

Sobolev Training for Implicit Neural Representations with Approximated Image Derivatives

Wentao Yuan^{1,2} Qingtian Zhu¹ Xiangyue Liu¹ Yikang Ding¹
Haotian Zhang^{1*} Chi Zhang¹

¹Megvii Research ²Peking University

Abstract. Recently, Implicit Neural Representations (INRs) parameterized by neural networks have emerged as a powerful and promising tool to represent different kinds of signals due to its continuous, differentiable properties, showing superiorities to classical discretized representations. However, the training of neural networks for INRs only utilizes input-output pairs, and the derivatives of the target output with respect to the input, which can be accessed in some cases, are usually ignored. In this paper, we propose a training paradigm for INRs whose target output is image pixels, to encode image derivatives in addition to image values in the neural network. Specifically, we use finite differences to approximate image derivatives. We show how the training paradigm can be leveraged to solve typical INRs problems, i.e., image regression and inverse rendering, and demonstrate this training paradigm can improve the data-efficiency and generalization capabilities of INRs. The code of our method is available at https://github.com/megvii-research/Sobolev_INRs.

Keywords: implicit neural representations, finite differences, Sobolev training

1 Introduction

A promising recent direction in computer vision and computer graphics is encoding complex signals implicitly by aid of multi-layer perceptron (MLP) named Implicit Neural Representations (INRs) [28,21,33,25,24,31,32,35]. The MLPs (more specifically, coordinate-based MLPs) take low-dimensional coordinates as inputs and are trained to output a representation of shape, density, and/or colors at each input location.

Compared to classical alternatives, i.e., discrete grid-based representation, INRs offer two main benefits. (a) Due to the fact that MLPs are continuous functions, they are significantly more memory efficient than discrete grid-based representations and theoretically, signals parameterized by MLPs can be presented in arbitrary resolution. (b) INRs are naturally differentiable, so gradient

This work is done by the first four authors as interns at Megvii Research.

* Corresponding author (zhanghaotian@megvii.com).

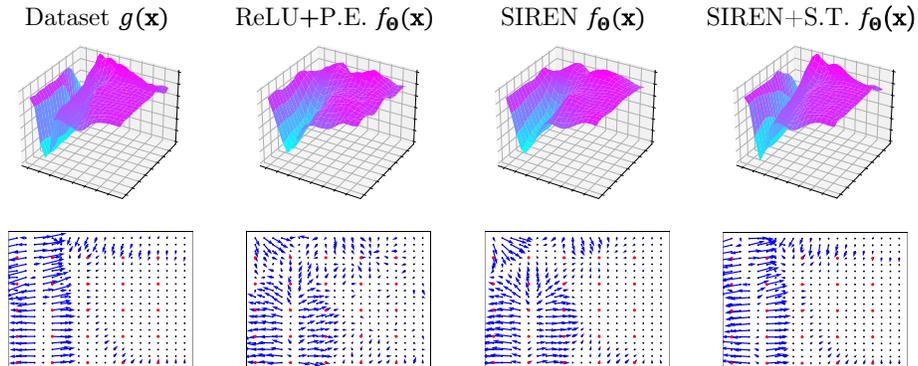


Fig. 1: Visualized approximation results of different training schemes and different activation functions. P.E. stands for positional encoding [24,36]; S.T. stands for Sobolev training. Top: dataset g (the target function to approximate) as well as approximated functions f_{θ} (functions parameterized by neural networks) represented by MLPs within an image patch of 20×20 pixels. Bottom: vector fields of the first-order derivatives, where red points are training samples (in a total proportion of 6.25%) and remaining black points are unseen samples. When activated by sine functions and trained with additional supervision on derivatives, the network demonstrates the best approximation of both function values and derivatives.

descent (GD) can be applied to optimize MLPs using modern auto-differentiation tools, e.g., PyTorch [30], TensorFlow [1], JAX [6].

However, as will be discussed below, the training paradigm adopted by most of the INRs only utilizes input-output pairs, which brings about bad generalization capabilities on unseen coordinates. In other words, if trained with low-resolution training data, MLPs will yield heavily degraded predictions when the target outputs are high-resolution. Based on this observation, an intuitive idea is that we can take advantage of the differentiability of INRs to improve the generalization capabilities. In this paper, we propose a training paradigm for INRs whose training data and the desired results are both images, with both supervision on values and derivatives enforced. Considering that the ground truth representation (the “continuous” image) is agnostic, finite differences are applied to approximate numerical first-order derivatives on images.

Further, we study some important peripheral problems relevant to the proposed training paradigm, e.g., the choice of activation functions, the number of layers in MLPs and the choice of image filters.

We find that in practice, most of the recent INRs build on ReLU-based MLPs [24,36,9] fail to represent well the derivatives of target outputs, even if the corresponding supervision is provided. For this reason, we use periodic activation functions [27,31], i.e., sine function, to improve MLPs’ convergence property of derivatives. As is shown in Fig. 1, INRs trained with traditional value-based

schemes fail to explicitly constraint the behavior at unseen coordinates, leading to poor generalization both in the image field and the vector field of the first-order derivatives. Besides, with sine-activated MLPs trained with both value loss and derivative loss, INRs are able to generalize well even if the training samples only take a very limited proportion (6.25% in Fig. 1).

We show how the training paradigm can be adapted to specific problems, namely direct image regression [31,36] and indirect inverse rendering [24,36]. Experiments results demonstrate that our training paradigm can improve the data-efficiency and generalization capabilities of INRs. The main contributions are summarized below.

- We propose a training paradigm for INRs whose training data as well as target outputs are image pixels. With the paradigm, image values and image derivatives are encoded into the INRs. Concretely, we use finite differences to approximate the derivatives of the ground truth signal. To the best of our knowledge, supervising network derivatives in INRs is minimally explored.
- We study the factor of activation functions in INRs when applying the proposed training paradigm. By experiments, We use sine functions for better convergence property of derivatives.
- Experiment results on two typical INRs problems demonstrate that with the proposed paradigm we are able to train better INRs - representations that require fewer training samples and generalise better on unseen inputs.

2 Related Work

2.1 Implicit Neural Representations

Implicit Neural Representations (INRs), which usually represent an object as a multi-layer perceptron model that maps low-dimensional coordinates to signal values, have drawn a lot of attention recently. Since this representation is continuous and can capture fine details of signals, it has been widely applied to novel view synthesis [24,20,19,29,42,41,37], shape representation [10,28,22,7,4,14,44,39], and multi-view 3D reconstruction [21,26,8,45]. As a milestone, Mildenhall et al. [24] propose a novel view synthesis method that learns an implicit representation for a specific 3D scene by using a set of multi-view calibrated images. Recent work from Park et al. [28] puts forward to learn a Signed Distance Function (SDF) to represent the shape of objects, and different SDFs can be inferred with different input latent codes. However, the (approximated) ground truth derivative information of the target signal has rarely been utilized to train INRs.

2.2 Derivative Supervision

Derivative information for neural networks has also been exploited in some previous work. Hornik [16] proves the universal approximation theorems for neural networks in Sobolev spaces - a vector space of functions equipped with a

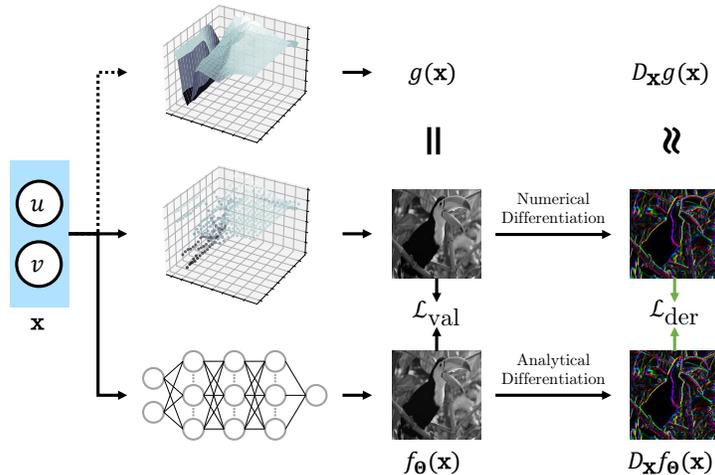


Fig. 2: An overview of the proposed training paradigm for INRs. Top: the underlying continuous signal g and its derivatives $D_{\mathbf{x}}g$. Middle: the pixels of image \mathbf{I} , which can be considered as discrete but exact sampling of g and its derivatives are obtained by numerical approximation at image space. Bottom: the MLP denoted as f_{θ} used as a general approximator whose derivatives are analytical thanks to auto-differentiation tools. In the proposed training paradigm, we enforce the supervision on both values and derivatives.

norm that is a combination of ℓ_p -norms of the function together with its derivatives up to a given order. Hornik [16] shows that neural networks with non-constant, bounded, continuous activation functions, with continuous derivatives up to order K are universal approximators in the Sobolev spaces of order K . Further, Czarnecki et al. [12] prove ReLU-based MLPs are universal approximators for function values and first-order derivatives theoretically and propose Sobolev training for neural networks. However, only the scenarios with analytical derivatives are covered and the ReLU-based MLPs show relatively poor convergence properties in practice. Gropp et al. [14] propose an *Eikonal term* which is a first-order derivative related penalty term to regularize INRs. The Eikonal term is different from our explicit derivative supervision with ground truth derivatives approximation. Sitzmann et al. [31] verify that INRs with periodic activation functions can fit derivatives robustly, while the supervision on both values and derivatives and the property of Sobolev training in INRs remain unexplored.

