1. **Stochastic approximation of scattering transform** through uniform sampling of paths.
2. **Path-wise adaptive moment estimation** ( $\mathcal{P}$ -Adam for short): an extension of the Adam algorithm (Kingma & Ba, 2014) which accounts for the non-i.i.d. nature of ST paths.
3. **Path-wise stochastic average gradient with acceleration** ( $\mathcal{P}$ -SAGA for short): a variant of the SAGA algorithm (Defazio et al., 2014) which keeps a memory of previous gradient values across all paths  $p$ .
4.  **$\theta$ -importance sampling**: a parallelizable initialization heuristic that supplies auxiliary information to the stochastic optimizer by sampling paths  $p$  in proportion to the typical rate of change of the gradient in the optimization landscape.

Our main empirical finding is that SCRAPL accomplishes a favorable tradeoff between goodness of fit and computational efficiency on unsupervised sound matching, i.e., a nonlinear inverse problem in which the forward operator implements an audio synthesizer. In the context of differentiable digital signal processing (DDSP), the state-of-the-art perceptual loss function for this task is multiscale spectral loss (MSS, Yamamoto et al. (2020), Engel et al. (2020)). However, the gradient of MSS is uninformative when input and reconstruction are misaligned or when the synthesizer controls involve spectrotemporal modulations (Vahidi et al., 2023). Taking advantage of the stability guarantees of JTFS, SCRAPL expands the class of synthesizers which can be effectively decoded via DDSP.

Figure 1 illustrates one of our experiments: unsupervised sound matching for a nondeterministic granular synthesizer. On one hand, models based on MSS and other state-of-the-art perceptual losses are computationally efficient but inaccurate. On the other hand, JTFS-based models are five times more accurate but twenty five times more costly. SCRAPL is a new point on this Pareto front: it is within a factor two of JTFS in terms of accuracy while being within a factor two of MSS in terms of runtime, making it suitable for large-scale DDSP. Relatedly, SCRAPL is also more memory-efficient than JTFS, thus reducing overhead between CPU/GPU cores and allowing for a larger batch size.

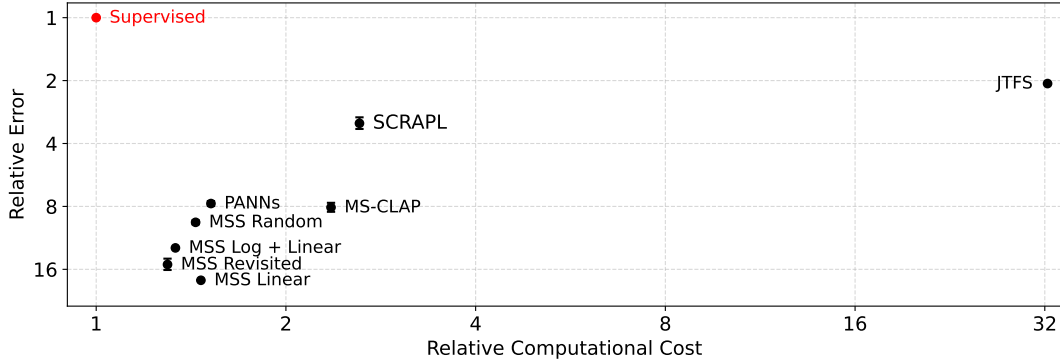


Figure 1: Mean average synthesizer parameter error (y-axis) versus computational cost (x-axis) of unsupervised sound matching models for the granular synthesis task. Both axes are rescaled by the performance of a supervised model with the same number of parameters. Whiskers denote 95% CI, estimated over 20 random seeds. Due to computational limitations, JTFS-based sound matching is evaluated only once.

## 2 RELATED WORK

The guiding intuition behind SCRAPL is that natural signals and images exhibit strong correlations across ST paths. This fact has been observed empirically since the onset of ST research (Bruna & Mallat, 2013; Andén & Mallat, 2011) and aligns with earlier work on texture modeling based on pairwise correlations between wavelet modulus coefficients (Portilla & Simoncelli, 2000).

Visual and auditory textures, understood as stationary random fields, play a key role in applied ST research. ST features outperform short-term Fourier features (e.g., MSS) in their ability to characterize intermittency in non-Gaussian textures (Muzy et al., 2015). Texture resynthesis by gradient descent of ST loss has been applied to such diverse settings as computer music creation (Lostanlen et al., 2019) and the study of the cosmic microwave background (Delouis et al., 2022).

The democratization of differentiable programming toolkits (e.g., TensorFlow, PyTorch, JAX) has greatly advanced the flexibility of gradient backpropagation in “hybrid” scattering–neural networks involving learnable and non-learnable modules. Angles & Mallat (2018) have built a hybrid scattering–GAN model for image generation, in which ST distance plays the role of a discriminator.

To our knowledge, the closest prior work to SCRAPL is the pruned graph scattering transform (pGST) of Ioannidis et al. (2020), a method which reduces the complexity of ST by pruning the path set  $\mathcal{P}$  down to a proper subset  $\mathcal{P}' \subset \mathcal{P}$ , based on a graph-spectrum-inspired criterion. Although both pGST and SCRAPL share a similar overarching goal, let us point out that pGST is a feature selection method: the cardinality of  $\mathcal{P}'$  is typically  $\sim 10\%$  that of  $\mathcal{P}$  and  $\mathcal{P}'$  is kept fixed across training examples and across epochs. In comparison, SCRAPL performs a more radical pruning, down to a single path ( $\text{card } \mathcal{P}' = 1$ ), while harnessing dedicated techniques in stochastic optimization ( $\mathcal{P}$ -Adam and  $\mathcal{P}$ -SAGA) to reduce the variance of ST loss during gradient backpropagation.

## 3 METHODS

### 3.1 STOCHASTIC APPROXIMATION OF SCATTERING TRANSFORM LOSS GRADIENT

The proposition below, proven in Appendix A, shows that if paths are drawn uniformly at random, the stochastic approximation in SCRAPL is unbiased: in other words, the expected value of the stochastic gradient of per-path loss is equal to the gradient of full ST loss.

**Proposition 3.1.** Let  $\Phi = (\phi_p)_0^{P-1}$  be a scattering transform with  $P$  paths. Given a signal or image  $x$ , let  $F_x$  be an autoencoder operating on  $x$  and let  $\mathcal{L}_x^\Phi$  be the associated ST reconstruction loss. Let  $\mathcal{U}_P$  be the uniform distribution over  $\mathcal{P} = \{0, \dots, P-1\}$ . One has, for every weight vector  $w$ :

$$\mathbb{E}_{z \sim \mathcal{U}_P} [\nabla(\mathcal{L}_x^{\phi_z} \circ F_x)(w)] = \nabla(\mathcal{L}_x^\Phi \circ F_x)(w). \quad (5)$$

Although a uniform sampling of paths matches the intuition of approximating the ST gradient in expectation, we will see that this may be suboptimal. The  $\theta$ -importance sampling method, which we will present in Section 3.4, does not satisfy the hypothesis of Proposition 3.1; yet, it consistently outperforms uniform sampling as part of SCRAPL. The design of biased stochastic approximation schemes is an active topic in machine learning research (Dieuleveut et al., 2023).

### 3.2 $\mathcal{P}$ -ADAM: PATH-WISE ADAPTIVE MOMENT ESTIMATION

The key idea behind the Adam optimizer is to smooth the successive realizations of the stochastic gradient, here denoted by  $g$ , via autoregressive estimates of its first- and second-order element-wise moments, denoted by  $m$  and  $v$  (Kingma & Ba, 2014). However, the smoothing technique in Adam is ineffective for SCRAPL because the gradients of path-wise ST losses are not identically distributed. Against this problem, we propose to maintain  $P$  estimates of path-wise moments ( $\mathcal{P}$ -Adam):

$$m_p \leftarrow \beta_1^{(k-\tau_p)/P} m_p + (1 - \beta_1^{(k-\tau_p)/P}) g \quad (6)$$

$$v_p \leftarrow \beta_2^{(k-\tau_p)/P} v_p + (1 - \beta_2^{(k-\tau_p)/P}) (g \odot g), \quad (7)$$

where  $k$  is the current iteration number,  $\tau_p$  is the iteration when path  $p$  was last drawn;  $\beta_1$  and  $\beta_2$  are hyperparameters; and the circled dot denotes element-wise multiplication of vectors. The exponent  $(k - \tau_p)/P$  adapts the time constant of smoothing to the recency of the previous estimate.

The second step in  $\mathcal{P}$ -Adam, following classical Adam, consists in bias correction and ratio of debiased first-order moment to stable square root of debiased second-order moment:

$$g_{\text{current}} = \frac{\frac{m_p}{1 - \beta_1^{k/P}}}{\sqrt{\varepsilon + \frac{v_p}{1 - \beta_2^{k/P}}}}, \quad (8)$$

where we have adapted the original exponents of Adam ( $\beta_1^k, \beta_2^k$ ) to account for the number of paths.

