

DINER: Disorder-Invariant Implicit Neural Representation

Shaowen Xie^{1,†}, Hao Zhu^{1,†}, Zhen Liu^{1,†}, Qi Zhang², You Zhou¹, Xun Cao^{1,*}, Zhan Ma^{1,*}
¹ School of Electronic Science and Engineering, Nanjing University, Nanjing 210023, China
² AI Lab, Tencent Company, Shenzhen 518054, China
[†] Equal contribution. * Corresponding authors: {caoxun, mazhan}@nju.edu.cn

Abstract

Implicit neural representation (INR) characterizes the attributes of a signal as a function of corresponding coordinates which emerges as a sharp weapon for solving inverse problems. However, the capacity of INR is limited by the spectral bias in the network training. In this paper, we find that such a frequency-related problem could be largely solved by re-arranging the coordinates of the input signal, for which we propose the disorder-invariant implicit neural representation (DINER) by augmenting a hash-table to a traditional INR backbone. Given discrete signals sharing the same histogram of attributes and different arrangement orders, the hash-table could project the coordinates into the same distribution for which the mapped signal can be better modeled using the subsequent INR network, leading to significantly alleviated spectral bias. Experiments not only reveal the generalization of the DINER for different INR backbones (MLP vs. SIREN) and various tasks (image/video representation, phase retrieval, and refractive index recovery) but also show the superiority over the state-of-the-art algorithms both in quality and speed.

1. Introduction

INR [34] continuously describes a signal, providing the advantages of Nyquist-sampling-free scaling, interpolation, and extrapolation without requiring the storage of additional samples [20]. By combining it with differentiable physical mechanisms such as the ray-marching rendering [13, 22], Fresnel diffraction propagation [44] and partial differential equations [12], INR becomes a universal and sharp weapon for solving inverse problems and has achieved significant progress in various scientific tasks, e.g., the novel view synthesis [39], intensity diffraction tomography [19] and multiphysics simulation [12].

However, the capacity of INR is often limited by the underlying network model itself. For example, the spectral bias [28] usually makes the INR easier to represent low-frequency signal components (see Fig. 1 and Fig. 2(c)).

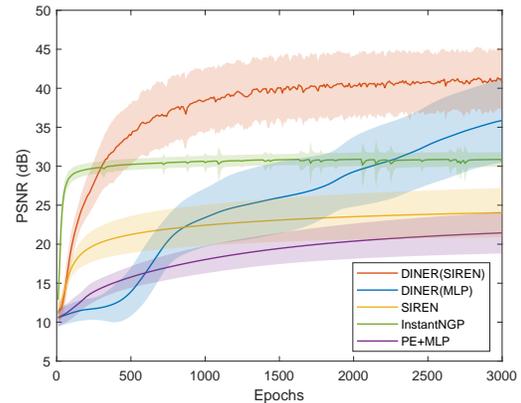


Figure 1. PSNR of various INRs on 2D image fitting over different training epochs.

To improve the representation capacity of the INR model, previous explorations mainly focus on encoding more frequency bases using either Fourier basis [22, 34, 38, 42] or wavelet basis [7, 17] into the network. However, the length of function expansion is infinite in theory, and a larger model with more frequency bases runs exceptionally slowly.

Such a problem is closely related to the input signal’s frequency spectrum. The signal’s frequency tells how fast the signal attribute changes following the intrinsic order of geometry coordinates. By properly re-arranging the order of coordinates of a discrete signal, we could modulate the signal’s frequency spectrum to possess more low-frequency components. Such a re-arranged signal can be better modeled using subsequent INR. Based on this observation, we propose the DINER. In DINER, the input coordinate is first mapped to another index using a hash-table and fed into a traditional INR backbone. We prove that no matter what orders the elements in the signal are arranged, the joint optimization of the hash-table and the network parameters guarantees the same mapped signal (Fig. 2(d)) with more low-frequency components. As a result, the representation capacity of the existing INR backbones and the task perfor-