3 Methodology

3.1 Formulation

Considering a continuous target signal $g : \mathbb{R}^D \rightarrow \mathbb{R}^C$, which is sampled discretely at $\{\mathbf{x}_i\}_{i=1}^N$, the goal of INRs is to approximate g with a neural network (typically

an MLP) $f_{\boldsymbol{\theta}} : \mathbb{R}^D \rightarrow \mathbb{R}^C$, parameterized by weights $\boldsymbol{\theta}$. The inputs to the MLP are typically low dimensional coordinates $\mathbf{x} \in \mathbb{R}^D$, e.g., $D = 2$ for image coordinates, and the corresponding outputs are signal values, e.g., RGB colors. Training pairs of most recent INRs are denoted as $\{(\mathbf{x}_i, g(\mathbf{x}_i))\}_{i=1}^N$, where N is the number of total training samples.

By measuring the distance between predictions of the MLP $\{f_{\boldsymbol{\theta}}(\mathbf{x}_i)\}_{i=1}^N$ and ground truth signal values $\{g(\mathbf{x}_i)\}_{i=1}^N$ with certain metrics, we actually obtain an optimization objective for tuning $\boldsymbol{\theta}$. Benefiting from the natural differentiability of MLPs, we are able to minimize the error using gradient descent (GD) implemented by auto differentiation frameworks.

MLPs are known as a powerful approximator of functions and with their continuous property, we are able to interpolate the discrete signal up to an arbitrarily higher resolution. However, given the traditional training scheme, the resulting ability to interpolate at unseen coordinates is usually not satisfactory. In this paper, we alleviate this problem by leveraging the supervision on derivatives, which is to say, we optimize $\boldsymbol{\theta}$ by supervising not only signal values, but also derivatives at given coordinates.

The overall training paradigm is illustrated in Fig. 2. We therefore construct a training dataset $\{(\mathbf{x}_i, g(\mathbf{x}_i), D_{\mathbf{x}}g(\mathbf{x}_i))\}_{i=1}^N$ for INRs. Since the analytical expression of g is unknown, its derivatives require to be approximated with finite differences, which will be elaborated in Sec. 3.2. Considering that the task is a standard regression task, we apply ℓ_2 error as the loss function. With the training set, the optimization objective evolves to a combination of both supervision on signal values and first-order derivatives:

$$\begin{aligned} \boldsymbol{\theta} &= \arg \min_{\boldsymbol{\theta}} \frac{1}{N} \sum_{i=1}^N [\mathcal{L}_{\text{val}}(g(\mathbf{x}_i), f_{\boldsymbol{\theta}}(\mathbf{x}_i)) + \lambda \mathcal{L}_{\text{der}}(D_{\mathbf{x}}g(\mathbf{x}_i), D_{\mathbf{x}}f_{\boldsymbol{\theta}}(\mathbf{x}_i))] \\ &= \arg \min_{\boldsymbol{\theta}} \frac{1}{N} \sum_{i=1}^N (\|g(\mathbf{x}_i) - f_{\boldsymbol{\theta}}(\mathbf{x}_i)\|_2^2 + \lambda \|D_{\mathbf{x}}g(\mathbf{x}_i) - D_{\mathbf{x}}f_{\boldsymbol{\theta}}(\mathbf{x}_i)\|_2^2). \end{aligned} \quad (1)$$

We additionally extend this paradigm to enable both pre-processing of \mathbf{x} and post-processing of $f_{\boldsymbol{\theta}}(\mathbf{x})$, for instances that INRs are used to be an intermediate representation where the inputs to the MLP are calculated from \mathbf{x} and the outputs of the MLP require further processing to obtain the final results. To formulate, we parameterize the MLP as $f_{\boldsymbol{\theta}}$, and denote the pre-processing and post-processing as p and q respectively. The target signal g is therefore approximated as $q(f_{\boldsymbol{\theta}}(p(\mathbf{x})))$. Assuming that both p and q are differentiable and non-parametric, we can integrate these two processes into $f_{\boldsymbol{\theta}}$ and follow the aforementioned training paradigm. We will carry out experiments on tasks satisfying this configuration in Sec. 4.2. Specially, if p and q are identical mappings, the case degenerates to the original form.

Though can be applied under more circumstances, in this paper, we set the scope of target signals to $g : \mathbb{R}^2 \rightarrow \mathbb{R}^3$, that $\{\mathbf{x}_i\}_{i=1}^N$ are image coordinates and $\{g(\mathbf{x}_i)\}_{i=1}^N$ are pixel values (colors). The loss supervision on derivatives \mathcal{L}_{der} in Eq. (1) is composed of partial derivatives w.r.t. u and v respectively.

3.2 Approximate Derivatives with Finite Differences

To enable Sobolev training on image coordinates, the first-order derivative of the target signal g is supposed to be known as a precondition. However, the only information accessible at this stage is discrete sampled values $\{g(\mathbf{x}_i)\}_{i=1}^N$ so we need to approximate derivatives from these values with finite differences, or namely numerical differentiation (Fig. 2). Since we only consider cases that $\mathbf{x} \in \mathbb{R}^2$, the partial derivatives are

$$\begin{aligned} D_u g(u, v) &= \frac{g(u+h, v) - g(u-h, v)}{2h}, \\ D_v g(u, v) &= \frac{g(u, v+h) - g(u, v-h)}{2h}. \end{aligned} \quad (2)$$

In this paper, we mainly apply the Sobel operator to obtain first-order derivatives of images. It defines a template that convolves the image \mathbf{I} to obtain the partial derivative w.r.t. u while the transposed template for the partial derivative w.r.t. v :

$$D_u = \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix} * \mathbf{I}, D_v = \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{bmatrix} * \mathbf{I}, \quad (3)$$

where $*$ denotes the 2D convolution operation. We also study the choice of image filters for derivatives in Sec. 4.3, the conclusion of which is that enforcing the supervision of derivatives matters, regardless of the specific choice of filters.

3.3 Activation Functions

Activation functions are usually non-linear functions applied after each fully-connected layer. ReLU (Rectified Linear Units) [2] is a universally applied choice for its simple design, strong biological motivations and mathematical justifications. Due to the fact that the ReLU function is piecewise linear and its second derivative is zero everywhere, ReLU can not fit derivatives well in practice even the positional encoding is applied, just as shown in Fig. 3a. On the other hand, [27,31] attempt to explore the potential of the sine function as activation functions due to its periodicity, boundedness, arbitrary-order differentiability. The superior convergence property of sine has been proved in [31] when supervising the derivative only. In our setting, we try to replace ReLU with sine for better convergence property in Sobolev space. As can be indicated from Fig. 3a, sine indeed converges to a smaller derivative difference (close to 0) than ReLU when the training in Sobolev space is enabled.

However, studying the properties of activation functions in-depth is a complicated problem and in this paper, we mainly study the training paradigm applied to MLPs with activation functions of (a) ReLU, whose first-order derivative is piece-wise smooth, and (b) sine, whose arbitrary-order derivative is itself (with phase shift). We also study other activation functions in the Supplementary Material.

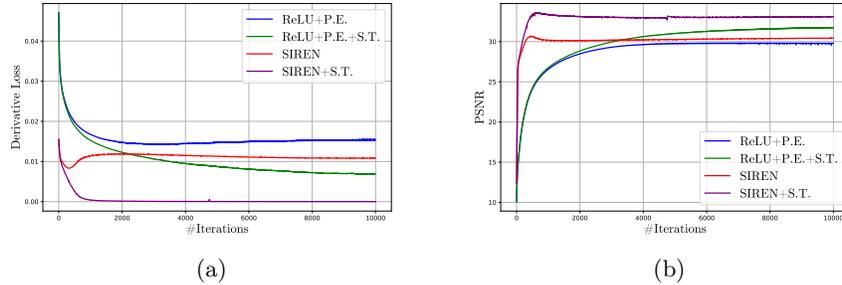


Fig. 3: Convergence situations of \mathcal{L}_{der} (a) and PSNR (b) w.r.t. iterations.

4 Experiments

In this section, we study the effectiveness of the proposed training paradigm on specific tasks that satisfy the preconditions. We sort tasks involving INRs into two categories, namely direct and indirect ones. Direct tasks are those where the supervision labels are in the same space as the network outputs, where p and q defined in Sec. 3.1 are both identical mappings. While in indirect tasks, the network outputs are processed afterwards so that the observations in the same space as the supervision labels are produced. We give an example of direct and indirect tasks respectively in Sec. 4.1 and Sec. 4.2.

4.1 Direct: Image Regression

A typical direct task is image regression. Given an RGB image \mathbf{I} , an MLP is appointed to regress the mapping from image coordinates to RGB colors $g : \mathbb{R}^2 \rightarrow \mathbb{R}^3$. This task is fundamental regarding demonstrating the representation ability of INRs.