### 3.3 $\mathcal{P}$ -SAGA: PATH-WISE STOCHASTIC AVERAGE GRADIENT WITH ACCELERATION

The stochastic average gradient (SAG) algorithm has the potential to accelerate stochastic gradient descent in the context of the minimization of finite sums (Schmidt et al., 2017). Although this sum is typically over training examples in neural network training, in SCRAPL, Equation 3 is a sum over paths for a given example  $x$ . With this observation in mind, we propose  $\mathcal{P}$ -SAGA, a path-wise version of SAG with acceleration (SAGA, Defazio et al. (2014)). We maintain a memory of the last  $\mathcal{P}$ -Adam updates over each path, denoted by  $(\hat{g}_p)_0^{P-1}$ ; and the set of paths previously visited, denoted by  $\Gamma$ . Given a learning rate  $\alpha_k$  at iteration  $k$ , the  $\mathcal{P}$ -SAGA update is:

$$w \leftarrow w - \alpha_k \left( g_{\text{current}} - \hat{g}_p + \frac{\sum_{\gamma \in \Gamma} \hat{g}_\gamma}{\max(1, \text{card}(\Gamma))} \right). \quad (9)$$

Unlike the original SAG and SAGA algorithms,  $\mathcal{P}$ -SAGA’s additional memory footprint is proportional to  $P$ , not the size of the dataset  $N$ , making it suitable for neural network training and real-world optimization tasks like our experiments in Section 4. We also note that  $\mathcal{P}$ -Adam and  $\mathcal{P}$ -SAGA introduce no additional hyperparameters over the standard Adam optimizer. Algorithm 1 summarizes SCRAPL with both  $\mathcal{P}$ -Adam and  $\mathcal{P}$ -SAGA enabled.

### 3.4 $\theta$ -IMPORTANCE SAMPLING

We now consider the important special case of differentiable digital signal processing (DDSP, see Section 1), in which the autoencoder composes a non-learnable decoder, typically a synthesizer, with a learned encoder: i.e.,  $F_x = (D \circ E_x)$  with  $E_x(w) = \tilde{\theta}$  and  $D(\tilde{\theta}) = \tilde{x}$  (Engel et al., 2020). We assume both  $D$  and  $E_x$  to be differentiable with respect to their inputs, but  $D$  is not necessarily deterministic. We denote by  $U$  the dimension of the parameter space  $\theta$ ; i.e., the output space of  $E_x$  and input space of  $D$ .

**Algorithm 1** “Scattering transform with **R**andom **P**aths for machine **L**earning” (SCRAPL). The pseudo-code below describes SCRAPL with a batch size equal to one, without loss of generality.

---

**Require:**  $\Phi = (\phi_p)_0^{P-1}$ : Scattering transform (ST) with  $P$  paths  
**Require:**  $\pi$ : Categorical distribution over the path set  $\mathcal{P} = \{0, \dots, P-1\}$   
**Require:**  $F$ : Autoencoder with trainable parameters  $w$   
**Require:**  $w$ : Neural network weights at initialization  
**Require:**  $\beta_1, \beta_2, \varepsilon$ : Adam hyperparameters  
**Require:**  $(\alpha_k)_1^K$ : Learning rate schedule

$\Gamma \leftarrow \emptyset$   
**for**  $p$  in  $\{0, \dots, P-1\}$  **do**  
 $\tau_p \leftarrow 0$   
 $m_p \leftarrow 0$   
 $v_p \leftarrow 0$   
 $\hat{g}_p \leftarrow 0$   
**end for**  
**for**  $k$  in  $\{1, \dots, K\}$  **do**  
 $n \leftarrow$  draw an integer uniformly at random in  $\{0, \dots, N-1\}$   
 $p \leftarrow$  draw an integer at random in  $\{0, \dots, P-1\}$  according to  $\pi$  {Stochastic approx.}  
 $\mathcal{L}(w) \leftarrow P \|\phi_p(x_n) - (\phi_p \circ F_w)(x_n)\|_2^2$   
 $g \leftarrow \nabla \mathcal{L}(w)$   
 $m_p \leftarrow \beta_1^{(k-\tau_p)/P} m_p + (1 - \beta_1^{(k-\tau_p)/P}) g$   
 $v_p \leftarrow \beta_2^{(k-\tau_p)/P} v_p + (1 - \beta_2^{(k-\tau_p)/P}) (g \odot g)$   
 $\hat{m} \leftarrow m_p / (1 - \beta_1^{k/P})$  { $\mathcal{P}$ -Adam}  
 $\hat{v} \leftarrow v_p / (1 - \beta_2^{k/P})$   
 $g_{\text{current}} \leftarrow \hat{m} / \sqrt{\varepsilon + \hat{v}}$   
 $\tau_p \leftarrow k$   
 $g_{\text{avg}} \leftarrow \frac{1}{\max(1, \text{card} \Gamma)} \sum_{\gamma \in \Gamma} \hat{g}_\gamma$   
 $g_{\text{SAGA}} \leftarrow g_{\text{current}} - \hat{g}_p + g_{\text{avg}}$  { $\mathcal{P}$ -SAGA}  
 $w \leftarrow w - \alpha_k g_{\text{SAGA}}$   
 $\hat{g}_p \leftarrow g_{\text{current}}$   
 $\Gamma \leftarrow \Gamma \cup \{p\}$   
**end for**  
**return**  $w$

---

A known drawback of DDSP is that the optimization landscape of spectral loss in parameter space (i.e., of  $\mathcal{L}_x^\Phi \circ D$ ) may not coincide with that of supervised parameter loss (i.e., Euclidean distance to  $\theta$ , also known as P-loss) (Hayes et al., 2024). Against this drawback, we propose a method named  $\theta$ -importance sampling ( $\theta$ -IS), which constructs a categorical distribution  $\pi$  over the path space  $\mathcal{P}$ . The key idea behind  $\theta$ -IS is to introduce bias in the stochastic approximation of spectral loss so as to bring it closer to P-loss. For lack of supervision, we are unable to construct the optimal distribution  $\pi$  but provide a heuristic of this form:

$$\pi_p = \frac{1}{U} \sum_{u=0}^{U-1} \frac{C_{u,p}}{\sum_{p=0}^{P-1} C_{u,p'}}, \quad (10)$$

where, intuitively,  $C_{u,p}$  represents the importance of parameter dimension  $\theta_u$  upon path  $p$ . We rescale this importance relative to all paths and average uniformly across parameters  $u$ , yielding an importance-weighted categorical distribution  $\pi$  over paths. We then use  $\pi$  instead of a uniform distribution for sampling paths in the SCRAPL algorithm (see Algorithm 1).

Let  $E_{x,u}(w)$  denote the  $u^{\text{th}}$  coordinate of  $E_x(w)$ . Given  $w$ , we measure the sensitivity of each ST path  $p$  to the parameter control  $u$  around the input  $x$  in terms of the following partial derivative:

$$s_{x,u,p} : w \mapsto \frac{\partial (\mathcal{L}_x^{\phi_p} \circ D)}{\partial \theta_u} (E_{x,u}(w)) \quad (11)$$

To convert the sensitivity function  $s_{x,u,p}$  into relative importance  $C_{u,p}$ , we analyze the curvature of the loss landscape  $\mathcal{L}_x^{\theta_p}$ . We multiply  $s_{x,u,p}$  by the gradient of  $E_{x,u}$ , yielding a linear transformation of neural network parameters. The coordinate-wise gradient of this linear transformation yields a positive definite matrix: we compute its largest eigenvalue. We repeat this process over a representative dataset  $\mathcal{X}$  of  $N_{\text{IS}}$  unlabeled signals from the training dataset. Formally:

$$C_{u,p} = \mathbb{E}_{x \sim \mathcal{X}} [\lambda_{\max}(\nabla_w(s_{x,u,p}(w)) \nabla E_{x,u}(w))], \quad (12)$$

where  $\lambda_{\max}(\mathbf{M})$  is the magnitude of the largest eigenvalue of a square matrix  $\mathbf{M}$ ;  $\nabla_w$  is the gradient with respect to  $w$ . In practice, we compute  $\lambda_{\max}(\mathbf{M})$  using a stochastic power iteration with deflation and the Hessian vector product, which has the same asymptotic runtime complexity as a backpropagation step.<sup>1</sup> Crucially, the computation required for  $\theta$ -IS can be trivially parallelized across  $p$  and  $u$ , and only needs to be computed once before training.

Our definition of  $\theta$ -IS is inspired by Schmidt et al. (2017), who propose a variant of the SAG algorithm in which mini-batches are sampled non-uniformly; more precisely, in proportion to the Lipschitz constant of the gradients, which we approximate using Equation 12 for each  $p$  and  $u$ . This heuristic relies on the argument that gradients which change quickly should be regarded as more important than gradients which change slowly.

## 4 EXPERIMENTS

We apply SCRAPL to a differentiable implementation of the joint time–frequency scattering transform (Muradeli et al., 2022). We conduct three unsupervised sound matching experiments under the DDSP autoencoder paradigm described in Section 3.4. The encoder,  $E_x$ , is a convolutional neural network which operates on a constant- $Q$  transform (Cheuk et al., 2020). We use relatively lightweight neural networks for our experiments, a choice made possible by the strong inductive bias DDSP provides and informed by prior work (Han et al., 2025) indicating that larger networks do not necessarily improve sound matching capabilities. We choose all hyperparameters in experiments heuristically.

To highlight the new kinds of perceptual quality assessment tasks SCRAPL enables, all three experiments investigate nondeterministic decoders that introduce random time shifts into the resulting reconstructed audio. While our experiments are for a discriminative and generative audio processing task, we emphasize that SCRAPL is a general algorithm and can be equally applied to deep inverse problems that leverage other scattering transforms.