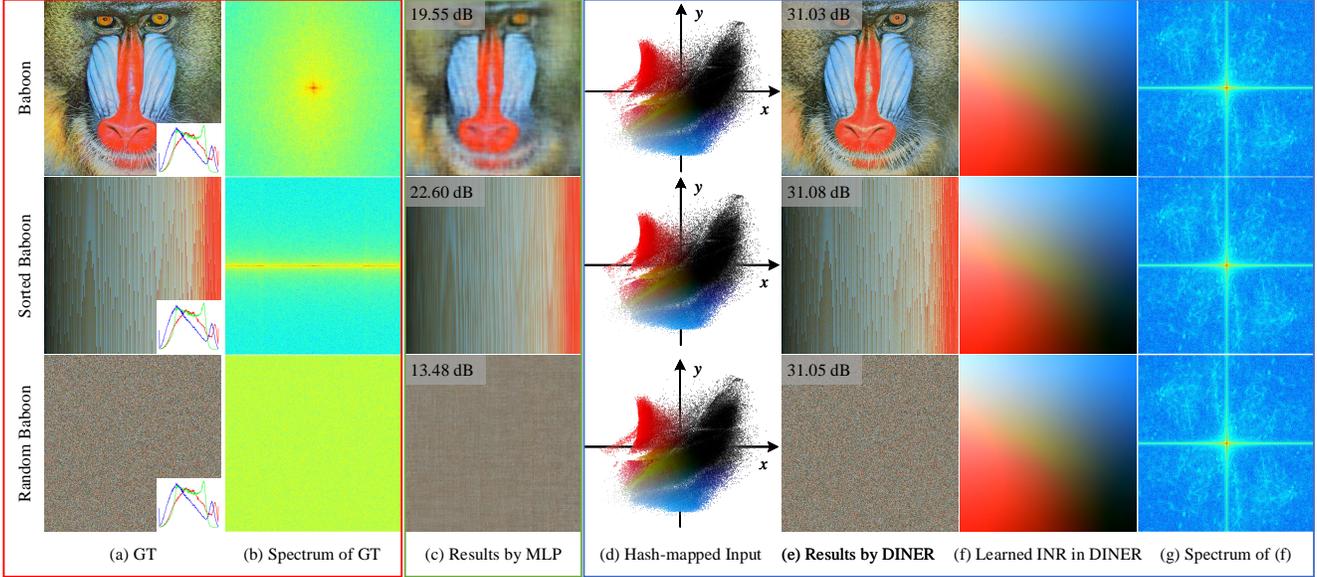


Figure 2. Comparisons of the existing INR and DINER for representing Baboon with different arrangements. From top to bottom, pixels in the Baboon are arranged in different geometric orders, while the histogram is not changed (the right-bottom panel in (a)). From left to right, (a) and (b) refer to the ground truth image and its Fourier spectrum. (c) contains results by an MLP with positional encoding (PE+MLP) [38] at a size of 2×64 . (d) refers to the hash-mapped coordinates. (e) refers to the results of the DINER that uses the same-size MLP as (c). (f) refers to the learned INR in DINER by directly feeding grid coordinates to the trained MLP (see Sec. 4.1 for more details). (g) is the Fourier spectrum of (f). (g) shares the same scale bar with (b).

mance are largely improved. As in Fig. 1 and Fig. 2(e), a tiny shallow and narrow MLP-based INR could well characterize the input signal with arbitrary arrangement orders. The use of hash-table trades the storage for fast computation where the caching of its learnable parameters for mapping is usually at a similar size as the input signal, but the computational cost is marginal since the back-propagation for the hash-table derivation is $\mathcal{O}(1)$.

Main contributions are summarized as follows:

1. The inferior representation capacity of the existing INR model is largely increased by the proposed DINER, in which a hash-table is augmented to map the coordinates of the original input for better characterization in the succeeding INR model.
2. The proposed DINER provides consistent mapping and representation capacity for signals sharing the same histogram of attributes and different arrangement orders.
3. The proposed DINER is generalized in various tasks, including the representation of 2D images and 3D videos, phase retrieval in lensless imaging, as well as 3D refractive index recovery in intensity diffraction tomography, reporting significant performance gains to the existing state-of-the-arts.

2. Related Work

2.1. INR and inverse problem optimization

INR (sometimes called coordinate neural network) builds the mapping between the coordinate and its signal value using a neural network, promising continuous and memory-efficient modeling for various signals such as the 1D audio [9], 2D image [38], 3D shape [26], 4D light field [35] and 5D radiance field [22]. Accurate INR for these signals could be supervised directly by comparing the network output with the ground truth, or an indirect way that calculates the loss between the output after differentiable operators and the variant of the ground truth signal. Thus, INR becomes a universal tool for solving inverse problems because the forward processes in these problems are often well-known. INR has been widely applied in the optimization of inverse problems in several disciplines, such as computer vision and graphics [39], computational physics [12], clinical medicine [41], biomedical engineering [19,44], material science [6] and fluid mechanics [29,30].

2.2. Encoding high-frequency components in INR

According to the approximation theory, an MLP network could approximate any function [15]. However, there is a spectral bias [28] in the network training, resulting in low performance of INR for high-frequency components. Several attempts have been explored to overcome this bias and

could be classified into two categories, *i.e.*, the function expansion and the parametric encodings.

The function expansion idea treats the INR fitting as a function approximation using different bases. Mildenhall et al. [22] encoded the input coordinates using a series of \sin / \cos functions with different frequencies and achieved great success in radiance field representation. The strategy of frequency-predefined \sin / \cos functions is further improved with random Fourier features, which has been proved to be effective in learning high-frequency components both in theory and in practical [38]. Sitzmann et al. [34] replaced the classical ReLU activation with periodic activation function (SIREN). The different layers of the SIREN could be viewed as increasing different frequency supports of a signal [42]. SIREN is suited for representing complex natural signals and their derivatives. Apart from the Fourier expansion, Fathony et al. [7] represented a complex signal by a linear combination of multiple wavelet functions (MFN), where the high-frequency components could be well modeled by modulating the frequency in the Gabor filter. Lindell et al. [17] developed MFN and proposed the band-limited coordinate networks, where the frequency at each network layer could be specified at initialization. These methods have achieved significant advantages in representing high-frequency components compared with the standard MLP network. However, *the performance of these INRs are limited by the frequency distribution of a signal-self, and often require a deeper or wider network architecture to improve the fitting accuracy.*

From the perspective of the parameter encoding [4, 11, 18, 37], each input coordinate is encoded using learned features which are fed into an MLP for fitting. Takikawa et al. [37] divided the 3D space using a sparse voxel octree structure where each point is represented using a learnable feature vector from its eight corners, achieving real-time rendering of high-quality signed distance functions. Martel et al. [20] divided the coordinate space iteratively during the INR training (ACORN), where the encoding features for each local block are obtained by a coordinate encoder network and are fed into a decoder network to obtain the attribute of a signal. ACORN achieves nearly 40 dB PSNR for fitting gigapixel images for the first time. Muller et al. [24] replaced the coordinate encoder network with a multi-resolution hash-table. Because the multi-resolution hash-table has higher freedom for characterizing coordinates' features, only a tiny network is used to map the features and the attribute values of a signal. Despite the superiority of faster convergence and higher accuracy in parameter encoding, one of the key questions is still not answered, *i.e.*, *what are the geometrical meanings of these features?*

Compared with these methods, the hash-table in the proposed DINER unambiguously projects the input coordinate into another, or in other words, mapping the signal into the

one with more low-frequency components. As a result, a tiny network could achieve very high accuracy compared with previous methods.

3. Performance of INR

3.1. Background of expressive power of INR

Following the form of Yuce et al. [42], an INR with a 1D input x could be modeled as a function $f_\theta(x)$ that maps the input coordinate x to its attribute, that

$$\begin{aligned} \mathbf{z}^0 &= \gamma(x), \\ \mathbf{z}^l &= \rho^l(\mathbf{W}^l \mathbf{z}^{l-1} + \mathbf{b}^l), \quad l = 1, \dots, L-1, \\ f_\theta(x) &= \mathbf{W}^L \mathbf{z}^{L-1} + \mathbf{b}^L \end{aligned} \quad (1)$$

where $\gamma(\cdot)$ is the preprocess function which is often used to encode more frequency bases in the network, \mathbf{z}^l is the output of the l -th layer in INR, ρ is the activation function, \mathbf{W}^l and \mathbf{b}^l are the weight and bias matrix in the l -th layer, L is the number of layers in INR, $\theta = \{\mathbf{W}^l, \mathbf{b}^l\}_1^L$ refers to the set of all training parameters in the network.