Implementation

Dataset We construct a dataset based on the Set 5 dataset [5] to demonstrate the great interpolatability of MLPs when trained with additional supervision on derivatives. We split the pixels within an image into two groups, namely training samples and evaluation samples, by performing nearest-neighbor downscaling by a factor of 4. This is to say, for an image patch of 4×4 , we have 1 training sample and the remaining 15 pixels as evaluation samples.

Network Architecture For network architecture, we implement a fully-connected MLP with 4 activated linear layers with 256 hidden units each and 1 output linear layer. For the ReLU-based implementation, we follow the conclusion in [43,36] and apply positional encoding (P.E.) to the inputs (additional 20 channels).

Table 1: Quantitative results on the task of image regression on Set 5 dataset [5]. When trained with additional supervision on derivatives and activated by sine and ReLU functions, the network achieves the best and second-best performance respectively in terms of both PSNR and SSIM. We also provided interpolated results by bilinear and bicubic interpolation. **best** **second-best**

Method	Mean	Baby	Bird	Butterfly	Head	Woman
PSNR \uparrow						
Bilinear	24.14	26.98	24.88	18.68	28.04	22.10
Bicubic	23.63	26.52	24.46	18.23	27.43	21.52
ReLU+P.E. [24]	26.59	29.80	28.18	20.54	29.47	24.95
ReLU+P.E.+S.T.	28.63	31.93	31.22	21.93	31.11	26.97
SIREN [31]	26.22	30.48	28.07	20.73	28.22	23.62
SIREN+S.T.	30.28	33.36	33.34	24.42	31.64	28.62
SSIM \uparrow						
Bilinear	0.716	0.778	0.754	0.636	0.674	0.736
Bicubic	0.705	0.773	0.746	0.625	0.655	0.727
ReLU+P.E. [24]	0.740	0.792	0.808	0.636	0.690	0.775
ReLU+P.E.+S.T.	0.809	0.867	0.888	0.708	0.757	0.827
SIREN [31]	0.722	0.849	0.804	0.698	0.635	0.625
SIREN+S.T.	0.890	0.918	0.955	0.870	0.795	0.910

Training & Evaluation We train the network with all training samples with Adam optimizer for 50k iterations, at a learning rate of 10^{-4} . Following the common practice in image super resolution [17,34,38], we measure PSNR and SSIM on the Y channel and ignore 4 pixels from the image border. For PSNR, we only take the evaluation samples into account.

Results The quantitative scores are shown in Tab. 1, where we compare the proposed training paradigm with the previous value-based training paradigm, and with different activation functions. When trained with \mathcal{L}_{val} alone, ReLU and sine (SIREN [31]) functions lead to comparable results; while with additional supervision on image derivatives, a huge performance improvement has been gained with both activations in both metrics. It is worth noting that in all experiments, INRs demonstrate greater interpolatability than rule-based interpolation methods, i.e., bilinear and bicubic interpolation. Fig. 4 gives some example regions where our method produces clearer images and sharper edges and also approximates better the derivatives. It can be observed that sine-based MLPs outperform ReLU-based ones when Sobolev training is enabled. Even so, Sobolev training still improves the results of ReLU-based INRs by a large margin. As additional evidences, Fig. 3b shows the growth of PSNR w.r.t. the number of iterations, where the combination of Sobolev training and sine leads to the fastest convergence and the best performance.

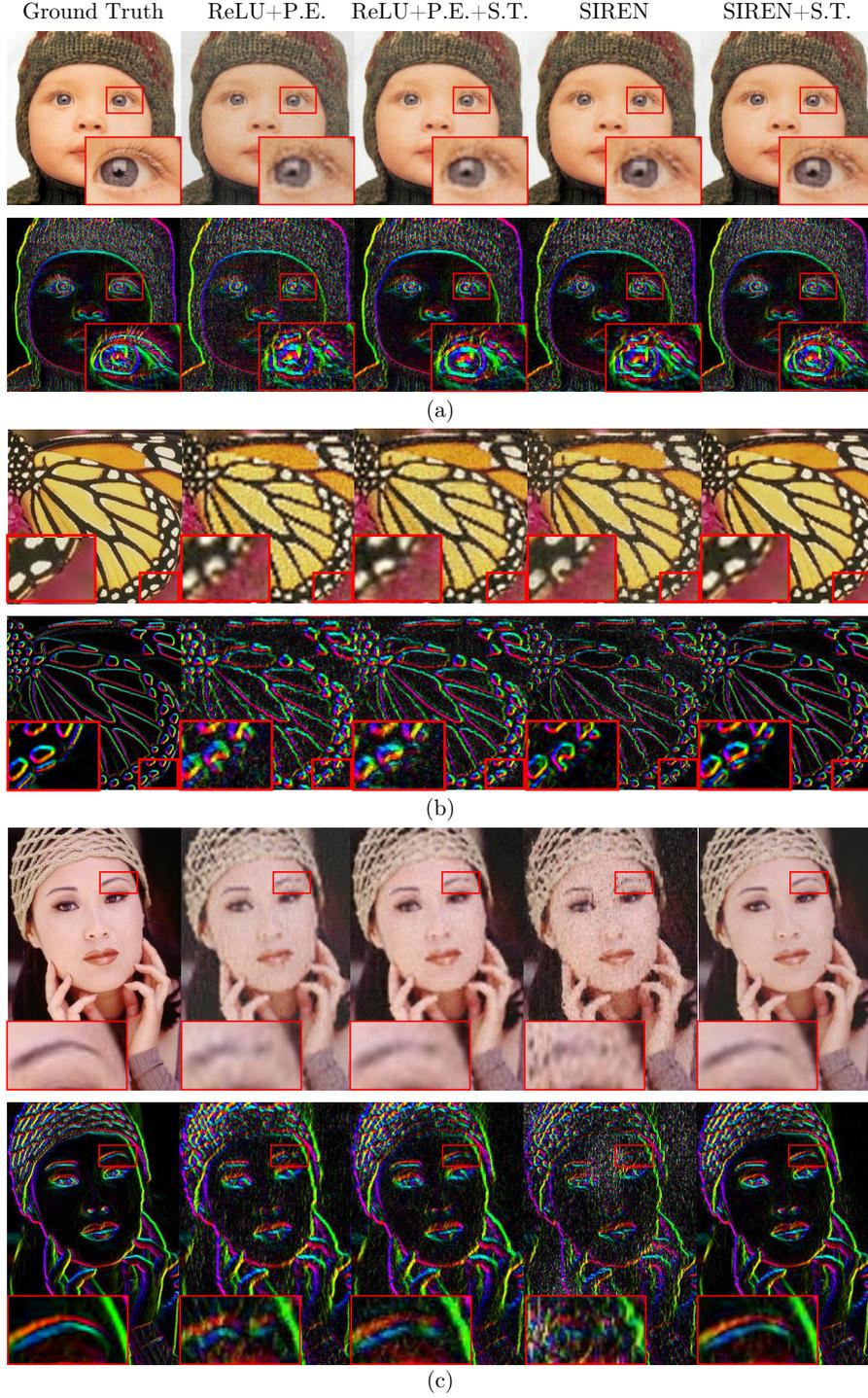


Fig. 4: Results of regressed images and derivatives of the set *Baby* (a), *Butterfly* (b) and *Woman* (c). When sine functions are adopted for activating linear layers and Sobolev training is enabled, the MLP yields the best visual effects, especially at high-frequency sharp edges, while ReLU-based MLPs fail to regress these details.

Table 2: Quantitative results on the task of inverse rendering on LLFF dataset [23]. When trained with additional supervision on derivatives and activated by sine functions, the network achieves the best performance in terms of both PSNR and SSIM. **best** **second-best**

Method	Mean	<i>Fern</i>	<i>Flower</i>	<i>Fortress</i>	<i>Horns</i>	<i>Leaves</i>	<i>Orchids</i>	<i>Room</i>	<i>T-Rex</i>
PSNR \uparrow									
ReLU+P.E. [24]	23.42	21.82	25.03	27.89	24.75	18.14	18.92	27.22	23.62
ReLU+P.E.+S.T.	23.74	22.99	25.27	27.90	24.50	19.30	19.40	27.10	23.49
SIREN[31]+P.E.	23.38	21.61	25.30	27.60	24.73	18.22	18.99	26.89	23.69
SIREN+P.E.+S.T.	24.38	23.63	25.99	28.30	25.24	20.02	19.74	27.76	24.40
SSIM \uparrow									
ReLU+P.E. [24]	0.706	0.651	0.755	0.769	0.729	0.535	0.544	0.877	0.791
ReLU+P.E.+S.T.	0.708	0.680	0.757	0.762	0.711	0.559	0.554	0.878	0.768
SIREN [31]+P.E.	0.682	0.596	0.739	0.763	0.720	0.495	0.525	0.851	0.768
SIREN+P.E.+S.T.	0.751	0.721	0.800	0.811	0.756	0.629	0.592	0.893	0.807

4.2 Indirect: Inverse Rendering

INRs are recently attention-drawing for their great power demonstrated when representing 3D scenes. NeRF [24], which is a typical example of the kind, leverages an MLP $f_{\theta} : \mathbb{R}^5 \rightarrow \mathbb{R}^4$, to regress the mapping from (x, y, z, θ, ϕ) to (r, g, b, σ) , where x, y, z are 3D coordinates and θ, ϕ are the view directions while the outputs are emitted colors $\mathbf{c} = (r, g, b)$ and volume density σ . In the pipeline of NeRF, the 3D sampling points (x, y, z) , as well as view directions (θ, ϕ) , of the MLP are determined from 2D image coordinates and corresponding calibrated camera parameters. After obtaining (r, g, b, σ) for each sampling point, a numerical approximation of volume rendering is applied to obtain the integral colors. Considering that both pre-processing and post-processing are differentiable, we can adapt the extended formulation of proposed training paradigm mentioned in Sec. 3.1 to perform inverse rendering of scenes.