### 4.1 JOINT TIME–FREQUENCY SCATTERING TRANSFORM (JTFS)

The joint time–frequency scattering transform (JTFS) is a nonlinear convolutional operator which extracts spectrotemporal modulations over the constant- $Q$  spectrogram (Andén et al., 2019). The multivariable filter  $\Psi_p$  comprises two stages: temporal scattering, i.e., 1-D band-pass filtering with Morlet wavelets over the time axis; and frequential scattering, i.e., idem over the log-frequency axis. The center frequencies of band-pass filters for temporal scattering, called *rates*, are measured in Hertz. The center frequencies for frequential scattering, called *scales*, are measured in cycles per octave. Thus, in the case of JTFS, the path index  $p$  is a rate–scale multiindex.

The JTFS has been shown to correlate with human perception (Lostanlen et al., 2021; Tian et al., 2025) and can provide an informative gradient for audio comparisons that are misaligned (Vahidi et al., 2023) or benefit from multi-resolution analysis like percussive sounds (Han et al., 2024), which is why we select it as the underlying ST of the SCRAPL algorithm in our unsupervised sound matching experiments. Additionally, due to its computational complexity, until now it has been used almost exclusively as a precomputed feature instead of a loss function for neural network training.

### 4.2 GRANULAR SYNTH SOUND MATCHING

Granular synthesis is an example of a new class of synths that can be effectively sound matched with SCRAPL and the JTFS, due to its inherently stochastic audio generation process with individual grains being misaligned in time at the micro-level, but still being perceived as a single texture. It has been extensively used in the production of electronic music since the late 1950s<sup>2</sup> and played a

<sup>1</sup><https://github.com/noahgolmant/pytorch-hessian-eigenthings>

<sup>2</sup><https://www.iannis-xenakis.org/en/granular-synthesis/>

fundamental role in the creation of contemporary electronic music genres. Our differentiable granular synth produces textures of chirplet grains with random temporal positions, center frequencies, and chirp rates, and has two continuous parameters: density ( $\theta_{\text{density}}$ ) which controls how many grains are produced, and slope ( $\theta_{\text{slope}}$ ) which controls their rate of frequency modulation.

We compare four MSS-based losses: linear, log + linear (Engel et al., 2020), random (Steinmetz & Reiss, 2020), and a SOTA hyperparameter-tuned revisited MSS loss (Schwär & Müller, 2023). Given their correlation with human perception (Kilgour et al., 2019; Tailleur et al., 2024), we also include the Euclidean distance of MS-CLAP (Elizalde et al., 2023) and PANNs Wavegram Logmel embeddings (Kong et al., 2020). In addition, we train with ordinary (i.e., full-tree) JTFS so as to put the speed and accuracy of SCRAPL into context. Lastly, as an estimate of best achievable performance with our encoder architecture and training configuration, we run a supervised version of sound matching, also known as “parameter loss” or P-loss for short (see Section 3.4). We summarize the implementation details in Appendix E.

#### 4.3 CHIRPLET SYNTH SOUND MATCHING

Similar to the unsupervised granular synth sound matching experiment, we evaluate our  $\theta$ -importance sampling initialization heuristic for SCRAPL on a differentiable chirplet synth (based on the implementation by Vahidi et al. (2023)) with two parameters:  $\theta_{\text{AM}}$  which controls the rate of amplitude modulation (Hz) and  $\theta_{\text{FM}}$  which controls the rate of frequency modulation (oct/s). Since the paths in the JTFS correspond to specific wavelet AM and FM center frequencies, given a chirplet synth configuration with bounded  $\theta_{\text{AM}}$  and  $\theta_{\text{FM}}$  ranges, we know approximately which paths of the JTFS should provide the most informative gradients for the synth parameters. After computing our initialization heuristic, we can analyze the resulting path probabilities and verify that the paths within the parameter ranges of the synth have been assigned a probability greater than uniform.

We evaluate four different synth configurations:

1. Slow AM ( $\theta_{\text{AM}} \in [1.0, 2.0]$  Hz), slow FM ( $\theta_{\text{FM}} \in [0.5, 1.0]$  oct/s);
2. Slow AM ( $\theta_{\text{AM}} \in [1.0, 2.0]$  Hz), moderate FM ( $\theta_{\text{FM}} \in [2.0, 4.0]$  oct/s);
3. Fast AM ( $\theta_{\text{AM}} \in [2.8, 8.4]$  Hz), moderate FM ( $\theta_{\text{FM}} \in [2.0, 4.0]$  oct/s);
4. Fast AM ( $\theta_{\text{AM}} \in [2.8, 8.4]$  Hz), fast FM ( $\theta_{\text{FM}} \in [4.0, 12.0]$  oct/s).

We compare SCRAPL training runs using uniform sampling and  $\theta$ -importance sampling calculated from a single training batch of 32 examples. We summarize the implementation details in Appendix E.

#### 4.4 ROLAND TR-808 SOUND MATCHING

As a real-world evaluation task, we sound match a DDSP implementation (Shier et al., 2024) of the Roland TR-808 Drum Machine, a historically meaningful synthesizer for the creation of Detroit techno, house, and hip-hop music.<sup>3</sup> Inharmonic transient sounds like percussion are a form of non-stationary signal that the JTFS is well suited for perceptual quality assessment (Han et al., 2024) because of its ability to extract spectrotemporal patterns at multiple scales and rates. Additionally, due to the transient nature of drum sounds, they are highly sensitive to even a few milliseconds of misalignment, thus further benefiting from the time invariance of JTFS.

We use a high fidelity, 100% analog dataset<sup>4</sup> of 681 bass drum, snare, tom, and hi-hat one-shot recordings of the TR-808 and repeat experiments 40 times on different train/validation/test splits and random seeds. Since the transient of analog drum recordings is rarely perfectly aligned, and no two analog TR-808 drum synths produce the same signal, we investigate both perfectly aligned (labeled *micro*) and unaligned (labeled *meso*) sound matching by up to  $\pm 46$  ms ( $\pm 2048$  samples at 44.1 kHz). Given its correlation with human perception, we employ the JTFS and Fréchet Audio Distance (FAD) (Kilgour et al., 2019; Défossez et al., 2023) as evaluation metrics. We also include MSS and mean frame-by-frame perceptual loudness and loudness-weighted perceptually-scaled spectral centroid and flatness for both the transient and decay portions of reconstructed signals (nine metrics in total). Additional context for these last six metrics can be found in Shier et al. (2024). We summarize the implementation details in Appendix E.

<sup>3</sup><https://www.ethanhein.com/wp/2016/beatmaking-fundamentals/>

<sup>4</sup><https://samplesfrommars.com/products/tr-808-samples>

Table 1: Evaluation results for the unsupervised granular synth sound matching task with two continuous  $\theta_{\text{synth}}$  parameters:  $\theta_{\text{density}}$  and  $\theta_{\text{slope}}$  (more details in Section 4.2). Uncertainties are 95% CI for 20 training runs using different random seeds. Due to computational limitations, the JTFS method is only evaluated once.

Method	$\theta_{\text{synth}} L_1 \%$ ↓	$\theta_{\text{density}} L_1 \%$ ↓	$\theta_{\text{slope}} L_1 \%$ ↓
JTFS	<b>42.4</b>	<b>65.8</b>	<b>19.0</b>
SCRAPL (no $\theta$ -IS)	73.8 ± 13	70.4 ± 8.8	77.2 ± 19
SCRAPL	65.7 ± 4.2	72.6 ± 6.3	58.7 ± 7.5
MSS Linear	370 ± 0.52	499 ± 0.84	241 ± 0.28
MSS Log + Linear	259 ± 1.7	277 ± 3.2	241 ± 0.42
MSS Revisited	311 ± 19	376 ± 40	246 ± 3.0
MSS Random	195 ± 4.2	149 ± 7.8	242 ± 1.0
MS-CLAP	166 ± 8.2	81.9 ± 9.0	250 ± 8.2
PANNs Wavegram-Logmel	159 ± 4.4	80.3 ± 4.2	238 ± 5.5
P-loss	20.5 ± 0.20	24.7 ± 0.31	16.3 ± 0.31

## 5 RESULTS

### 5.1 GRANULAR SYNTH SOUND MATCHING

We benchmark all loss function computational costs (see Appendix B, Table 5) and plot them in Figure 1 against their evaluation accuracy (see Table 1) on  $\theta_{\text{synth}} L_1$  relative to supervised training (i.e., P-loss). We observe that SCRAPL comes within factor two of JTFS in terms of accuracy, and within factor two of MSS in terms of runtime, striking a notable balance between them. The significant difference in runtime and convergence between JTFS and SCRAPL is further illustrated in Figure 2 where we plot validation accuracy against wall-clock time, instead of optimization steps. We also note that MSS is unable to sound match the synth at all, and the SOTA embedding losses are only able to optimize  $\theta_{\text{density}}$ , albeit not as well as SCRAPL and JTFS. Validation accuracy curves for all methods are also provided in Figure 2. We provide a variation of Figure 1 plotting the test JTFS audio distance on the y-axis in Appendix B, Figure 3.

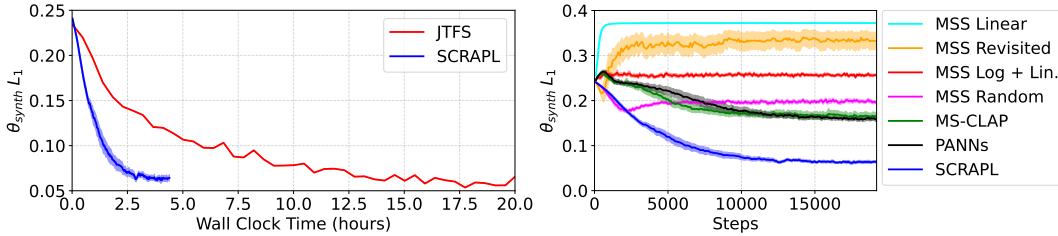


Figure 2: Left: JTFS vs. SCRAPL wall-clock training times on a single NVIDIA RTX A5000 GPU. Due to computational limitations, the JTFS method is only evaluated once. Right: Validation convergence graphs for the unsupervised granular synth sound matching task. Both: Shaded areas are 95% CI for 20 training runs using different random seeds.