This INR could represent signals with functions following the form

$$f(x) = \sum_{\omega' \in \mathcal{H}_\Omega} c_{\omega'} \sin(\langle \omega', x \rangle + \phi_{\omega'}), \quad (2)$$

where \mathcal{H}_Ω is the frequency set [42] determined by the frequency selected in the preprocess function $\gamma(\cdot)$, *e.g.*, the Fourier encoding [38], or the \sin activation [34]. In other words, *the expressive power of INR is restricted to functions that can be represented using a linear combination of certain harmonics of the $\gamma(\cdot)$ [42].*

3.2. Arrangement order of a signal determines the capacity of INR

According to the expressive power of an INR (Eqn. 2), a signal could be well learned when the encoded frequencies in the INR are consistent with the signal's frequency distribution. However, there are two problems in applying this conclusion,

1. The frequency distribution of a signal could not be known in advance, especially in inverse problems. Thus proper frequencies could not be well set in designing the architecture of an INR.
2. Due to the spectral bias in network training [28], the low-frequency components in a signal will be learned first, while the high-frequency components are learned in an extremely slow convergence [3, 10, 32, 38].

We notice that most of the signals recorded or to be inversely solved today are discrete signals. The frequency distribution of a discrete signal could be changed by arranging elements in different orders at the cost of additional storage for the arrangement rule, resulting in different satisfactions

of Eqn. 2. Consequently, the capacity of an INR for representing a signal changes with different arrangement orders.

Fig. 2 gives an intuitive demonstration. The Baboon¹ image is arranged in different orders in Fig. 2(a), and the corresponding Fourier spectrum images are shown in Fig. 2(b). The original image contains rich low-, intermediate- and high-frequency information. By sorting the Baboon according to the intensities of pixels, the high-frequency information in y -axis almost disappeared. We then arrange the pixels using a random order. Currently, the Baboon contains much high-frequency information. Then a PE+MLP (2×64 , *i.e.*, 2 hidden layers and 64 neurons per layer with ReLU activation) network is applied to learn the mapping between the coordinates and intensities of these three images (Fig. 2(c)). The fitting results differ significantly. The PE+MLP gets the best performance in the sorted image, which contains the most low-frequency information, while the worst results appear in the random sorted image, which contains the most high-frequency information. In summary:

Proposition 1. *Different arrangements of a signal have different frequency distributions, resulting in different capacities of INR for representing the signal-self.*

4. Disorder-invariant INR

4.1. Hash-mapping for INR

Given a paired discrete signal $Y = \{(\vec{x}_i, \vec{y}_i)\}_{i=1}^N$, where \vec{x}_i be the i -th d_{in} -dimensional coordinate, and \vec{y}_i be the corresponding d_{out} -dimensional signal attribute. Following the analysis mentioned above, an ideal arrangement rule $M^* : \mathbb{R}^{d_{in}} \rightarrow \mathbb{R}^{d_{in}}$ should meet the following rule,

$$M^* = \arg \min_M \sum_{k=1}^{K_M} |\omega_k| \quad (3a)$$

$$\Omega_M \subseteq \mathcal{H}_\Omega, \quad \Omega_M = \{\omega_k\}_{k=1}^{K_M}, \quad (3b)$$

where Ω_M is the set of frequency by mapping the signal following the rule M , \mathcal{H}_Ω is the supported frequency set of the INR network (Eqn. 2), K_M is the number of frequency in the arranged signal, $|\cdot|$ returns the absolute value of \cdot . The signal with this arrangement could be well learned since both the problems of improper frequency setting (Eqn. 3b) and the spectral bias (Eqn. 3a) are taken into account.

However, this strategy requires prior knowledge of the signal distribution, which is only suitable for the compression task. In contrast, it losses the ability to optimize inverse problems where the signal distribution to be optimized could not be achieved in advance. In this subsection, we detail the proposed DINER to handle this problem.

We specifically design a full-resolution hash-table \mathcal{HM} to model the mapping mentioned above. The hash-table \mathcal{HM} is set as the same length N as the number of elements

¹Baboon is a self-contained image in the Matlab by MathWorks[©].

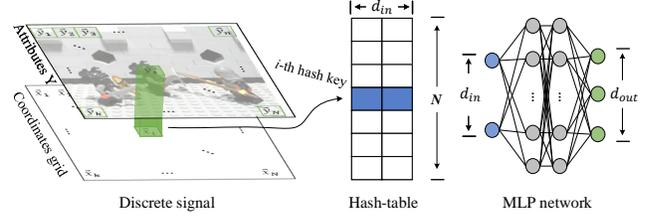


Figure 3. Pipeline of the DINER.

in Y , meanwhile the width of \mathcal{HM} is set as d_{in} (the dimensions of \vec{x}_i). Firstly, the input coordinate \vec{x}_i is used to query the i -th hash key $M(\vec{x}_i)$ in the \mathcal{HM} . Then the mapped coordinate $M(\vec{x}_i)$ is fed into a standard MLP. All parameters in the \mathcal{HM} are set as learnable, *i.e.*, parameters in \mathcal{HM} and the network parameters will be jointly optimized during the training process. Fig. 3 demonstrates the above process (The Lego Knight used comes from [1]).

Due to the hash-table, the MLP network actually learns the mapped signal. Fig. 2(d) shows the mapped pixels of the Baboon after hash-table, it is noticed that the original grid coordinates (Fig. 2(a)) are projected into irregular points (Fig. 2(d)). We sample a mesh evenly according to the minimum and the maximum values in the mapped coordinates, and feed them into the trained MLP, *i.e.*, the Fig. 2(f). For simplicity, the image in Fig. 2(f) is later called ‘learned INR’. The learned INR differs significantly from the Baboon in that the former is much smoother than the latter and has many low-frequency components (Fig. 2(g)). As a result, high accuracy of MLP fitting could be achieved using the hash-table (Fig. 2(e)).