Implementation

Dataset We carry out experiments on LLFF dataset [23]. Similarly, we down-sample the training images to 189×252 while the resolution of target images to render is 756×1008 .

Network Architecture We implement a simplified version of NeRF [24] $(x, y, z) \rightarrow (r, g, b, \sigma)$, removing the skip connection, the strategy of hierarchical sampling and view dependence. For network architecture, we implement a fully-connected MLP with 8 activated linear layers with 256 hidden units each and 1 output layer. We apply positional encoding (additional 60 channels) to both sine-activated and ReLU-activated MLPs.

Table 3: Results for comparisons on the task of inverse rendering on LLFF dataset [23] when the training resolution is larger. It is worth noting that even only trained with 1/16 samples, the MLP trained with the proposed paradigm is able to outperform the ReLU-activated network on both metrics in the set of *Flower*.

$H \times W$	Method	Mean	<i>Fern</i>	<i>Flower</i>	<i>Fortress</i>	<i>Horns</i>	<i>Leaves</i>	<i>Orchids</i>	<i>Room</i>	<i>T-Rex</i>
PSNR \uparrow										
756×1008	ReLU+P.E.	24.69	24.26	25.92	28.58	25.32	20.39	19.99	28.54	24.52
378×504	ReLU+P.E.	24.55	23.93	25.90	28.46	25.22	20.20	19.97	28.30	24.42
378×504	SIREN+P.E.+S.T.	24.67	24.27	26.16	28.37	25.37	20.64	19.94	27.99	24.61
189×252	SIREN+P.E.+S.T.	24.38	23.63	25.99	28.30	25.24	20.02	19.74	27.76	24.40
SSIM \uparrow										
756×1008	ReLU+P.E.	0.755	0.744	0.788	0.792	0.746	0.646	0.606	0.900	0.816
378×504	ReLU+P.E.	0.750	0.734	0.785	0.790	0.744	0.637	0.603	0.897	0.813
378×504	SIREN+P.E.+S.T.	0.766	0.753	0.805	0.815	0.762	0.666	0.614	0.899	0.815
189×252	SIREN+P.E.+S.T.	0.751	0.721	0.800	0.811	0.756	0.629	0.596	0.893	0.807

Training & Evaluation At the stage of training, we train the network with a batch size of 128 with Adam optimizer for 400k iterations, at a learning rate of 5×10^{-4} . For evaluation metrics, we report PSNR and SSIM on all channels (RGB) of the rendered images.

Results Tab. 2 demonstrates the PSNR and SSIM on all scenes of LLFF dataset [23] and Fig. 5 provides some results of novel view synthesis. The conclusions drawn are similar to the task of image regression.

Data-efficiency To provide more insights of the data-efficiency brought about by the proposed training paradigm, we also conduct experiments with different image resolution, namely different amount of training samples under different training paradigms. For the value-based paradigm, we train MLPs with images of resolution 378×504 and 756×1008 respectively and render novel views of 756×1008 . For the proposed paradigm, we additionally train MLPs with images of 378×504 and render novel views of 756×1008 .

The quantitative comparisons demonstrating the data-efficiency of Sobolev training are shown in Tab. 3. The following two points can be observed and indicated. (a) Even with only 1/16 training samples, the Sobolev trained MLP sometimes exceeds the ReLU-activated MLP with the value-based training paradigm, e.g., on *Flower* in terms of PSNR and SSIM, on *Fortress* and *Horns* in terms of SSIM. (b) The Sobolev trained MLP using images of 378×504 gets similar PSNR and even higher SSIM compared to the ReLU-activated MLP trained with the value-based paradigm using images of 756×1008 . Note that the former case is trained only with 1/4 samples as the latter.

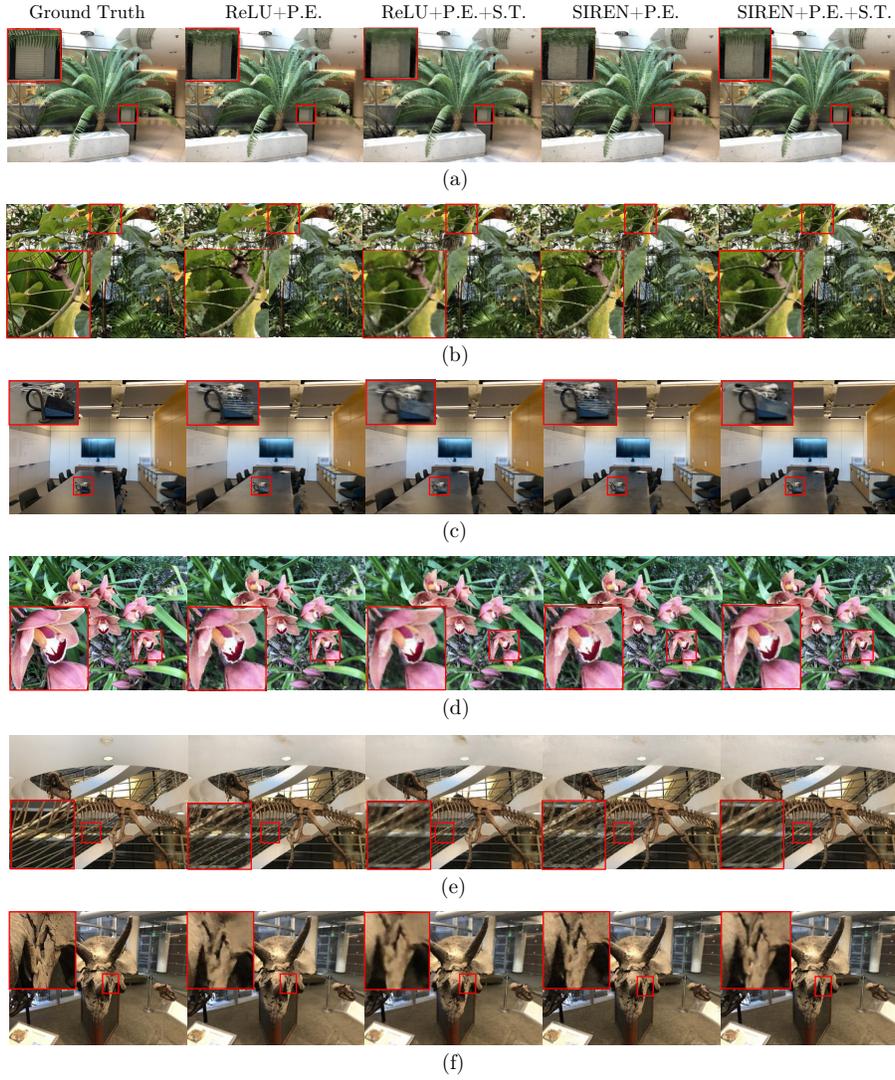


Fig. 5: Results of novel view synthesis by INRs in the task of inverse rendering. Top to bottom: *Fern* (a), *Leaves* (b), *Room* (c), *Orchids* (d), *T-Rex* (e), *Horns* (f). Supervision on derivatives, applied with sine-activated MLPs, helps to render clear images and sharp edges at unseen views.

4.3 Ablation Study

Image Derivative Filters We compare the network performance with another image filter as finite differences approximating derivatives, namely vanilla image

Table 4: Quantitative comparison on the choice of image filters when approximating derivatives of g . Experiments are carried out with MLPs activated by sine functions. PSNR and SSIM stand for mean values over all scenes.

Task	Filter	PSNR	SSIM
Image	Vanilla	30.24	0.886
Regression	Sobel	30.28	0.890
Inverse	Vanilla	24.28	0.746
Rendering	Sobel	24.38	0.751

Table 5: Quantitative comparison on the number of layers in MLPs under the task of inverse rendering. PSNR and SSIM stand for mean values over all scenes.

#layers	Method	PSNR	SSIM
4	ReLU+P.E.	23.29	0.685
	SIREN+P.E.	23.17	0.654
	ReLU+P.E.+S.T.	23.35	0.689
	SIREN+P.E.+S.T.	23.66	0.697
8	ReLU+P.E.	23.42	0.706
	SIREN+P.E.	23.38	0.682
	ReLU+P.E.+S.T.	23.74	0.708
	SIREN+P.E.+S.T.	24.38	0.751

derivative filter whose templates are

$$T_u = \begin{bmatrix} 0 & 0 & 0 \\ -1 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}, T_v = \begin{bmatrix} 0 & -1 & 0 \\ 0 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}. \quad (4)$$

It considers the same neighborhood as Sobel operator but removes the spatial smoothing operation included in Sobel operator. The quantitative comparisons are shown in Tab. 4.