Table 2 summarizes the results of an ablation of SCRAPL and its  $\mathcal{P}$ -Adam,  $\mathcal{P}$ -SAGA, and  $\theta$ -IS optimization techniques for the granular synth sound matching task, with Appendix B, Table 6 providing additional information about statistical significance. There is a clear monotonic improvement in accuracy and convergence time for each technique, as well as a reduction in variance provided by  $\mathcal{P}$ -SAGA, and  $\theta$ -IS. We emphasize that SCRAPL without any extra optimization techniques still outperforms all other non-JTFS methods in terms of accuracy, demonstrating its ability to optimize a new class of problems and making just stochastic sampling of scattering transforms a viable approach if the additional memory and computational requirements of  $\mathcal{P}$ -Adam,  $\mathcal{P}$ -SAGA, and  $\theta$ -IS are undesirable. Finally, from Table 1, we see that  $\theta$ -IS results in a better overall accuracy of  $\theta_{\text{synth}}$  than uniform sampling (despite  $\theta_{\text{density}}$  now being slightly worse), which is consistent with our hypothesis from Section 3.4 that  $\theta$ -IS results in more balanced convergence of all synth parameters. Validation accuracy curves for all ablations are provided in Appendix B, Figure 4.

Table 2: Ablation table for SCRAPL with test results and validation  $\theta_{\text{synth}} L_1$  total variation and convergence steps for the unsupervised granular synth sound matching task. Convergence is defined as  $\theta_{\text{synth}} L_1 < 100 \%$ . Statistical significance results for each additional optimization technique are presented in Appendix B, Table 6. Uncertainties are 95% CI for 20 training runs using different random seeds. Due to computational limitations, the JTFS method is only evaluated once.

Method	$\mathcal{P}$ -Adam	$\mathcal{P}$ -SAGA	$\theta$ -IS	Test $\theta_{\text{synth}} L_1 \%$ ↓	Validation		
					Total Var. ↓	Conv. Steps ↓	
SCRAPL				99.7 ± 8.2	5.30 ± 0.25	10 906 ± 1170	
	✓			87.4 ± 15	6.98 ± 0.25	8006 ± 697	
	✓	✓		73.8 ± 13	3.46 ± 0.15	7296 ± 683	
	✓	✓	✓	<b>65.7 ± 4.2</b>	<b>3.27 ± 0.12</b>	<b>6014 ± 642</b>	
JTFS				42.4	5.66	1442	
P-loss				20.5 ± 0.20	1.83 ± 0.025	672 ± 23	

In summary, this experiment demonstrates that the variance of the gradient elicited by the stochastic approximation of a ST with the SCRAPL algorithm is manageably small in the context of deep neural network (DNN) training, resulting in a favorable tradeoff between computational speed and convergence rate when compared to full-tree scattering (i.e., the JTFS). Training with SCRAPL is nearly equivalent to training with the gradient of full-tree scattering in terms of JTFS loss on unseen test data: see Appendix B, Figure 3. In terms of synthesizer parameter error, training with SCRAPL is not as accurate as training with full-tree ST, but outperforms prior work: see Table 1. This SOTA result paves the way towards a new kind of DNN training for deep inverse problems, in which the forward operator (synth) produces nondeterministic time–frequency patterns.

Understanding the convergence properties of algorithms like SCRAPL for convex and non-convex tasks is an active area of research (Reddi et al., 2016; 2018; Défossez et al., 2022; Kim & Oh, 2025). Our proof of Proposition 3.1 in Appendix A that SCRAPL without any additional optimization techniques is an unbiased estimator of full-tree ST is an important first step in this direction. We believe that further convergence analysis of SCRAPL remains a promising avenue for future work.

## 5.2 CHIRPLET SYNTH SOUND MATCHING

Table 3 summarizes the chirplet synth evaluation results, with Appendix C, Figure 5 showing validation accuracy curves for uniform and  $\theta$ -importance sampling on the four synth configurations.  $\theta$ -IS improves the prediction of  $\theta_{\text{AM}}$  by 25–55% and of  $\theta_{\text{FM}}$  by 14–80%, while reducing time to convergence by 23–50%: see Appendix C, Table 7. Of course, these improvements are for synth configurations that have been designed to showcase the benefit of nonuniform sampling of paths; however, this overall trend remains true, albeit not as pronounced, for the granular synth (Table 2) and real-world sound matching task (Table 4). Finally, we plot the path  $\theta$ -IS probabilities in Appendix C, Figure 6 and observe that indeed, a unique distribution is learned for each synth, and the greater than uniform probabilities appear to roughly correspond to each configuration’s limited AM/FM range.

## 5.3 ROLAND TR-808 SOUND MATCHING

Table 4 and Appendix D, Tables 8, and 9 summarize the unsupervised Roland TR-808 synth sound matching audio distance, transient, and decay perceptual similarity results. Overall, we observe that JTFS dominates almost all metrics in both micro and meso environments, showcasing its suitability for transient percussive sounds and temporal invariance. After JTFS, MSS performs best when samples are perfectly aligned (micro), but performs worse in the unaligned (meso) setting and is unable to match the transient, which is the most salient part of a drum hit. In contrast, SCRAPL shows consistent sound matching performance in both micro and meso environments and is able to preserve the transient even when audio is misaligned. However, SCRAPL fails to recover the less audible decay portion of the signal. We hypothesize this is due to informative, low-frequency paths for the decay being sparse and underrepresented in the categorical distribution over paths, even after accounting for  $\theta$ -IS. We provide listening samples at the accompanying website<sup>5</sup> and encourage readers to evaluate the results directly.

<sup>5</sup>Anonymous companion website: <https://icewithfrosting.github.io/scrapl/>

Table 3: Evaluation results for SCRAPL with and without the  $\theta$ -importance sampling initialization heuristic on unsupervised sound matching of four different AM / FM chirplet synths with two continuous  $\theta_{\text{synth}}$  parameters:  $\theta_{\text{AM}}$  and  $\theta_{\text{FM}}$  (more details in Section 4.3). Uncertainties are 95% CI for 20 training runs using different random seeds.

Sampling Method ( $\pi$ )	Synth Configuration		$\theta_{\text{AM}} L_1 \text{ } \%$ ↓		$\theta_{\text{FM}} L_1 \text{ } \%$ ↓	
	$\theta_{\text{AM}}$ (Hz)	$\theta_{\text{FM}}$ (oct/s)				
Uniform $\theta$ -IS	1.0 – 2.0	0.5 – 1.0	124 ± 10 <b>77.7</b> ± 6.7		155 ± 18 <b>78.4</b> ± 11	
Uniform $\theta$ -IS	1.0 – 2.0	2.0 – 4.0	111 ± 20 <b>55.5</b> ± 4.1		68.6 ± 11 <b>44.4</b> ± 2.8	
Uniform $\theta$ -IS	2.8 – 8.4	2.0 – 4.0	122 ± 22 <b>54.9</b> ± 3.5		238 ± 21 <b>48.5</b> ± 4.7	
Uniform $\theta$ -IS	2.8 – 8.4	4.0 – 12.0	108 ± 12 <b>81.5</b> ± 12		95.6 ± 20 <b>82.1</b> ± 11	

Table 4: Audio distance evaluation results for the unsupervised Roland TR-808 DDSP synth sound matching task with 14 continuous  $\theta_{\text{synth}}$  parameters (more details in Section 4.4). Uncertainties are 95% CI for 40 training runs using different random seeds and dataset splits. Due to computational limitations, the JTFS method is only trained and evaluated for 4 random seeds.

Method	MSS Log. + Linear ↓		JTFS ↓		FAD (EnCodec) ↓	
	Micro	Meso	Micro	Meso	Micro	Meso
JTFS	617 ± 46	622 ± 45	<b>490</b> ± 28	<b>523</b> ± 17	<b>0.781</b> ± 0.069	<b>1.04</b> ± 0.15
SCRAPL (no $\theta$ -IS)	862 ± 36	944 ± 48	1140 ± 48	1250 ± 51	2.75 ± 0.39	2.85 ± 0.38
SCRAPL	857 ± 42	879 ± 42	1050 ± 50	1110 ± 52	2.43 ± 0.22	2.42 ± 0.22
MSS Lin.	611 ± 15	724 ± 37	779 ± 31	1470 ± 83	1.22 ± 0.082	3.33 ± 0.46
MSS L+L	<b>596</b> ± 19	<b>615</b> ± 18	1260 ± 58	1390 ± 49	2.14 ± 0.39	3.01 ± 0.40
MSS Rev.	637 ± 16	797 ± 20	870 ± 23	1250 ± 27	2.02 ± 0.37	2.21 ± 0.34
MSS Rand.	682 ± 25	700 ± 26	1410 ± 87	1500 ± 59	7.03 ± 2.2	6.65 ± 1.7

## 6 CONCLUSION

Differentiable similarity measures have the potential to enhance the perceptual quality of generative models and deep inverse problem solvers. In spite of their mathematical guarantees and neuro-physiological plausibility, scattering transforms (ST) have not been able to realize this potential, for lack of tractable optimization algorithms. To fill this gap, SCRAPL takes advantage of the tree-like structure of ST to save computation at each backward pass. Our numerical simulations show the value of SCRAPL for unsupervised sound matching, particularly when the synthesizer of interest is nondeterministic. Although our ST architecture of choice is joint time–frequency scattering (JTFS), we stress that SCRAPL is agnostic to the specifics of multivariable filterbank design: beyond wavelet scattering, it extends to learnable scattering-like architectures (Lattner et al., 2019; Cotter & Kingsbury, 2019; Gauthier et al., 2022). We consider investigating SCRAPL’s generalizability to other ST architectures, different audio tasks such as speech enhancement and automatic mixing, and additional modalities like adversarial image generation and texture synthesis as promising directions for future work. As a longer-term perspective, the success of our architecture-informed importance sampling heuristic highlights the opportunity to meta-learn the relative importance of each ST path for the task at hand over the course of neural network training (Yamaguchi et al., 2023).