4.2. Analysis of disorder-invariance

When applying the INR to tasks of signal representation and inverse problem optimization, the training or the optimization of the traditional INR and DINER could be modeled as Eqns. 4a and 4b, respectively.

$$\theta^* = \arg \min_{\theta} \mathcal{L} \left(\mathcal{P} \left(\{f_{\theta}(\vec{x}_i)\}_1^N \right), \mathcal{P} \left(\{\vec{y}_i\}_1^N \right) \right) \quad (4a)$$

$$\begin{aligned} \theta^*, \mathcal{HM}^* &= \arg \min_{\theta, \mathcal{HM}} \mathcal{L} \left(\mathcal{P} \left(\{f_{\theta}(\mathcal{HM}(\vec{x}_i))\}_1^N \right), \mathcal{P} \left(\{\vec{y}_i\}_1^N \right) \right) \\ &= \arg \min_{\theta, \mathcal{HM}} \mathcal{L} \left(\mathcal{P} \left(\{f_{\theta}(\mathcal{HM}_i)\}_1^N \right), \mathcal{P} \left(\{\vec{y}_i\}_1^N \right) \right), \end{aligned} \quad (4b)$$

where \mathcal{P} is a physical process and is an identical transformation for the representation task, \mathcal{L} is the loss function according to the measurements and the reconstructed results, \mathcal{HM}_i is the i -th key in the hash-table. Because the hash-index operation has no gradient, the first equation in Eqn. 4b could be simplified to the second one in Eqn. 4b.

It is noticed that paired relationship between the coordinate \vec{x}_i and value \vec{y}_i in Eqn. 4a is broken in the Eqn. 4b. There is only one independent variable \vec{y}_i in the loss function Eqn. 4b. Assuming parameters in θ are initialized with

the same values and all keys in \mathcal{HM} are set with the same one (e.g., 0) in every experiment, when applying a different order to the signal Y , e.g., $Y' = \{\vec{x}_j, \vec{y}_i\}_{j=1, i=N}^{j=N, i=1}$, the training of θ and \mathcal{HM} for signals Y and Y' share the same optimization progress in every gradient update since all parameters in Eqn. 4b for Y are equivalent to ones for Y' . As a result, the same θ^* will be optimized while the \mathcal{HM}' of Y' is also an inverse arrangement of the \mathcal{HM} of Y . This equivalence is not limited to the Y' with an inverse order; actually, it could be easily proved that the equivalence holds for Y with an arbitrary order.

Fig. 2(d), (e), and (f) illustrate this equivalence. Although the Baboon is arranged with different orders, the hash-table maps them into the same signal (Fig. 2(d), (f)), and DINER optimizes them with similar PSNR values 31.03, 31.08, 31.05². In summary:

Proposition 2. *The DINER is disorder-invariant, and signals with the same histogram distribution of attributes share an optimized network with the same parameter values.*

4.3. Discussion

Backbone network. The backbone of the proposed DINER is not limited to the standard MLP used above; actually, other network structures such as the SIREN could also be integrated with the hash-table and get better performance than the original structure. Please refer to the experimental section for more details.

Complexity. Although the number of parameters of the hash-table is much larger than the network, the training cost is very small because only one hash-key needs to be updated for training an MLP with batchsize 1. As a result, the computational complexity of training hash-table is $\mathcal{O}(1)$ in each iteration of training.

5. Experiments

To verify the performance of the proposed DINER, four separate tasks are conducted: 2D image fitting, neural representation for 3D video, phase retrieval in lensless imaging, and 3D Refractive Index recovery in intensity diffraction tomography.

5.1. 2D Image Fitting

5.1.1 Dataset and Algorithm Setup

The 2D image fitting task is adopted to test the performance of the proposed DINER and to illustrate the change of frequency distribution of the signal. We use 30 high-resolution images with 1200×1200 resolution from the SAMPLING category of the TESTIMAGES dataset [2] to

²The slight difference in values comes from the floating point errors of GPU for summing matrix with same histogram and different arrangement orders.

Freq. bands				
Original Image	0.4426	0.2484	0.1753	0.1337
Learned INR (MLP)	0.6784	0.1218	0.0973	0.1025
Learned INR (SIREN)	0.6354	0.1220	0.1149	0.1276

Table 1. Comparisons of the ratio of frequency distributions between the original image and the learned INR. Two backbones, i.e., the MLP and SIREN, are both compared.

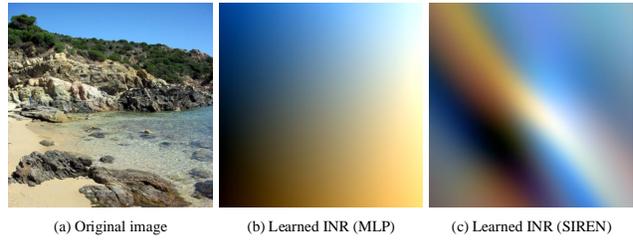


Figure 4. Comparisons of learned INRs in the DINER with MLP and SIREN backbones, respectively.

evaluate the performance of various algorithms. Each image is generated using custom Octave/MATLAB software scripts specifically written to guarantee the precise positioning and value of every pixel and contains rich low-, intermediate- and high-frequency information. In the following experiments, all our results are obtained with the same network configuration 2×64 unless otherwise stated.

5.1.2 Comparisons of frequency distribution with and without hash-table

As noticed in the Sec. 4.1 and the Figs. 2, the hash-table maps the original signal with more low-frequency contents. Tab. 1 provides statistics of the mean frequency distributions of the original image and the learned INR over 30 images. The ratios of the intermediate- and high-frequency information are all reduced after mapping, while the low-frequency information are increased.