Though the results with Sobel operator are slightly better than the vanilla image derivative filter, the performance gap is relatively small, indicating the specific choice of filters is not crucial to the performance.

Network Capacity A simple intuition towards the fitting ability, including both values and derivatives, of a neural network is that the more parameters a network possesses, the more complex functions it can represent. We present the quantitative results obtained with networks with different number of layers (4 and 8 respectively) in the task of inverse rendering in Tab. 5. Other experiment settings are aligned with Sec. 4.2. It can be indicated that when the layer number increases to 8 from 4, the representation ability of INRs is stronger so it can obtain more performance gain from derivatives.

5 Discussions

5.1 Future Work

To restate, the preconditions of the proposed training paradigm are that (a) the raw dataset is in the form of $\{(\mathbf{x}_i, g(\mathbf{x}_i))\}_{i=1}^N$, where $\{\mathbf{x}_i\}_{i=1}^N$ are spatial coordinates; (b) the derivatives of g can be obtained or approximated with limited $\{g(\mathbf{x}_i)\}$; (c) $\{\mathbf{x}_i\}$ do not necessarily be the direct inputs of the MLP, and the

target values $\{g(\mathbf{x}_i)\}$ do not necessarily be the direct outputs if and only if the pre-processing p is differentiable w.r.t. to \mathbf{x} and the post-processing q is differentiable w.r.t. $f_{\Theta}(p(\mathbf{x}))$.

Theoretically, referring to the formulation in Sec. 3.1, the training paradigm can further generalize to more tasks, e.g., audio regression (please refer to the Supplementary Material) and video regression, as long as their partial derivatives can be obtained by numerical approximation. Take the task of video regression $(u, v, t) \rightarrow (r, g, b)$ as an example, whose partial derivatives w.r.t. u and v are the same on images. We can approximate the partial derivative w.r.t. t by subtracting consecutive frames.

5.2 Limitations

We summarize the limitations of the proposed training paradigm as below.

- The computation of analytical derivatives is computationally expensive and occupies a considerable amount of memory. For example, both PyTorch [30] and TensorFlow [1] cannot calculate derivatives for each sample in a batch separately, leading to redundant computation.
- Though the proposed training paradigm can be applied to various INRs-based tasks, the scope of application is still limited. Since the derivatives are approximated with finite differences, the distribution of $\{\mathbf{x}_i\}$ is supposed to be as uniform as possible at its domain, which is not always satisfied.

6 Conclusion

In this paper, we put forward a novel training paradigm for INRs in Sobolev space, where supervision on both signal values and first-order derivatives is enforced for network training. The formulation of the task is generalized to any circumstance that the inputs are 2D image coordinates and the outputs are image values (e.g., RGB) as long as the pre-processing is differentiable w.r.t. the coordinates and the post-processing is differentiable w.r.t. the MLP’s outputs. We obtain approximated image derivatives by the Sobel operator.

Experiments are carried out on the task of image regression (direct) and inverse rendering (indirect) and results show that the training paradigm brings remarkable improvement to the quality of reconstructed images, especially when applied together with an MLP whose activation functions are sine.

In addition, we study some important and interesting peripheral problems relevant to the proposed training paradigm and give insights based on our observations.

Acknowledgements The authors would like to thank Zhongtian Zheng from Peking University for his valuable advice.

References

1. Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., et al.: {TensorFlow}: a system for {Large-Scale} machine learning. In: 12th USENIX symposium on operating systems design and implementation (OSDI 16). pp. 265–283 (2016) [2](#), [14](#)
2. Agarap, A.F.: Deep learning using rectified linear units (relu). arXiv preprint arXiv:1803.08375 (2018) [6](#), [21](#)
3. Agustsson, E., Timofte, R.: Ntire 2017 challenge on single image super-resolution: Dataset and study. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops (July 2017) [21](#)
4. Atzmon, M., Lipman, Y.: Sal: Sign agnostic learning of shapes from raw data. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 2565–2574 (2020) [3](#)
5. Bevilacqua, M., Roumy, A., Guillemot, C., Alberi-Morel, M.L.: Low-complexity single-image super-resolution based on nonnegative neighbor embedding (2012) [7](#), [8](#), [18](#), [20](#)
6. Bradbury, J., Frostig, R., Hawkins, P., Johnson, M.J., Leary, C., Maclaurin, D., Necula, G., Paszke, A., VanderPlas, J., Wanderman-Milne, S., Zhang, Q.: JAX: composable transformations of Python+NumPy programs (2018), <http://github.com/google/jax> [2](#)
7. Chabra, R., Lenssen, J.E., Ilg, E., Schmidt, T., Straub, J., Lovegrove, S., Newcombe, R.: Deep local shapes: Learning local sdf priors for detailed 3d reconstruction. In: European Conference on Computer Vision. pp. 608–625. Springer (2020) [3](#)
8. Chen, A., Xu, Z., Zhao, F., Zhang, X., Xiang, F., Yu, J., Su, H.: Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 14124–14133 (2021) [3](#)
9. Chen, Y., Liu, S., Wang, X.: Learning continuous image representation with local implicit image function. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8628–8638 (2021) [2](#)
10. Chen, Z., Zhang, H.: Learning implicit fields for generative shape modeling. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5939–5948 (2019) [3](#)
11. Clevert, D.A., Unterthiner, T., Hochreiter, S.: Fast and accurate deep network learning by exponential linear units (elus). arXiv preprint arXiv:1511.07289 (2015) [21](#)
12. Czarnecki, W.M., Osindero, S., Jaderberg, M., Swirszcz, G., Pascanu, R.: Sobolev training for neural networks. *Advances in Neural Information Processing Systems* **30** (2017) [4](#)
13. Glorot, X., Bengio, Y.: Understanding the difficulty of training deep feedforward neural networks. In: Proceedings of the thirteenth international conference on artificial intelligence and statistics. pp. 249–256. JMLR Workshop and Conference Proceedings (2010) [21](#)
14. Gropp, A., Yariv, L., Haim, N., Atzmon, M., Lipman, Y.: Implicit geometric regularization for learning shapes. In: International Conference on Machine Learning. pp. 3789–3799. PMLR (2020) [3](#), [4](#)
15. He, K., Zhang, X., Ren, S., Sun, J.: Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In: Proceedings of the IEEE international conference on computer vision. pp. 1026–1034 (2015) [18](#), [21](#)

16. Hornik, K.: Approximation capabilities of multilayer feedforward networks. *Neural networks* **4**(2), 251–257 (1991) [3](#), [4](#)
17. Jo, Y., Kim, S.J.: Practical single-image super-resolution using look-up table. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 691–700 (2021) [8](#)
18. Klambauer, G., Unterthiner, T., Mayr, A., Hochreiter, S.: Self-normalizing neural networks. *Advances in neural information processing systems* **30** (2017) [21](#)
19. Lin, C.H., Ma, W.C., Torralba, A., Lucey, S.: Barf: Bundle-adjusting neural radiance fields. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 5741–5751 (2021) [3](#)
20. Martin-Brualla, R., Radwan, N., Sajjadi, M.S., Barron, J.T., Dosovitskiy, A., Duckworth, D.: Nerf in the wild: Neural radiance fields for unconstrained photo collections. In: *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. pp. 7206–7215. IEEE Computer Society (2021) [3](#)
21. Mescheder, L., Oechsle, M., Niemeyer, M., Nowozin, S., Geiger, A.: Occupancy networks: Learning 3d reconstruction in function space. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 4460–4470 (2019) [1](#), [3](#)
22. Michalkiewicz, M., Pontes, J.K., Jack, D., Baktashmotlagh, M., Eriksson, A.: Implicit surface representations as layers in neural networks. In: *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*. pp. 4742–4751. IEEE Computer Society (2019) [3](#)
23. Mildenhall, B., Srinivasan, P.P., Ortiz-Cayon, R., Kalantari, N.K., Ramamoorthi, R., Ng, R., Kar, A.: Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)* **38**(4), 1–14 (2019) [10](#), [11](#), [18](#), [21](#), [22](#)
24. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. In: *European conference on computer vision*. pp. 405–421. Springer (2020) [1](#), [2](#), [3](#), [8](#), [10](#), [18](#), [19](#), [21](#), [22](#), [23](#)
25. Niemeyer, M., Mescheder, L., Oechsle, M., Geiger, A.: Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 3504–3515 (2020) [1](#)
26. Oechsle, M., Peng, S., Geiger, A.: Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In: *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. pp. 5589–5599 (2021) [3](#)
27. Parascandolo, G., Huttunen, H., Virtanen, T.: Taming the waves: sine as activation function in deep neural networks (2017) [2](#), [6](#)
28. Park, J.J., Florence, P., Straub, J., Newcombe, R., Lovegrove, S.: DeepSDF: Learning continuous signed distance functions for shape representation. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 165–174 (2019) [1](#), [3](#)
29. Park, K., Sinha, U., Barron, J.T., Bouaziz, S., Goldman, D.B., Seitz, S.M., Martin-Brualla, R.: Nerfies: Deformable neural radiance fields. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 5865–5874 (2021) [3](#)
30. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* **32** (2019) [2](#), [14](#)