## REPRODUCIBILITY STATEMENT

Appendix E contains all hyperparameters and training details for each of the three experiments in this paper. We also provide listening samples, anonymized source code, configuration files, and instructions to reproduce our experiments at the anonymous companion website: <https://icewithfrosting.github.io/scrap/>

## REFERENCES

- Joakim Andén and Stéphane Mallat. Multiscale scattering for audio classification. In *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*. Miami, Florida, 2011.
- Joakim Andén and Stéphane Mallat. Deep scattering spectrum. *IEEE Transactions on Signal Processing*, 62(16):4114–4128, 2014.
- Joakim Andén, Vincent Lostanlen, and Stéphane Mallat. Joint time–frequency scattering. *IEEE Transactions on Signal Processing*, 67(14):3704–3718, 2019.
- Tomás Angles and Stéphane Mallat. Generative networks as inverse problems with scattering transforms. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.
- Albert Benveniste, Michel Métivier, and Pierre Priouret. *Adaptive algorithms and stochastic approximations*, volume 22. Springer Science & Business Media, 2012.
- Joan Bruna and Stéphane Mallat. Invariant scattering convolution networks. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1872–1886, 2013.
- Kin Wai Cheuk, Hans Anderson, Kat Agres, and Dorien Herremans. nnaudio: An on-the-fly gpu audio to spectrogram conversion toolbox using 1d convolutional neural networks. *IEEE Access*, 8: 161981–162003, 2020. doi: 10.1109/ACCESS.2020.3019084.
- Fergal Cotter and Nick Kingsbury. A learnable ScatterNet: Locally invariant convolutional layers. In *Proceedings of the IEEE International Conference on Image Processing (ICIP)*. IEEE, 2019.
- Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives. *Advances in Neural Information Processing Systems (NeurIPS)*, 2014.
- Alexandre Défossez, Leon Bottou, Francis Bach, and Nicolas Usunier. A simple convergence proof of adam and adagrad. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856. URL <https://openreview.net/forum?id=ZPQhzTSA7>.
- Alexandre Défossez, Jade Copet, Gabriel Synnaeve, and Yossi Adi. High fidelity neural audio compression. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. <https://openreview.net/forum?id=ivCd8z8zR2>.
- J-M Delouis, Erwan Allys, Edouard Gauvrit, and François Boulanger. Non-gaussian modelling and statistical denoising of planck dust polarisation full-sky maps using scattering transforms. *Astronomy & Astrophysics*, 668:A122, 2022.
- Aymeric Dieuleveut, Gersende Fort, Eric Moulines, and Hoi-To Wai. Stochastic approximation beyond gradient for signal processing and machine learning. *IEEE Transactions on Signal Processing*, 71:3117–3148, 2023.
- Benjamin Elizalde, Soham Deshmukh, Mahmoud Al Ismail, and Huaming Wang. Clap learning audio concepts from natural language supervision. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5, 2023. doi: 10.1109/ICASSP49357.2023.10095889.
- Jesse Engel, Lamtharn Hantrakul, Chenjie Gu, and Adam Roberts. DDSP: Differentiable digital signal processing. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020.

- Shanel Gauthier, Benjamin Thérien, Laurent Alsene-Racicot, Muawiz Chaudhary, Irina Rish, Eugene Belilovsky, Michael Eickenberg, and Guy Wolf. Parametric scattering networks. In *Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- Han Han, Vincent Lostanlen, and Mathieu Lagrange. Learning to solve inverse problems for perceptual sound matching. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 32:2605–2615, 2024.
- Han Han, Vincent Lostanlen, and Mathieu Lagrange. Gradient Clipping Improves Neural Network Optimization for Perceptual Sound Matching. In *Proceedings of the European Signal Processing Conference (EUSIPCO)*, Palermo, Italy, September 2025. URL <https://hal.science/hal-05124224>.
- Ben Hayes, Jordie Shier, György Fazekas, Andrew McPherson, and Charalampos Saitis. A review of differentiable digital signal processing for music and speech synthesis. *Frontiers in Signal Processing*, 3:1284100, 2024.
- Vassilis N. Ioannidis, Siheng Chen, and Georgios B. Giannakis. Pruned graph scattering transforms. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020.
- Kevin Kilgour, Mauricio Zuluaga, Dominik Roblek, and Matthew Sharifi. Fréchet audio distance: A reference-free metric for evaluating music enhancement algorithms. In *Proceedings of the International Speech Communication Association Conference (Interspeech)*, pp. 2350–2354, 2019. doi: 10.21437/Interspeech.2019-2219.
- Gyu Yeol Kim and Min-hwan Oh. ADAM optimization with adaptive batch selection. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=BZrSCv2SBq>.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2014.
- Qiuqiang Kong, Yin Cao, Turab Iqbal, Yuxuan Wang, Wenwu Wang, and Mark D. Plumbley. Panns: Large-scale pretrained audio neural networks for audio pattern recognition. *IEEE/ACM Trans. Audio, Speech and Lang. Proc.*, 28:2880–2894, November 2020. ISSN 2329-9290. doi: 10.1109/TASLP.2020.3030497. URL <https://doi.org/10.1109/TASLP.2020.3030497>.
- Nina Kowalski, Didier A Depireux, and Shihab A Shamma. Analysis of dynamic spectra in ferret primary auditory cortex. i. characteristics of single-unit responses to moving ripple spectra. *Journal of Neurophysiology*, 76(5):3503–3523, 1996.
- Stefan Lattner, Monika Dörfler, and Andreas Arzt. Learning complex basis functions for invariant representations of audio. In *Proceedings of the International Society Society for Music Information Retrieval Conference (ISMIR)*, 2019.
- Vincent Lostanlen and Stéphane Mallat. Wavelet scattering on the pitch spiral. In *Proceedings of the Digital Audio Effects Conference (DAFx)*, 2016.
- Vincent Lostanlen, Joakim Andén, and Mathieu Lagrange. Fourier at the heart of computer music: From harmonic sounds to texture. *Comptes Rendus. Physique*, 20(5):461–473, 2019.
- 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*, 2021(1):3, 2021.
- Stéphane Mallat. Group invariant scattering. *Communications on Pure and Applied Mathematics*, 65(10):1331–1398, 2012.
- Stéphane Mallat. Understanding deep convolutional networks. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 374(2065):20150203, 2016.
- John Muradeli, Cyrus Vahidi, Changhong Wang, Han Han, Vincent Lostanlen, Mathieu Lagrange, and George Fazekas. Differentiable time-frequency scattering on gpu. In *Proceedings of the Digital Audio Effects Conference (DAFx)*, 2022.

- Jean-François Muzy, Emmanuel Bacry, Stéphane Mallat, and Joan Bruna. Intermittent process analysis with scattering moments. *Annals of Statistics*, 43(1):323, 2015.
- Sam V Norman-Haignere and Josh H McDermott. Neural responses to natural and model-matched stimuli reveal distinct computations in primary and nonprimary auditory cortex. *PLoS biology*, 16(12):e2005127, 2018.
- Edouard Oyallon, Stéphane Mallat, and Laurent Sifre. Generic deep networks with wavelet scattering. In *Proceedings of the International Conference on Learning Representations (ICLR) Workshop Track*, 2014.
- Edouard Oyallon, Sergey Zagoruyko, Gabriel Huang, Nikos Komodakis, Simon Lacoste-Julien, Matthew Blaschko, and Eugene Belilovsky. Scattering networks for hybrid representation learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(9):2208–2221, 2018.
- Kailash Patil, Daniel Pressnitzer, Shihab Shamma, and Mounya Elhilali. Music in our ears: The biological bases of musical timbre perception. *PLoS computational biology*, 8(11):e1002759, 2012.
- Javier Portilla and Eero P Simoncelli. A parametric texture model based on joint statistics of complex wavelet coefficients. *International Journal of Computer Vision*, 40(1):49–70, 2000.
- Sashank J. Reddi, Suvrit Sra, Barnabás Póczos, and Alex Smola. Fast incremental method for smooth nonconvex optimization. In *2016 IEEE 55th Conference on Decision and Control (CDC)*, pp. 1971–1977. IEEE Press, 2016. doi: 10.1109/CDC.2016.7798553. URL <https://doi.org/10.1109/CDC.2016.7798553>.
- Sashank J. Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=ryQu7f-RZ>.
- Mark Schmidt, Nicolas Le Roux, and Francis Bach. Minimizing finite sums with the stochastic average gradient. *Math. Program.*, 162(1–2):83–112, March 2017. ISSN 0025-5610. doi: 10.1007/s10107-016-1030-6. URL <https://doi.org/10.1007/s10107-016-1030-6>.
- Simon Schwär and Meinard Müller. Multi-scale spectral loss revisited. *IEEE Signal Processing Letters*, 30:1712–1716, 2023. doi: 10.1109/LSP.2023.3333205.
- Jordie Shier, Charalampos Saitis, Andrew Robertson, and Andrew McPherson. Real-time timbre remapping with differentiable dsp. In *Proceedings of the International Conference on New Interfaces for Musical Expression NIME*, 2024. URL <https://arxiv.org/abs/2407.04547>.
- Laurent Sifre and Stéphane Mallat. Rotation, scaling and deformation invariant scattering for texture discrimination. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1233–1240, 2013.
- 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.
- Modan Tailleur, Junwon Lee, Mathieu Lagrange, Keunwoo Choi, Laurie M. Heller, Keisuke Imoto, and Yuki Okamoto. Correlation of fréchet audio distance with human perception of environmental audio is embedding dependant. In *Proceedings of the European Signal Processing Conference (EUSIPCO)*, 2024.
- Haokun Tian, Stefan Lattner, and Charalampos Saitis. Assessing the alignment of audio representations with timbre similarity ratings. *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, 2025.
- Cyrus Vahidi, Han Han, Changhong Wang, Mathieu Lagrange, György Fazekas, and Vincent Lostanlen. Mesostructures: Beyond spectrogram loss in differentiable time-frequency analysis. *Journal of the Audio Engineering Society*, 71(9):577–585, 2023.