Compared with the supported frequency set $\mathcal{H}_\Omega^{SIREN}$ of the SIREN structure where a default frequency 30 is used in activation, the \mathcal{H}_Ω^{MLP} of the MLP structure without any frequency encoding contains more low frequencies and less high frequencies. Accordingly, there are more low-frequency information in the mapped image of DINER with MLP backbone than the one with SIREN backbone (0.6784 vs 0.6354, Fig. 4). Please refer to the supplementary material for more qualitative comparisons.



Figure 5. Qualitative comparisons of various methods on 2D image fitting after 3000 epochs.

5.1.3 Comparisons with the State-of-the-arts

We compare the proposed DINER with the Fourier feature positional encoding (PE+MLP) [38], SIREN [34] and InstantNGP [24]. Noting that, two backbones, *i.e.*, the standard MLP with ReLU activation and the SIREN with periodic function activation, are all combined with the proposed hash-table to better evaluate the performance. We control the size of the hash-table used in the InstantNGP to guarantee the similar parameters with ours, *e.g.*, 2^{21} in the 2D image fitting task while ours has a length of $1200^2 < 2^{21}$. Apart from this, all 5 methods are trained with the same L_2 loss between the predicted value and the ground truth, and other parameters are set with the default values by authors.

Fig. 1 shows the PSNR of various methods at different epochs. It is noticed that the SIREN and PE+MLP convergences quickly at the early stage and reaches about 24dB and 21dB finally. On the contrary, the proposed two methods both provide higher accuracy than backbones. The PSNRs of two backbones for image fitting are increased 14dB and 16dB using the hash-table, respectively. Additionally, although the InstantNGP converges very fast to about 30dB at about 200 epochs, the curve tends to be stable at the last 2800 epochs. The proposed DINER with MLP backbone achieves an advantage of 5dB than the InstantNGP.

Fig. 5 shows the qualitative results at 3000 epochs. The proposed methods outperform the SIREN and PE+MLP. Our methods provide more clear details especially in high-frequency boundaries, such as the bell tower (yellow box) in Fig. 5. The fitted image of the InstantNGP is very similar to the GT at first sight, however many noises appear in the zoom-in results, for example there are a lot of noisy points in the wall (red box) and the sky (yellow box) of Fig. 5, results in a lower PSNR metric.

Tab. 2 lists the training time of 5 methods. The InstantNGP is implemented with the tiny-cuda-nn [23], while other 4 methods are implemented with the Pytorch. All 5 methods are trained on a NVIDIA A100 40GB GPU. The optimization of hash-table requires additional 3 seconds on the

	DINER (SIREN)	DINER (MLP)	SIREN	InstantNGP	PE +MLP
Time	80.5s	58.5s	77.6s	38s	78.8s

Table 2. Comparisons of training a 2D image with 3000 epochs.

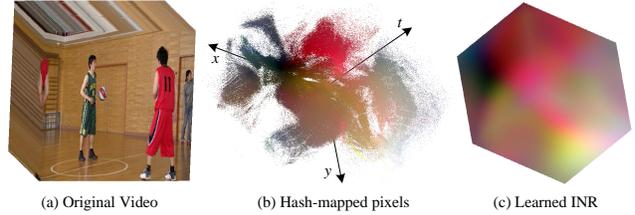


Figure 6. Visualization of learned INR on the 3D video ‘BasketballPass’ [36]. (a) and (b) compare the coordinates with and without hash-table. (c) shows the learned INR of the SIREN after the mapping by hash-table.

SIREN architecture and reduces 20 seconds compared with the classical PE+MLP architecture, verifying the low complexity of optimizing hash-table.

5.2. Neural Representation for 3D video

Video describes a dynamic 3D scene $I(t, x, y)$ composed of multiple frames. Accurate representation for 3D video is becoming a popular task [5, 14, 16, 31] in the community of INR. We compare the proposed DINER with the SIREN [34] and the state-of-the-art INR NeRV [5]. Noting only the Hash+SIREN architecture is evaluated in this task. The NeRV is implemented using their default parameters, while the SIREN and the proposed Hash+SIREN both use the same network structure with the size 4×64 . All three methods are evaluated on the videos ‘ReadySetGo’ and ‘ShakeNDry’ of the UVG dataset [21] and are trained with 500 epochs. The first 30 frames with 1920×1080 resolution are used in our experiment.

Fig. 6(a) and (b) show the mapping of the coordinates with and without mapping. Fig. 6(c) illustrates the learned INR. It is noticed that the low-frequency property also ap-



Figure 7. Qualitative comparisons of various methods on 3D video representation after 500 epochs.

Methods	Network Para.	Training time	PSNR
SIREN	8.77K	706s	21.22 dB
NeRV	97.24M	7445s	36.08 dB
DINNER	8.77K	1309s	50.74 dB

Table 3. Comparisons of training 3D videos 'ReadySetGo' and 'ShakeNDry' with 500 epochs.

pears in the learned INR of the 3D video in our DINER. Tab. 3 shows the quantitative comparisons. The proposed method outperforms the NeRV both in quality and speed with 14 dB and $5\times$ improvements, respectively. Because a large hash-table is used in DINER, more time are taken in the transmission between the memory and the cache in the GPU, resulting more training time of DINER than the SIREN. Fig. 7 shows the qualitative comparisons. Noting that the SIREN with a tiny network could not provide reasonable representation for the 'ReadySetGo' data with $30\times 1920\times 1080$ pixels, thus all pixels are smoothed. NeRV provides better results than the ones by SIREN, however the high-frequency details are lost such as the character 'U' (left-top corner) and the red logo of horsehead (right-top corner) in the purple box, as well as the folds of the trousers in the orange box. On the contrary, the original video is mapped with little high-frequency component in the proposed DINER (Fig. 6(c)). As a result, the details mentioned above could be well represented.