31. Sitzmann, V., Martel, J., Bergman, A., Lindell, D., Wetzstein, G.: Implicit neural representations with periodic activation functions. *Advances in Neural Information Processing Systems* **33**, 7462–7473 (2020) [1](#), [2](#), [3](#), [4](#), [6](#), [8](#), [10](#), [18](#), [19](#), [20](#), [21](#), [22](#), [23](#), [25](#), [26](#)
32. Sitzmann, V., Rezkikov, S., Freeman, B., Tenenbaum, J., Durand, F.: Light field networks: Neural scene representations with single-evaluation rendering. *Advances in Neural Information Processing Systems* **34**, 19313–19325 (2021) [1](#)
33. Sitzmann, V., Zollhöfer, M., Wetzstein, G.: Scene representation networks: Continuous 3d-structure-aware neural scene representations. *Advances in Neural Information Processing Systems* **32** (2019) [1](#)
34. Soh, J.W., Park, G.Y., Jo, J., Cho, N.I.: Natural and realistic single image super-resolution with explicit natural manifold discrimination. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 8122–8131 (2019) [8](#)
35. Suhail, M., Esteves, C., Sigal, L., Makadia, A.: Light field neural rendering. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 8269–8279 (2022) [1](#)
36. Tancik, M., Srinivasan, P., Mildenhall, B., Fridovich-Keil, S., Raghavan, N., Singhal, U., Ramamoorthi, R., Barron, J., Ng, R.: Fourier features let networks learn high frequency functions in low dimensional domains. *Advances in Neural Information Processing Systems* **33**, 7537–7547 (2020) [2](#), [3](#), [7](#), [20](#)
37. Tewari, A., Fried, O., Thies, J., Sitzmann, V., Lombardi, S., Xu, Z., Simon, T., Nießner, M., Tretschk, E., Liu, L., et al.: Advances in neural rendering. In: *ACM SIGGRAPH 2021*. pp. 1–320. ACM (2021) [3](#)
38. Yang, W., Zhang, X., Tian, Y., Wang, W., Xue, J.H., Liao, Q.: Deep learning for single image super-resolution: A brief review. *IEEE Transactions on Multimedia* **21**(12), 3106–3121 (2019) [8](#)
39. Yariv, L., Kasten, Y., Moran, D., Galun, M., Atzmon, M., Ronen, B., Lipman, Y.: Multiview neural surface reconstruction by disentangling geometry and appearance. *Advances in Neural Information Processing Systems* **33**, 2492–2502 (2020) [3](#)
40. Yen-Chen, L.: Nerf-pytorch. <https://github.com/yenchenlin/nerf-pytorch/> (2020) [22](#)
41. Yu, A., Li, R., Tancik, M., Li, H., Ng, R., Kanazawa, A.: Plenotrees for real-time rendering of neural radiance fields. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 5752–5761 (2021) [3](#)
42. Yu, A., Ye, V., Tancik, M., Kanazawa, A.: pixelnerf: Neural radiance fields from one or few images. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 4578–4587 (2021) [3](#)
43. Yüce, G., Ortiz-Jiménez, G., Besbinar, B., Frossard, P.: A structured dictionary perspective on implicit neural representations. *arXiv preprint arXiv:2112.01917* (2021) [7](#), [20](#)
44. Zhang, J., Yang, G., Tulsiani, S., Ramanan, D.: Ners: Neural reflectance surfaces for sparse-view 3d reconstruction in the wild. *Advances in Neural Information Processing Systems* **34** (2021) [3](#)
45. Zhang, J., Yao, Y., Quan, L.: Learning signed distance field for multi-view surface reconstruction. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 6525–6534 (2021) [3](#)

A Additional Details & Results

A.1 Image Regression

For image regression task, we set the angular velocity of the sinusoidal function to 30 and initialize the weights of MLPs following [31]. We provide other results of Set 5 dataset [5] in Fig. 6. We also report the memory consumption as well as training time required in Tab. 6. The proposed training paradigm requires more memory and time at the training phase for additional computation and storage of derivatives but it is worth noting that the time, as well as the memory consumption at the inference phase, is the same with or without the derivative supervision.

A.2 Inverse Rendering

For inverse rendering task, we set the angular velocity of the sinusoidal function to 1 and initialize the weights of MLPs according to the method described in [15]. LLFF dataset [23,24] consists of 8 scenes captured with a handheld cellphone, captured with 20 to 62 images. We follow [24] to hold out $\frac{1}{8}$ of images for the evaluation set. All of the training images are 756×1008 , but the dataset also provides the raw cellphone images of 3024×4032 , which will be used to evaluate high resolution rendering results in Appendix D.

We provide other results of LLFF dataset [23,24] in Fig. 7. Tab. 7 reports the training-time memory consumption and time. Similar to image regression, the additional cost is only at training-time.

B Audio Regression

As mentioned in the discussions of the main paper, the training paradigm we proposed can further generalize to more tasks, as long as the tasks satisfy the formulation in Sec. 3.1 of the main paper. Here we follow [31] to conduct the task of audio signal representation.

Implementation

Dataset Following [31], we use two different audio signals, one for music and one for speech. For music data, we use the first 7 seconds from Bach’s Cello Suite No.1 (*Bach*)*, and for the speech, we use stock audio of a male actor counting from 0 to 9 (*Counting***). Both audio signals share a sampling rate of $44100Hz$, and in total there are 308207 samples in *Bach* and 537936 samples in *Counting*. We normalize signal values to the range of $[-1, 1]$. Same as other tasks, we separate

* Audio file available at <https://www.yourclassical.org/episode/2017/04/04/daily-download-js-bach--cello-suite-no-1-prelude>.

** Audio file available at <http://soundbible.com/2008-0-9-Male-Vocalized.html>.

Table 6: GPU memory consumption at training phase and training time of 1000 epochs in image regression task of *Baby*. P.E. stands for positional encoding; S.T. stands for Sobolev training.

Method	Memory(MB)	Time(s)
ReLU+P.E. [24]	130	4.678
ReLU+P.E.+S.T.	316	20.548
ReLU+P.E.+S.T.*	401	14.081
SIREN [31]	176	3.977
SIREN+S.T.	753	22.963
SIREN+S.T.*	704	16.556

* The new PyTorch 1.11.0 supports computation of per-sample derivatives, so we also report the statistics obtained.

all samples within an audio clip into two groups, namely training samples and evaluation samples, by performing nearest-neighbor downscaling by a factor of 5. The approximated derivatives are obtained by two-sided differences method, similar to the vanilla derivative filter of the task on images.

Network Architecture We implement a fully-connected MLP with 4 activated linear layers with 256 hidden units and 1 output linear layer. We set the angular velocity of the sinusoidal function to 30 and follow the initialization strategy in [31] to initialize the weights of MLPs.

Training & Evaluation At the stage of training, we use the entire set of training samples to train the network as a batch with Adam optimizer, at a learning rate of 5×10^{-5} . For evaluation, we measure PSNR only on the evaluation samples.

Results The quantitative results in PSNR of regressing these two audio clips are shown in Tab. 8. Sobolev training significantly improves the regression quality of both audio clips. We further illustrate the visualized comparisons respectively in Fig. 8 and Fig. 9, showing the regression error of the derivative-trained network is smaller.

As our results show ReLU-activated networks can not fit audio signals well, which is consistent with the result of [31], we only report the results of Sine-activated networks.

C Different Activation Functions

Recent ReLU-based MLPs normally have poor convergence property under derivative supervision, we mainly investigate the performance difference between ReLU and periodic activation function, i.e., sine, in the main paper. In this part, an experiment is designed to study the convergence properties of derivatives with different activation functions. The scope of the investigation is shown in Tab. 9.

Table 7: GPU memory consumption at training phase and training time of 1000 iterations in inverse rendering task of *Fern*.

Method	Memory(MB)	Time(s)
ReLU+P.E.	228	87.264
ReLU+P.E.+S.T.	912	190.434
SIREN+P.E.	348	86.642
SIREN+P.E.+S.T.	1827	204.748

Table 8: Quantitative PSNR results on the task of audio regression on *Bach* and *Counting*. 1/5 samples are used for training, and rest samples are used for evaluation. When trained with additional supervision on derivatives, PSNR is significantly improved.

Method	Mean	<i>Bach</i>	<i>Counting</i>
SIREN [31]	35.89	41.50	30.28
SIREN+S.T.	41.23	48.89	33.57

Dataset We convert the RGB image *Bird* of Set 5 [5] to grayscale as our dataset. The original image shape is $288 \times 288 \times 1$, we use the Sobel operator to get approximate image partial derivatives, resulting the derivative image with shape $288 \times 288 \times 2$, w.r.t. u and v respectively. We perform nearest-neighbor down-sampling on the original image and derivative image by a factor of 4.

Network Architecture The network architecture is the same as the image regression task in the main paper, except for the channel number of the last linear layer is 1 instead of 3. The angular velocity of the sinusoidal function and weight initialization method are both the same as image regression task. Taking the conclusion in [36,43] into consideration, we do two parallel experiments with and without positional encoding for all activation functions except for sine.