Shin'ya Yamaguchi, Daiki Chijiwa, Sekitoshi Kanai, Atsutoshi Kumagai, and Hisashi Kashima. Regularizing neural networks with meta-learning generative models. *Advances in Neural Information Processing Systems (NeurIPS)*, 36:27315–27331, 2023.

Ryuichi Yamamoto, Eunwoo Song, and Jae-Min Kim. Parallel wavegan: A fast waveform generation model based on generative adversarial networks with multi-resolution spectrogram. *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020.

## A PROOF OF PROPOSITION 3.1

Let us write  $\tilde{x} = F_x(w)$ . By linearity of the gradient, we may decompose  $\nabla \mathcal{L}_x^\Phi(\tilde{x})$  over paths:

$$\nabla \mathcal{L}_x^\Phi(\tilde{x}) = \frac{1}{P} \sum_{p=0}^{P-1} \nabla \mathcal{L}_x^{\phi_p}(\tilde{x}). \quad (13)$$

Let us denote the Jacobian of  $F_x$  at  $w$  by  $\mathbf{J}_{F_x}(w)$ . For each  $p \in \mathcal{P}$ , we apply the chain rule:

$$\nabla(\mathcal{L}_x^{\phi_p} \circ F_x)(w) = \mathbf{J}_{F_x}(w)^\top \nabla \mathcal{L}_x^{\phi_p}(\tilde{x}). \quad (14)$$

Plugging the identity above into Equation 13 yields:

$$\begin{aligned} \nabla(\mathcal{L}_x^\Phi \circ F_x)(w) &= \frac{1}{P} \sum_{p=0}^{P-1} \left( \mathbf{J}_{F_x}(w)^\top \nabla \mathcal{L}_x^{\phi_p}(\tilde{x}) \right) \\ &= \mathbf{J}_{F_x}(w)^\top \left( \frac{1}{P} \sum_{p=0}^{P-1} \nabla \mathcal{L}_x^{\phi_p}(\tilde{x}) \right), \end{aligned} \quad (15)$$

where the latter equation holds by associativity of matrix multiplication.

We now compute the expected value of  $\nabla \mathcal{L}_x^{\phi_z}(\tilde{x})$  for  $z \sim \mathcal{U}_P$ , i.e., a uniform distribution over  $\mathcal{P}$ :

$$\mathbb{E}_{z \sim \mathcal{U}_P} [\nabla \mathcal{L}_x^{\phi_z}(\tilde{x})] = \frac{1}{P} \sum_{p=0}^{P-1} \nabla \mathcal{L}_x^{\phi_p}(\tilde{x}). \quad (16)$$

We recognize the column vector on the right-hand side of Equation 15. Thus:

$$\begin{aligned} \nabla(\mathcal{L}_x^\Phi \circ F_x)(w) &= \mathbf{J}_{F_x}(w)^\top \mathbb{E}_{z \sim \mathcal{U}_P} [\nabla \mathcal{L}_x^{\phi_z}(\tilde{x})] \\ &= \mathbb{E}_{z \sim \mathcal{U}_P} \left[ \mathbf{J}_{F_x}(w)^\top \nabla \mathcal{L}_x^{\phi_z}(\tilde{x}) \right], \end{aligned} \quad (17)$$

where the latter equation holds by linearity of the expected value. Finally, we use the reverse form of the chain rule in Equation 14 to identify the expected SCRAPL gradient:

$$\nabla(\mathcal{L}_x^\Phi \circ F_x)(w) = \mathbb{E}_{z \sim \mathcal{U}_P} [\nabla(\mathcal{L}_x^{\phi_z} \circ F_x)(w)] \quad (18)$$

concluding the proof.

## B ADDITIONAL GRANULAR SYNTH EVALUATION RESULTS

Table 5: Loss function benchmark results for one optimization step (forward + backward, 32768 samples of audio, batch size 4, 1 thread, single precision, 1 NVIDIA RTX A5000 GPU, CUDA 12.4, PyTorch 2.8.0). SCRAPL paths are benchmarked individually and then aggregated across all paths using the median for time and interquartile range (IQR), and maximum for memory usage.

Method	Median Time (ms) ↓	IQR (ms) ↓	Max. Memory (MB) ↓
JTFS	1730	23.9	12 967
SCRAPL	89.8	3.62	2503
MSS Linear	26.3	1.12	694
MSS Log + Linear	19.1	0.696	702
MSS Revisited	<b>17.0</b>	<b>0.210</b>	<b>663</b>
MSS Random	24.7	5.81	706
MS-CLAP	75.6	1.69	2032
PANNs Wavegram-Logmel	29.3	5.92	1360
P-loss ( $\dim(\theta_{\text{synth}}) = 2$ )	0.516	0.108	625

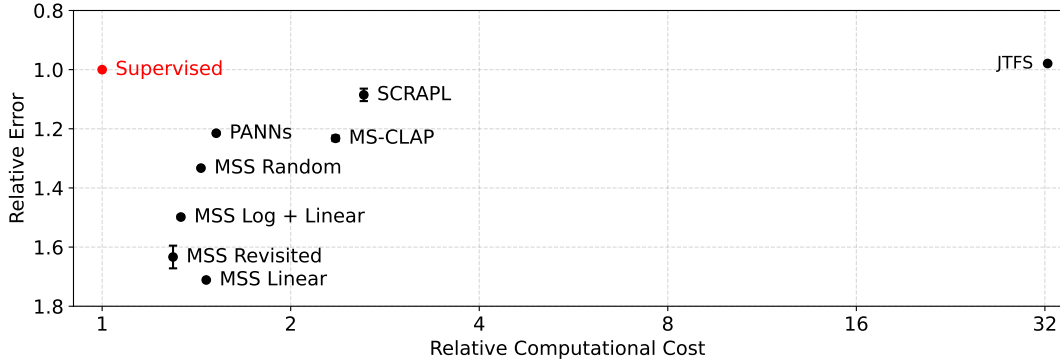


Figure 3: Mean average JTFS perceptual audio distance (y-axis) versus computational cost (x-axis) of unsupervised sound matching models for the granular synthesis task. Both axes are rescaled by the performance of a supervised model with the same number of parameters. Whiskers denote 95% CI, estimated over 20 random seeds. Due to computational limitations, JTFS-based sound matching is evaluated only once.

Table 6: Statistical significance and relative improvement of each additional SCRAPL optimization technique for the unsupervised granular synth sound matching task. Convergence is defined as  $\theta_{\text{synth}} L_1 < 100\%$ . See Table 2 for absolute results and comparisons to JTFS and P-loss. Uncertainties are 95% CI for 20 training runs using different random seeds. Due to computational limitations, the JTFS method is only evaluated once.