5.3. Phase Recovery in Lensless Imaging

Lensless imaging [25] observes specimen in a very close distance without any optical lens. By directly recording the diffractive measurements, it provides the advantage of wide field of view observation and has become an attractive microscopic technique [43] for analyzing the properties of the specimen. We take the classic multi-height lensless imaging as an example, where N measurements $\{I_z\}_{z=z_1}^{z_N}$ are captured under different specimen-to-sensor distances z for the specimen's amplitude and phase imaging recovery. In multi-height lensless imaging, I_z could be modeled as ap-

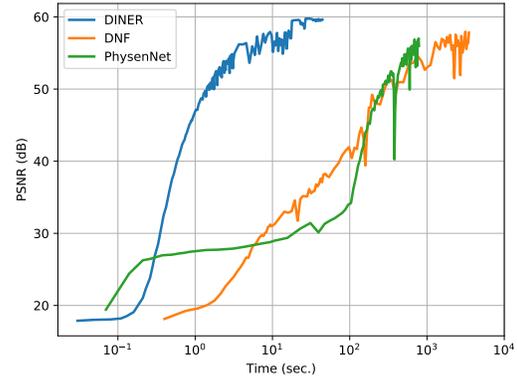


Figure 8. PSNR of reconstructed measurements over training time on lensless imaging.

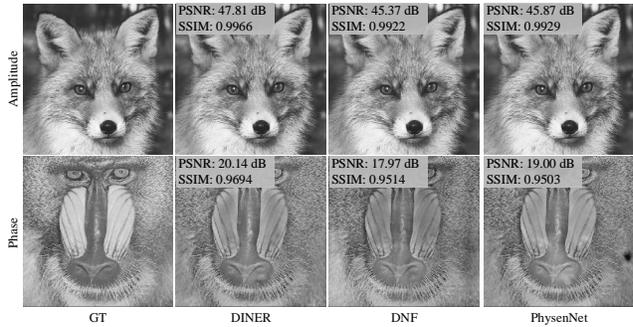


Figure 9. Comparisons on lensless imaging.

plying Fresnel propagation to the complex field $O(x, y)$ of the specimen, *i.e.*,

$$I_z = |PSF_z * [P(x, y) \cdot O(x, y)]|^2, \quad (5)$$

where $P(x, y)$ is the illumination pattern, PSF_z is the point spread function of Fresnel propagation over distance z between the specimen and the sensor.

We model $O(x, y)$ using the proposed method with the SIREN backbone and the network size is 2×64 . The loss function is built by comparing the measurements with the results from applying the Eqn. 5 to the network output. We compare our method with the current SOTAs, *i.e.*, the

diffractive neural field (DNF) [44] and the PhysenNet [40]. Due to the lack of the ground truth in real images, we only provide comparisons on the open-source synthetic data [44] here, please refer the supplementary material for more qualitative comparisons on real data.

Fig. 8 shows the PSNR of reconstructed measurements over training time. The proposed method has $18\times$ and $80\times$ advantages on the convergence speed over the PhysenNet and DNF, respectively. Fig. 9 provides qualitative comparisons. Although only 1.5% network parameters (2×64 vs. 8×256) are used, the proposed DINER could provide better results than the DNF thanks to the hash-table.

5.4. 3D Refractive Index Recovery in Intensity Diffraction Tomography

The 3D refractive index characterizes the interaction between light and matter within a specimen. It is an endogenous source of optical contrast for imaging specimen without staining or labelling, and plays an important role in many areas, *e.g.* the morphogenesis, cellular pathophysiology, biochemistry [27]. Intensity diffraction tomography measures the squared amplitude of the light scattered from the specimen at different angles multiple times and has become a popular technique for recovering the 3D refractive index.

Given the 3D refractive index $\mathbf{n} = (\mathbf{n}_{re} + j\mathbf{n}_{im})$ of a specimen, where \mathbf{n}_{re} and \mathbf{n}_{im} are the real and imaginary parts of the specimen’s refractive index, respectively. The forward imaging process of sensor placed at location ρ could be modeled as

$$I_\rho = \mathbf{A}_\rho \Delta\epsilon, \tag{6}$$

where \mathbf{A}_ρ records the sample-intensity mapping with the illuminations. $\Delta\epsilon = \Delta\epsilon_{re} + j\Delta\epsilon_{im}$ is the complex-valued permittivity contrast and could be obtained by solving

$$\begin{aligned} \mathbf{n}_{re} &= \sqrt{\frac{1}{2} \left((\mathbf{n}_0^2 + \Delta\epsilon_{re}) + \sqrt{(\mathbf{n}_0^2 + \Delta\epsilon_{re})^2 + \Delta\epsilon_{im}^2} \right)}, \\ \mathbf{n}_{im} &= \frac{\Delta\epsilon_{im}}{2 \cdot \mathbf{n}_{re}} \end{aligned} \tag{7}$$

where \mathbf{n}_0 is the refractive index of the background medium.

We model the $(\Delta\epsilon_{re}, \Delta\epsilon_{im})$ using the DINER with network size 2×64 . We compare our method with the SOTA, *i.e.*, DeCAF [19] which uses a combination of the standard MLP structure with network size 10×208 and positional+radial encodings.

Fig. 10 compares our method with the DeCAF on the 3D Granulocyte Phantom data using the Fiji software [33]. Because the ground truth is not open-source until the submission, we take a screenshot of the ground truth in their paper as a reference. The MSE values labelled in Fig. 10 are computed by comparing the reconstructed measurements using the network output and the real measurements. Since the hash-table could map a high-frequency signal in a low-frequency way, our results provide more details on the sur-

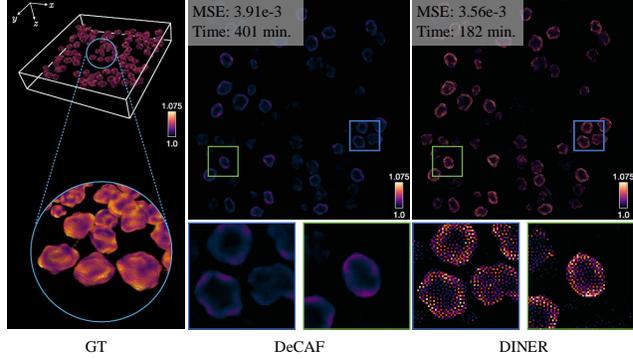


Figure 10. Comparisons on 3D refractive index recovery. DINER takes less training time and could reconstruct more surface details of the Granulocyte.