Training & Evaluation We use the downsampled 1/16 pixels as our training data and the remaining pixels as evaluation samples. The network is optimized for 10k iterations at a learning rate of 10^{-4} .

Results Fig. 10 and Fig. 11 respectively show the convergence curve of derivative loss with and without positional encoding. Sine demonstrates its great power in fitting the network’s derivative compared with other activation functions. It is worth noting that positional encoding is helpful to improve the MLPs’ ability of approximating derivatives when applied with other activation functions. The corresponding quantitative results are shown in Tab. 10.

Table 9: Different activation functions and their definitions and weight initialization

Activation	Definition	Derivative	Initialization
ReLU [2]	$f(x) = \begin{cases} x, & x > 0 \\ 0, & x \leq 0 \end{cases}$	$f'(x) = \begin{cases} 1, & x > 0 \\ 0, & x < 0 \end{cases}$	Kaiming normal[15]
ELU [11]	$f(x) = \begin{cases} x, & x > 0 \\ \alpha(e^x - 1), & x \leq 0 \end{cases}$	$f'(x) = \begin{cases} 1, & x > 0 \\ \alpha e^x, & x < 0 \end{cases}$	Normal
SELU [18]	$f(x) = \begin{cases} \lambda x, & x > 0 \\ \lambda \alpha(e^x - 1), & x \leq 0 \end{cases}$	$f'(x) = \begin{cases} \lambda, & x > 0 \\ \lambda \alpha e^x, & x \leq 0 \end{cases}$	Normal
Sigmoid	$f(x) = \frac{1}{1+e^{-x}}$	$f'(x) = f(x)(1 - f(x))$	Xavier normal [13]
Softplus	$f(x) = \ln(1 + e^x)$	$f'(x) = \frac{1}{1+e^{-x}}$	Kaiming normal[15]
Tanh	$f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$	$f'(x) = 1 - (f(x))^2$	Xavier normal [13]
Sine [31]	$f(x) = \sin(x)$	$f'(x) = \cos(x)$	Specific uniform [31]

D Precision of Approximate Image Derivatives

In all aforementioned experiments, the partial derivatives we leverage for supervision are obtained by applying Sobel filters on the images before downsampling and for training, we perform nearest-neighbor downsampling for both image values and approximated derivatives. This is generally reasonable since in real applications, the original images are usually of very high resolution, such as LLFF dataset [23,24], whose original resolution is 3024×4032 . Without getting deteriorated rendering results, we would like the training data of INRs to be as few as possible for faster training and a direct solution is to downscale the images. As our results in the main paper show, by keeping the derivatives obtained at a high resolution, the problem of downgraded performance will get alleviated.

We also conduct experiments of image regression and inverse rendering where derivatives are obtained from the downsampled images, abandoning the data dependence of the raw data before downsampling, which is to say, we will use all samples in hand for training.

D.1 Image Regression

The dataset we used consists of 5 images of 1356×2040 , which are generated by uniformly sampling images from DIV2K [3] validation set. We downsample the original images by nearest-neighbor interpolation with a factor of 4 and use Sobel filters to calculate the approximate derivatives on downsampled images. The network architecture and training & evaluation settings are consistent with the image regression task of the main paper.

Results The quantitative results are shown in Tab. 11. From Tab. 11, we can see training with derivatives from downsampled images still gives a great performance improvement than value-based training paradigm.

Table 10: Quantitative comparisons of different activation functions on the task of (grayscale) image regression.

Activation	PSNR \uparrow	SSIM \uparrow
ReLU	25.06	0.711
ReLU P.E.	29.36	0.856
ELU	21.13	0.537
ELU P.E.	29.76	0.879
SELU	20.20	0.433
SELU P.E.	28.55	0.837
Sigmoid	18.06	0.408
Sigmoid P.E.	24.27	0.638
Softplus	18.38	0.417
Softplus P.E.	24.24	0.639
Tanh	23.27	0.635
Tanh P.E.	31.14	0.902
Sine	33.13	0.960

D.2 Inverse Rendering

The experiment settings are different from the main paper. Here we do not perform downsampling to images; instead, we use images of 756×1008 for training and render novel views of 756×1008 and 3024×4032 . The approximated image derivatives are calculated from images of 756×1008 . As mentioned earlier, LLFF dataset [23] provides raw images of 3024×4032 so we can evaluate rendering results at a higher resolution.

Results The quantitative results on rendering views of different resolutions are shown in Tab. 12. As can be proven by the results, we can still get a slight performance improvement when training with approximated derivatives from downsampled training images. Also, the performance gap is enlarged, when the rendering views’ resolutions are higher, indicating the great generalizability of Sobolev trained MLPs.

The difference between the two sources of image derivatives is basically the precision of approximated derivatives. As shown in the results of both tasks, the precision of approximated derivatives does not affect much. It is using the derivatives for supervision that matters.

E Use of Existing Assets

Some codes of image regression task and audio regression task are borrowed from SIREN [31]. The implementation of inverse rendering task are based on NeRF-PyTorch [40], which is a PyTorch version of original NeRF [24].

Table 11: Quantitative results on the task of image regression on DIV2K validation set.

Method	Mean	DIV2K Validation Set				
		0820	0840	0860	0880	0900
PSNR \uparrow						
ReLU+P.E. [24]	23.34	20.44	25.75	18.62	29.36	22.51
ReLU+P.E.+S.T.	23.91	20.90	26.30	19.17	30.00	23.18
SIREN [31]	22.69	19.70	24.87	17.23	29.21	22.45
SIREN+S.T.	24.44	21.58	26.87	19.35	30.64	23.74
SSIM \uparrow						
ReLU+P.E. [24]	0.577	0.497	0.651	0.371	0.817	0.547
ReLU+P.E.+S.T.	0.625	0.559	0.692	0.401	0.843	0.629
SIREN [31]	0.594	0.508	0.650	0.356	0.842	0.616
SIREN+S.T.	0.703	0.663	0.754	0.488	0.875	0.733

Table 12: Quantitative results of using different paradigms when training with images of resolution 756×1008 and rendering novel views of 756×1008 and 3024×4032 .

Method	Render $H \times W$	Mean	<i>Fern</i>	<i>Flower</i>	<i>Fortress</i>	<i>Horns</i>	<i>Leaves</i>	<i>Orchids</i>	<i>Room</i>	<i>T-Rex</i>
PSNR \uparrow										
ReLU+P.E. [24]	756×1008	24.69	24.26	25.92	28.58	25.32	20.39	19.991	28.54	24.52
	3024×4032	23.12	22.24	24.93	27.22	23.66	18.90	19.258	26.24	22.54
SIREN+P.E.+S.T.	756×1008	24.72	24.33	26.06	28.66	25.40	20.65	19.985	28.03	24.68
	3024×4032	23.22	22.32	25.07	27.26	23.72	19.10	19.261	26.14	22.88
SSIM \uparrow										
ReLU+P.E. [24]	756×1008	0.755	0.744	0.788	0.792	0.746	0.646	0.606	0.900	0.816
	3024×4032	0.723	0.714	0.773	0.819	0.703	0.570	0.616	0.860	0.732
SIREN+P.E.+S.T.	756×1008	0.767	0.755	0.807	0.819	0.760	0.666	0.614	0.899	0.818
	3024×4032	0.729	0.717	0.781	0.830	0.705	0.581	0.622	0.859	0.736

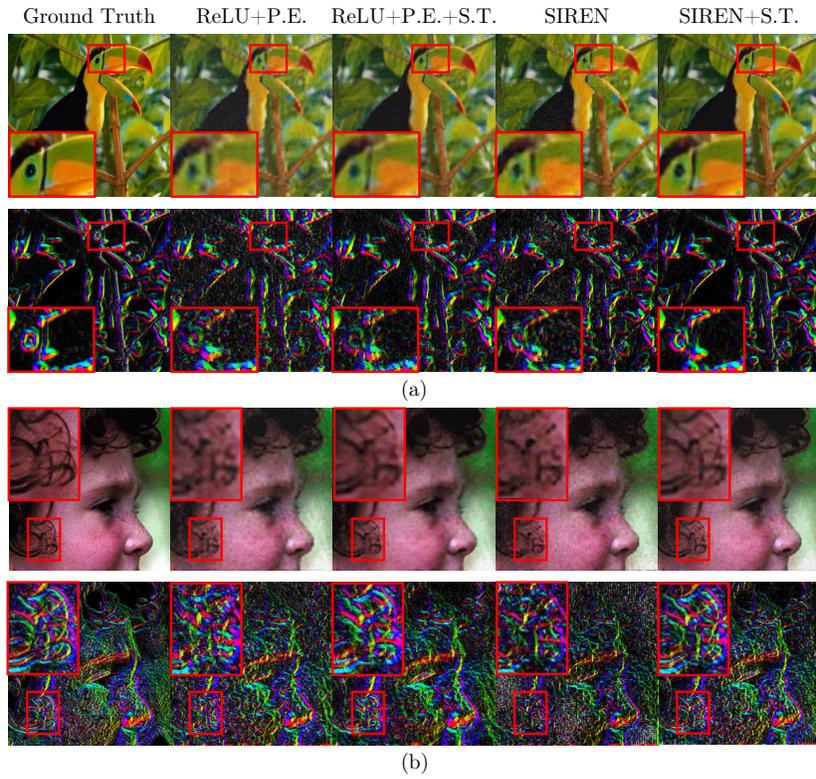


Fig. 6: Additional results of image regression. *Bird* (a), *Head* (b).