Opt. Technique	Test	Validation	
	$\Delta \theta_{\text{synth}} L_1 \%$ ↓	$\Delta$ Total Var. ↓	$\Delta$ Conv. Steps ↓
+ $\mathcal{P}$ -Adam	$-12.3 \pm 17$ ( $p = 0.15$ )	$1.68 \pm 0.35$ ( $p < 0.01$ )	<b>-2900</b> $\pm 1800$ ( $p < 0.01$ )
+ $\mathcal{P}$ -SAGA	<b>-13.6</b> $\pm 5.6$ ( $p < 0.01$ )	<b>-3.52</b> $\pm 0.18$ ( $p < 0.01$ )	-710 $\pm 450$ ( $p < 0.01$ )
+ $\theta$ -IS	$-8.13 \pm 15$ ( $p = 0.26$ )	$-0.188 \pm 0.092$ ( $p < 0.01$ )	-1280 $\pm 410$ ( $p < 0.01$ )

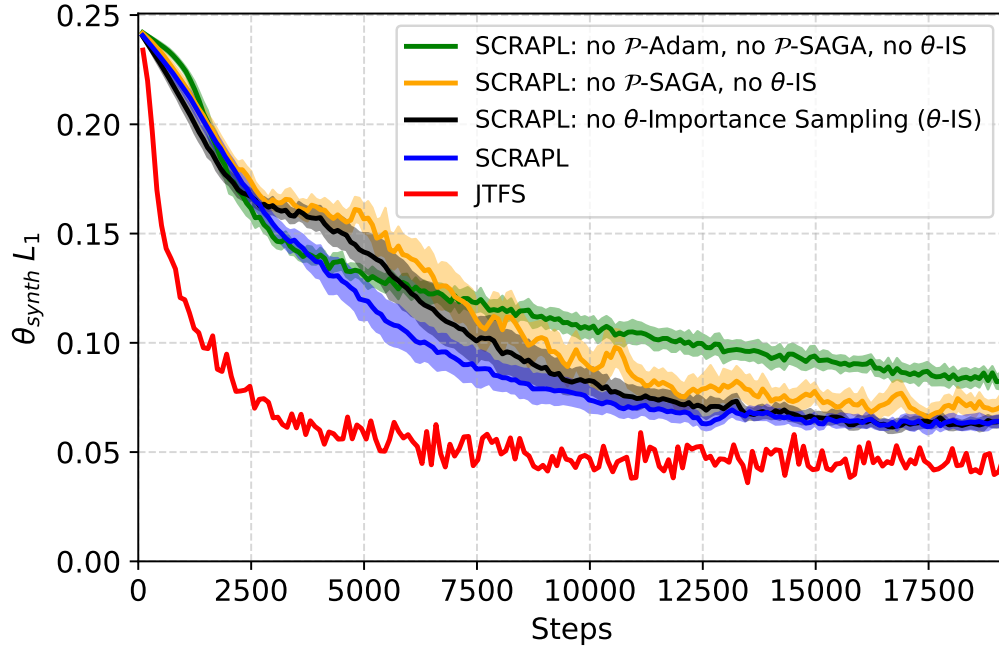


Figure 4: Validation convergence graphs of SCRAPL ablations and the JTFS for the unsupervised granular synth sound matching task. Shaded areas are 95% CI for 20 training runs using different random seeds. Due to computational limitations, the JTFS method is only evaluated once.

## C ADDITIONAL CHIRPLET SYNTH EVALUATION RESULTS

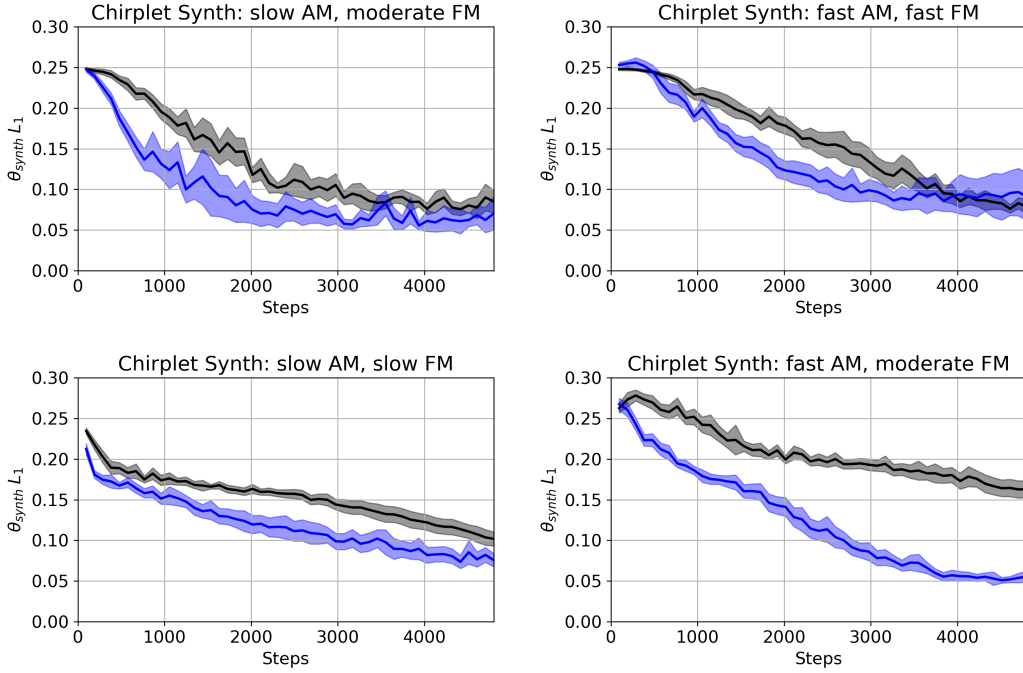


Figure 5: SCRAPL  $\theta_{\text{synth}} L_1$  validation values during training for four different AM / FM chirplet synths with two continuous  $\theta_{\text{synth}}$  parameters:  $\theta_{\text{AM}}$  and  $\theta_{\text{FM}}$  (more details in Section 4.3). Blue is using the  $\theta$ -importance sampling initialization heuristic, and black is using uniform sampling. Shaded areas are 95% CI for 20 training runs using different random seeds.

Table 7: Convergence rate (CR) and steps for SCRAPL with and without the  $\theta$ -importance sampling initialization heuristic on unsupervised sound matching of four different AM / FM chirplet synths with two continuous  $\theta_{\text{synth}}$  parameters:  $\theta_{\text{AM}}$  and  $\theta_{\text{FM}}$  (more details in Section 4.3). Convergence is defined as  $L_1 < 100\%$  for  $\theta_{\text{AM}}$  or  $\theta_{\text{FM}}$ . Uncertainties are 95% CI for 20 training runs using different random seeds.

Sampling Method ( $\pi$ )	Synth Configuration		$\theta_{\text{AM}}$		$\theta_{\text{FM}}$	
	$\theta_{\text{AM}}$ (Hz)	$\theta_{\text{FM}}$ (oct/s)	CR $\uparrow$	Conv. Steps $\downarrow$	CR $\uparrow$	Conv. Steps $\downarrow$
Uniform $\theta$ -IS	1.0 – 2.0	0.5 – 1.0	60%	3944 $\pm$ 342	45%	4064 $\pm$ 372
			<b>100%</b>	<b>2002 <math>\pm</math> 324</b>	<b>100%</b>	<b>3134 <math>\pm</math> 492</b>
Uniform $\theta$ -IS	1.0 – 2.0	2.0 – 4.0	<b>100%</b>	2203 $\pm$ 135	<b>100%</b>	1536 $\pm$ 194
			<b>100%</b>	<b>1099 <math>\pm</math> 173</b>	<b>100%</b>	<b>768 <math>\pm</math> 118</b>
Uniform $\theta$ -IS	2.8 – 8.4	2.0 – 4.0	95%	3254 $\pm$ 250	0%	N/A
			<b>100%</b>	<b>1925 <math>\pm</math> 165</b>	<b>100%</b>	<b>2966 <math>\pm</math> 210</b>
Uniform $\theta$ -IS	2.8 – 8.4	4.0 – 12.0	<b>100%</b>	3096 $\pm$ 334	<b>95%</b>	3208 $\pm$ 235
			95%	<b>2253 <math>\pm</math> 218</b>	<b>95%</b>	<b>2178 <math>\pm</math> 173</b>

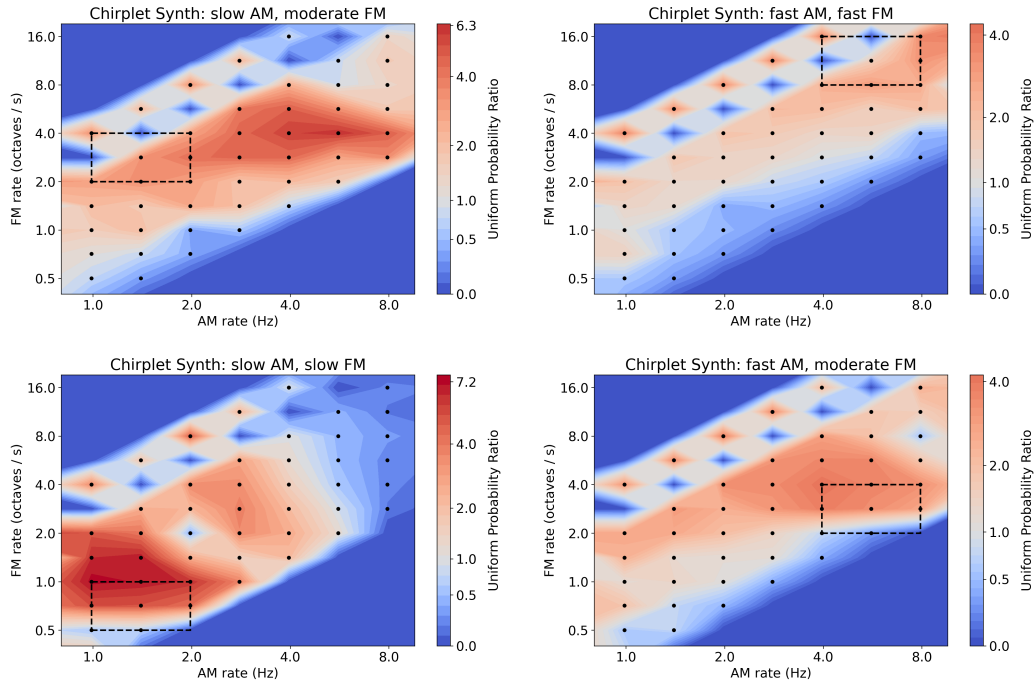


Figure 6: SCRAPL path  $\theta$ -importance sampling probabilities for four different AM / FM chirplet synths calculated from 1 batch of 32 log-uniformly randomly sampled  $\theta_{\text{synth}}$  values. Black dots are individual path (wavelet) AM / FM center frequency locations and each dashed rectangle is the  $\theta_{\text{synth}}$  range for each synth configuration. A uniform probability ratio of 1.0 means a path is sampled with probability  $1/P$ .