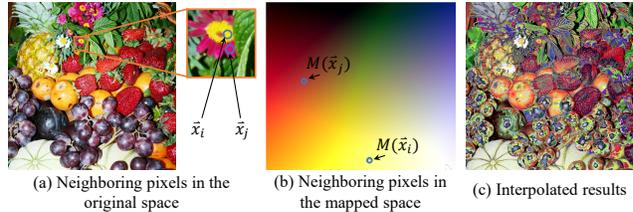


Figure 11. Analysis of feeding interpolated hash-key to the DINER. (a) Two neighboring coordinates \vec{x}_i and \vec{x}_j are labelled in the original image. (b) The distance between the mapped coordinates $\mathcal{H}\mathcal{M}(\vec{x}_i)$ and $\mathcal{H}\mathcal{M}(\vec{x}_j)$ is larger than the one in the original space. (c) Results by feeding the interpolated mapped coordinates to the trained MLP.

face of the Granulocyte. While the surface of the Granulocyte is over-smoothed in the results of the DeCAF since the PE+MLP could not accurately model the high-frequency components (Please refer to the supplemented video on different heights for better comparison).

5.5. Discussion

The aforementioned experiments are all focused on discrete signals. To query an unseen coordinate in a continuous signal, it is suggested to apply a post-interpolation operation to the network output instead of feeding interpolated hash-key to the network (see Fig. 11), such as the exploration of Plenoxels which interpolates the grid of density and harmonic coefficients [8] instead of feeding unseen position and direction coordinates to the network directly [22].

6. Conclusions

In this work, we have proposed the DINER which could greatly improve the accuracy of current INR backbones by introducing an additional hash-table. We have pointed out that the performance of INR for representing a signal is determined by the arrangement order of elements in it. The

proposed DINER could map the input discrete signal into a low-frequency one, which is invariant if only the arrangement order changes while the histogram of attributes is not changed. For this reason, the accuracy of different INR backbones could be greatly improved. Extensive experiments have verified the high accuracy and efficiency of the proposed DINER for tasks of signal fitting and inverse problem optimization.

However, the current DINER could only process discrete signals. In the future, we will focus on continuous mapping methods instead of discrete hash-table-based mapping to extend the advantages for continuous signals such as the signed distance function [26].

References

- [1] The new stanford light field archive. <http://lightfield.stanford.edu/lfs.html>. 4
- [2] Nicola Asuni and Andrea Giachetti. TESTIMAGES: a large-scale archive for testing visual devices and basic image processing algorithms. In *STAG*, pages 63–70, 2014. 5
- [3] Alberto Bietti and Julien Mairal. On the inductive bias of neural tangent kernels. *Advances in Neural Information Processing Systems*, 32, 2019. 3
- [4] Rohan Chabra, Jan E Lenssen, Eddy Ilg, Tanner Schmidt, Julian Straub, Steven Lovegrove, and Richard Newcombe. Deep local shapes: Learning local sdf priors for detailed 3d reconstruction. In *European Conference on Computer Vision*, pages 608–625. Springer, 2020. 3
- [5] Hao Chen, Bo He, Hanyu Wang, Yixuan Ren, Ser Nam Lim, and Abhinav Shrivastava. NeRV: Neural representations for videos. *Advances in Neural Information Processing Systems*, 34:21557–21568, 2021. 6
- [6] Yuyao Chen, Lu Lu, George Em Karniadakis, and Luca Dal Negro. Physics-informed neural networks for inverse problems in nano-optics and metamaterials. *Optics Express*, 28(8):11618–11633, 2020. 2
- [7] Rizal Fathony, Anit Kumar Sahu, Devin Willmott, and J Zico Kolter. Multiplicative filter networks. In *International Conference on Learning Representations*, 2020. 1, 3
- [8] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinzhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5501–5510, 2022. 8
- [9] Ruohan Gao, Zilin Si, Yen-Yu Chang, Samuel Clarke, Jeanette Bohg, Li Fei-Fei, Wenzhen Yuan, and Jiajun Wu. Objectfolder 2.0: A multisensory object dataset for sim2real transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10598–10608, 2022. 2
- [10] Reinhard Heckel and Mahdi Soltanolkotabi. Compressive sensing with un-trained neural networks: Gradient descent finds a smooth approximation. In *International Conference on Machine Learning*, pages 4149–4158. PMLR, 2020. 3
- [11] Chiyu Jiang, Avneesh Sud, Ameesh Makadia, Jingwei Huang, Matthias Nießner, Thomas Funkhouser, et al. Local implicit grid representations for 3d scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6001–6010, 2020. 3
- [12] George Em Karniadakis, Ioannis G Kevrekidis, Lu Lu, Paris Perdikaris, Sifan Wang, and Liu Yang. Physics-informed machine learning. *Nature Reviews Physics*, 3(6):422–440, 2021. 1, 2
- [13] Petr Kellnhofer, Lars C Jebe, Andrew Jones, Ryan Spicer, Kari Pulli, and Gordon Wetzstein. Neural lumigraph rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4287–4297, 2021. 1
- [14] Subin Kim, Sihyun Yu, Jaeho Lee, and Jinwoo Shin. Scalable neural video representations with learnable positional features. *arXiv preprint arXiv:2210.06823*, 2022. 6
- [15] Moshe Leshno, Vladimir Ya Lin, Allan Pinkus, and Shimon Schocken. Multilayer feedforward networks with a nonpolynomial activation function can approximate any function. *Neural Networks*, 6(6):861–867, 1993. 2
- [16] Zizhang Li, Mengmeng Wang, Huaijin Pi, Kechun Xu, Jianbiao Mei, and Yong Liu. E-NeRV: Expedite neural video representation with disentangled spatial-temporal context. In *European Conference on Computer Vision*, pages 267–284. Springer, 2022. 6
- [17] David B Lindell, Dave Van Veen, Jeong Joon Park, and Gordon Wetzstein. BACON: Band-limited coordinate networks for multiscale scene representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021. 1, 3
- [18] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *Advances in Neural Information Processing Systems*, 33:15651–15663, 2020. 3
- [19] Renhao Liu, Yu Sun, Jiabei Zhu, Lei Tian, and Ulugbek S Kamilov. Recovery of continuous 3d refractive index maps from discrete intensity-only measurements using neural fields. *Nature Machine Intelligence*, 4(9):781–791, 2022. 1, 2, 8
- [20] Julien NP Martel, David B Lindell, Connor Z Lin, Eric R Chan, Marco Monteiro, and Gordon Wetzstein. Acorn: adaptive coordinate networks for neural scene representation. *ACM Transactions on Graphics (TOG)*, 40(4):1–13, 2021. 1, 3
- [21] Alexandre Mercat, Marko Viitanen, and Jarno Vanne. UVG dataset: 50/120fps 4k sequences for video codec analysis and development. In *Proceedings of the 11th ACM Multimedia Systems Conference*, pages 297–302, 2020. 6
- [22] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. NeRF: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*, pages 405–421. Springer, 2020. 1, 2, 3, 8
- [23] Thomas Müller. Tiny cuda neural network framework. <https://github.com/nvmlabs/tiny-cuda-nn>, 2021. 6
- [24] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a mul-