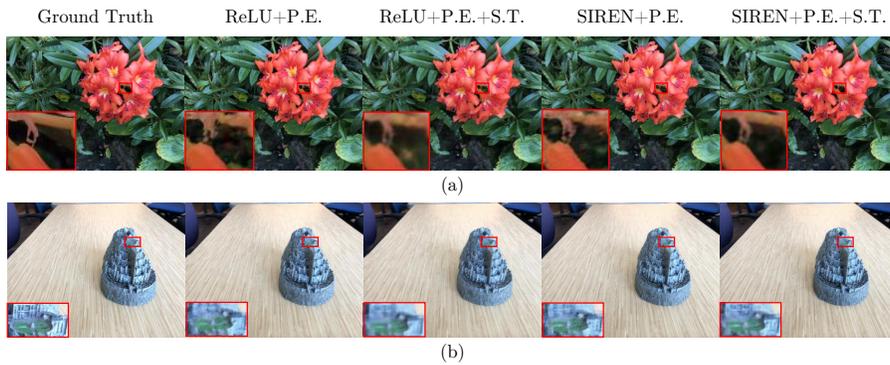
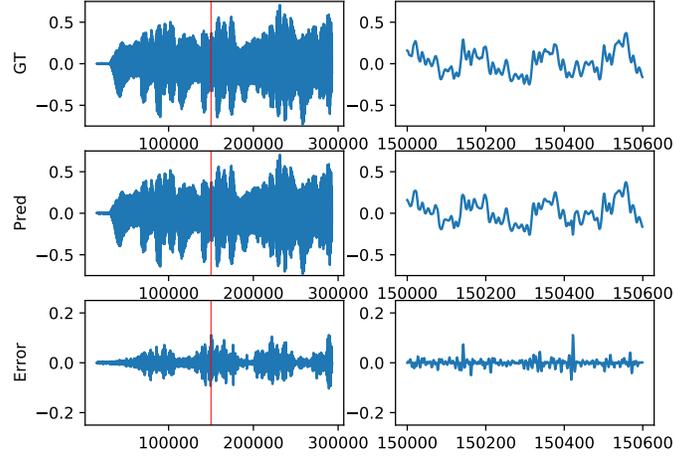
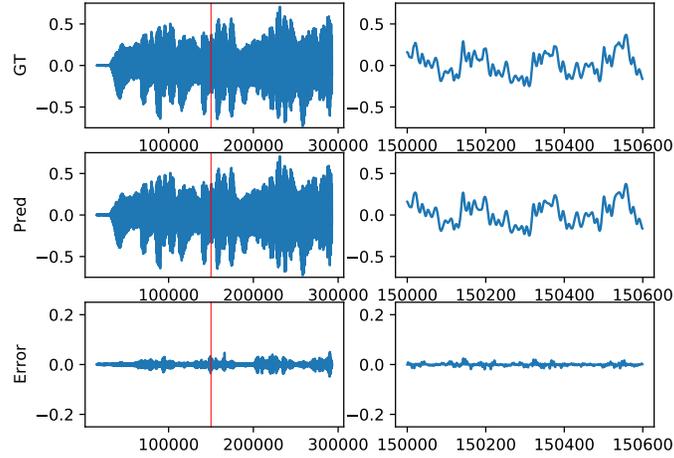


Fig. 7: Additional results of inverse rendering. *Flower* (a), *Fortress* (b).

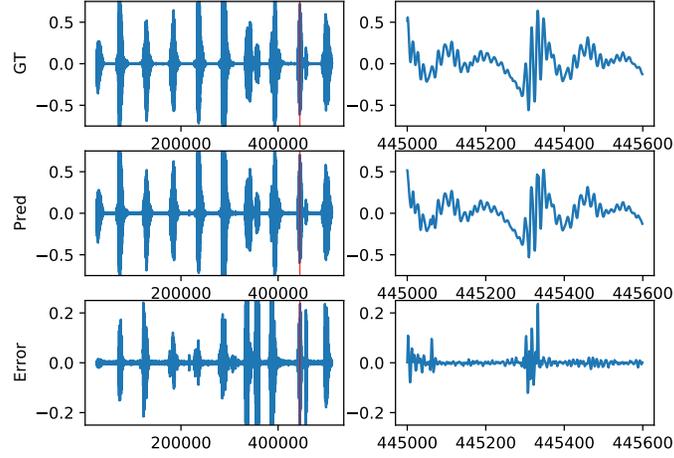


(a) SIREN [31]

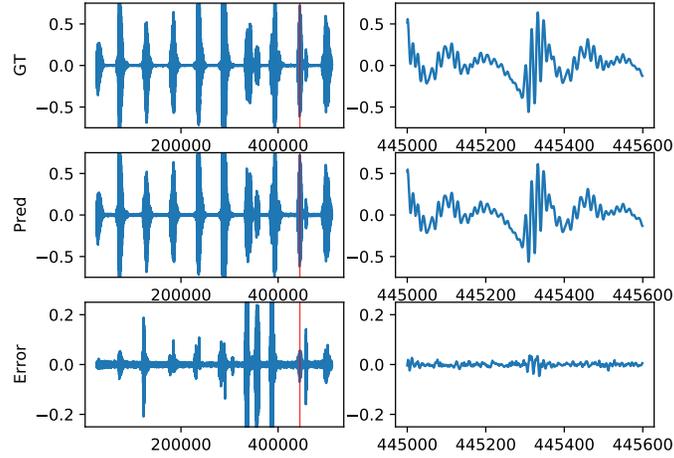


(b) SIREN+S.T.

Fig. 8: Regressed waveforms of the INRs trained with and without derivative supervision on the clip of *Bach*. From top to bottom: ground truth waveform, regressed waveform and error detected. We zoom in the waveform marked with red at the right column accordingly.



(a) SIREN [31]



(b) SIREN+S.T.

Fig. 9: Regressed waveforms of the INRs trained with and without derivative supervision on the clip of *Counting*. From top to bottom: ground truth waveform, regressed waveform and error detected. We zoom in the waveform marked with red at the right column accordingly.

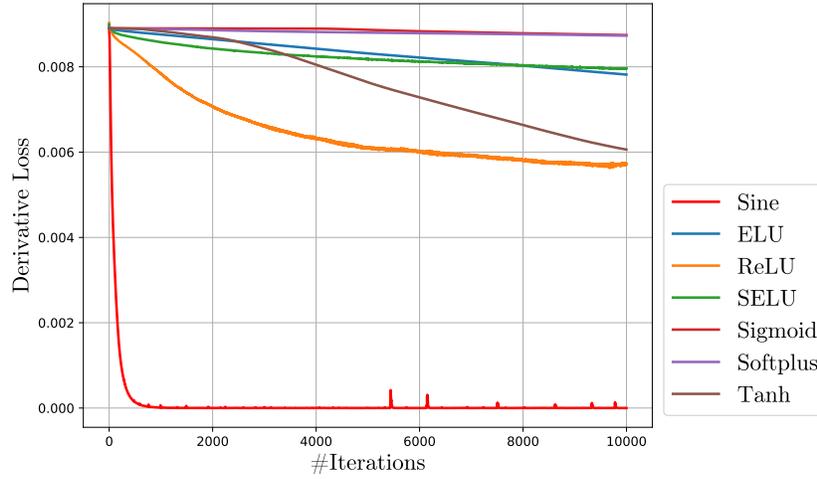


Fig. 10: Derivative loss of different activation functions.

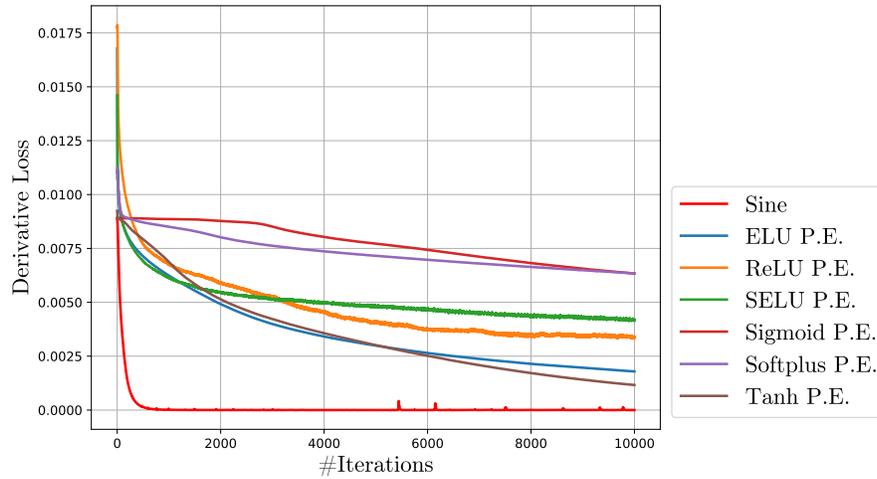


Fig. 11: Derivative loss of different activation functions with positional encoding. P.E. means positional encoding.