## D ADDITIONAL ROLAND TR-808 EVALUATION RESULTS

Table 8: Drum transient evaluation results for the unsupervised Roland TR-808 DDSP synth sound matching task with 14 continuous  $\theta_{\text{synth}}$  parameters (more details in Section 4.4). Uncertainties are 95% CI for 40 training runs using different random seeds and dataset splits. Due to computational limitations, the JTFS method is only trained and evaluated for 4 random seeds.

Method	Loudness $L_1 \downarrow$		Spectral Centroid $L_1 \downarrow$		Spectral Flatness $L_1 \downarrow$	
	Micro	Meso	Micro	Meso	Micro	Meso
JTFS	<b>137</b> $\pm$ 10	<b>158</b> $\pm$ 20	<b>859</b> $\pm$ 68	<b>819</b> $\pm$ 96	1200 $\pm$ 87	<b>1090</b> $\pm$ 110
SCRAPL						
(no $\theta$ -IS)	389 $\pm$ 41	460 $\pm$ 60	1020 $\pm$ 100	992 $\pm$ 110	1800 $\pm$ 170	1990 $\pm$ 180
SCRAPL	374 $\pm$ 39	377 $\pm$ 32	1000 $\pm$ 120	1080 $\pm$ 120	1750 $\pm$ 270	1820 $\pm$ 220
MSS Lin.	381 $\pm$ 12	2510 $\pm$ 480	902 $\pm$ 27	2350 $\pm$ 310	962 $\pm$ 43	3620 $\pm$ 680
MSS L+L	492 $\pm$ 44	1080 $\pm$ 91	928 $\pm$ 65	1380 $\pm$ 76	<b>916</b> $\pm$ 50	1320 $\pm$ 150
MSS Rev.	330 $\pm$ 21	808 $\pm$ 40	1070 $\pm$ 49	1540 $\pm$ 53	1390 $\pm$ 62	2640 $\pm$ 96
MSS Rand.	584 $\pm$ 75	1030 $\pm$ 89	1200 $\pm$ 100	1350 $\pm$ 99	1690 $\pm$ 120	1950 $\pm$ 140

Table 9: Drum decay evaluation results for the unsupervised Roland TR-808 DDSP synth sound matching task with 14 continuous  $\theta_{\text{synth}}$  parameters (more details in Section 4.4). Uncertainties are 95% CI for 40 training runs using different random seeds and dataset splits. Due to computational limitations, the JTFS method is only trained and evaluated for 4 random seeds.

Method	Loudness $L_1 \downarrow$		Spectral Centroid $L_1 \downarrow$		Spectral Flatness $L_1 \downarrow$	
	Micro	Meso	Micro	Meso	Micro	Meso
JTFS	315 $\pm$ 22	<b>355</b> $\pm$ 110	614 $\pm$ 51	617 $\pm$ 71	527 $\pm$ 31	718 $\pm$ 190
SCRAPL						
(no $\theta$ -IS)	1810 $\pm$ 190	2210 $\pm$ 210	1530 $\pm$ 170	1860 $\pm$ 170	2620 $\pm$ 310	3300 $\pm$ 370
SCRAPL	1810 $\pm$ 160	1740 $\pm$ 170	1490 $\pm$ 120	1470 $\pm$ 140	2540 $\pm$ 290	2480 $\pm$ 290
MSS Lin.	357 $\pm$ 12	1120 $\pm$ 260	654 $\pm$ 18	1110 $\pm$ 160	<b>472</b> $\pm$ 17	1500 $\pm$ 350
MSS L+L	389 $\pm$ 42	466 $\pm$ 45	<b>563</b> $\pm$ 22	<b>597</b> $\pm$ 24	565 $\pm$ 29	<b>644</b> $\pm$ 51
MSS Rev.	<b>279</b> $\pm$ 12	494 $\pm$ 22	589 $\pm$ 21	801 $\pm$ 29	552 $\pm$ 22	846 $\pm$ 29
MSS Rand.	453 $\pm$ 21	485 $\pm$ 24	660 $\pm$ 27	640 $\pm$ 35	594 $\pm$ 30	658 $\pm$ 33

## E EXPERIMENT TRAINING DETAILS AND HYPERPARAMETERS

Table 10: Unsupervised granular synth sound matching task hyperparameters.

Category	Hyperparameter Name	Value
Data	$N$ (# of examples)	5120
	train / val / test split	60% / 20% / 20%
Encoder	# of parameters	604 K
	CQT # of octaves	5
	CQT bins / octave	12
	CQT hop length	256
	CQT postprocessing	log1p
	CNN # of conv. blocks	5
	CNN kernel size	$5 \times (3, 3)$
	CNN stride	$5 \times (1, 1)$
	CNN pooling	$5 \times (2, 2)$
	CNN conv. block channels	128
	CNN activation function	PReLU
	CNN embedding dim.	64
	CNN dense layer dropout prob.	0.5
Decoder (Synth)	$\dim(\theta_{\text{synth}})$	2
	sampling rate	8192 Hz
	$T$ (# of samples)	32768
	max. # of grains	64
	grain # of samples	4096
	min. grain pitch	256 Hz
	max. grain pitch	2048 Hz
SCRAPL & JTFS	$J$	12
	$Q_1$	8
	$Q_2$	2
	$J_{\text{fr}}$	3
	$Q_{\text{fr}}$	2
	$T$	4096
	$F$	8
	$\rho$	identity function
	$P$ (# of paths)	315
$\theta$ -Importance Sampling	$N_{\text{IS}}$ (# of examples)	320
	# of deflated power iterations	20
Training	# of random seed training runs	20
	epochs	200
	batch size	32
	starting learning rate	$1 \times 10^{-5}$
	learning rate scheduler	none
	Adam $\beta_1$	0.9
	Adam $\beta_2$	0.999
	weight decay	0.01

Table 11: Unsupervised AM/FM chirplet synth sound matching task hyperparameters.

Category	Hyperparameter Name	Value
Data	$N$ (# of examples)	5120
	train / val / test split	60% / 20% / 20%
Encoder	# of parameters	604 K
	CQT # of octaves	5
	CQT bins / octave	12
	CQT hop length	256
	CQT postprocessing	log1p
	CNN # of conv. blocks	5
	CNN kernel size	$5 \times (3, 3)$
	CNN stride	$5 \times (1, 1)$
	CNN pooling	$5 \times (2, 2)$
	CNN conv. block channels	128
	CNN activation function	PReLU
	CNN embedding dim.	64
	CNN dense layer dropout prob.	0.5
Decoder (Synth)	$\dim(\theta_{\text{synth}})$	2
	sampling rate	8192 Hz
	$T$ (# of samples)	32768
	chirplet center frequency	512 Hz
	chirplet bandwidth	2 octaves
	min. time shift	-2048 samples
	max. time shift	+2048 samples
SCRAPL & JTFS	$J$	12
	$Q_1$	8
	$Q_2$	2
	$J_{fr}$	3
	$Q_{fr}$	2
	$T$	4096
	$F$	8
	$\rho$	identity function
	$P$ (# of paths)	315
$\theta$ -Importance Sampling	$N_{\text{IS}}$ (# of examples)	32
	# of deflated power iterations	20
Training	# of random seed training runs	20
	epochs	50
	batch size	32
	starting learning rate	$1 \times 10^{-4}$
	learning rate scheduler	none
	Adam $\beta_1$	0.9
	Adam $\beta_2$	0.999
	weight decay	0.01

Table 12: Unsupervised Roland TR-808 synth sound matching task hyperparameters.

Category	Hyperparameter Name	Value
Data	$N$ (# of examples)	681
	$N_{\text{bass drum}}$	215
	$N_{\text{snare}}$	240
	$N_{\text{tom}}$	189
	$N_{\text{hi-hat}}$	37
	$N_{\text{train}}$	425
	$N_{\text{val}}$	128
	$N_{\text{test}}$	128
Encoder	# of parameters	724 K
	CQT # of octaves	9
	CQT bins / octave	12
	CQT hop length	256
	CQT postprocessing	loglp
	CNN # of conv. blocks	5
	CNN kernel size	$5 \times (3, 3)$
	CNN stride	$5 \times (1, 1)$
	CNN pooling	(2, 2), (2, 2), (2, 3), (3, 3), (3, 3)
	CNN conv. block channels	128
	CNN activation function	PReLU
	CNN embedding dim.	128
	CNN dense layer dropout prob.	0.25
Decoder (Synth)	$\dim(\theta_{\text{synth}})$	14
	sampling rate	44100 Hz
	$T$ (# of samples)	44100
	min. time shift	-2048 samples
	max. time shift	+2048 samples
SCRAPL & JTFS	$J$	12
	$Q_1$	8
	$Q_2$	2
	$J_{fr}$	5
	$Q_{fr}$	2
	$T$	2048
	$F$	1
	$\rho$	loglp
	$P$ (# of paths)	483
$\theta$ -Importance Sampling	$N_{\text{IS}}$ (# of examples)	16
	# of deflated power iterations	20
Training	# of random seed training runs	40
	epochs	50
	batch size	8
	starting learning rate	$1 \times 10^{-5}$
	learning rate scheduler	linearly decreasing until $1 \times 10^{-4}$
	Adam $\beta_1$	0.9
	Adam $\beta_2$	0.999
	weight decay	0.01