- tiresolution hash encoding. *ACM Transactions on Graphics (TOG)*, 41(4):102:1–102:15, 2022. 3, 6
- [25] Aydogan Ozcan and Euan McLeod. Lensless imaging and sensing. *Annual Review of Biomedical Engineering*, 18:77–102, 2016. 7
- [26] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 165–174, 2019. 2, 9
- [27] YongKeun Park, Christian Depeursinge, and Gabriel Popescu. Quantitative phase imaging in biomedicine. *Nature Photonics*, 12(10):578–589, 2018. 8
- [28] Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In *International Conference on Machine Learning*, pages 5301–5310. PMLR, 2019. 1, 2, 3
- [29] Maziar Raissi, Alireza Yazdani, and George Em Karniadakis. Hidden fluid mechanics: Learning velocity and pressure fields from flow visualizations. *Science*, 367(6481):1026–1030, 2020. 2
- [30] Brandon Reyes, Amanda A Howard, Paris Perdikaris, and Alexandre M Tartakovsky. Learning unknown physics of non-newtonian fluids. *Physical Review Fluids*, 6(7):073301, 2021. 2
- [31] Daniel Rho, Junwoo Cho, Jong Hwan Ko, and Eunbyung Park. Neural residual flow fields for efficient video representations. *arXiv preprint arXiv:2201.04329*, 2022. 6
- [32] Basri Ronen, David Jacobs, Yoni Kasten, and Shira Kritchman. The convergence rate of neural networks for learned functions of different frequencies. *Advances in Neural Information Processing Systems*, 32, 2019. 3
- [33] Johannes Schindelin, Ignacio Arganda-Carreras, Erwin Frise, Verena Kaynig, Mark Longair, Tobias Pietzsch, Stephan Preibisch, Curtis Rueden, Stephan Saalfeld, Benjamin Schmid, et al. Fiji: an open-source platform for biological-image analysis. *Nature Methods*, 9(7):676–682, 2012. 8
- [34] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. *Advances in Neural Information Processing Systems*, 33:7462–7473, 2020. 1, 3, 6
- [35] Vincent Sitzmann, Semon Rezchikov, Bill Freeman, Josh Tenenbaum, and Fredo Durand. Light field networks: Neural scene representations with single-evaluation rendering. *Advances in Neural Information Processing Systems*, 34:19313–19325, 2021. 2
- [36] Gary J Sullivan, Jens-Rainer Ohm, Woo-Jin Han, and Thomas Wiegand. Overview of the high efficiency video coding (hevc) standard. *IEEE Transactions on circuits and systems for video technology*, 22(12):1649–1668, 2012. 6
- [37] Towaki Takikawa, Joey Litalien, Kangxue Yin, Karsten Kreis, Charles Loop, Derek Nowrouzezahrai, Alec Jacobson, Morgan McGuire, and Sanja Fidler. Neural geometric level of detail: Real-time rendering with implicit 3d shapes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11358–11367, 2021. 3
- [38] Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *Advances in Neural Information Processing Systems*, 33:7537–7547, 2020. 1, 2, 3, 6
- [39] Ayush Tewari, Justus Thies, Ben Mildenhall, Pratul Srinivasan, Edgar Tretschk, W Yifan, Christoph Lassner, Vincent Sitzmann, Ricardo Martin-Brualla, Stephen Lombardi, et al. Advances in neural rendering. In *Computer Graphics Forum*, volume 41, pages 703–735. Wiley Online Library, 2022. 1, 2
- [40] Fei Wang, Yaoming Bian, Haichao Wang, Meng Lyu, Giancarlo Pedrini, Wolfgang Osten, George Barbastathis, and Guohai Situ. Phase imaging with an untrained neural network. *Light: Science & Applications*, 9(1):1–7, 2020. 8
- [41] Pak-Hei Yeung, Linde Hesse, Moska Aliasi, Monique Haak, Weidi Xie, Ana IL Namburete, et al. ImplicitVol: Sensorless 3d ultrasound reconstruction with deep implicit representation. *arXiv preprint arXiv:2109.12108*, 2021. 2
- [42] Gizem Yüce, Guillermo Ortiz-Jiménez, Beril Besbinar, and Pascal Frossard. A structured dictionary perspective on implicit neural representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19228–19238, 2022. 1, 3
- [43] You Zhou, Xia Hua, Zibang Zhang, Xuemei Hu, Krishna Dixit, Jingang Zhong, Guoan Zheng, and Xun Cao. Wirtinger gradient descent optimization for reducing gaussian noise in lensless microscopy. *Optics and Lasers in Engineering*, 134:106131, 2020. 7
- [44] Hao Zhu, Zhen Liu, You Zhou, Zhan Ma, and Xun Cao. DNF: Diffractive neural field for lensless microscopic imaging. *Optics Express*, 30(11):18168–18178, 2022. 1, 2, 8