# From MLP to NeoMLP: Leveraging Self-Attention for Neural Fields

**Anonymous Author(s)** 

Affiliation Address email

# **Abstract**

Neural fields (NeFs) have recently emerged as a state-of-the-art method for encoding spatio-temporal signals of various modalities. Despite the success of NeFs in reconstructing individual signals, their use as representations in downstream tasks, such as classification or segmentation, is hindered by the complexity of the parameter space and its underlying symmetries, in addition to the lack of powerful and scalable conditioning mechanisms. In this work, we draw inspiration from the principles of connectionism to design a new architecture based on MLPs, which we term NeoMLP. We start from an MLP, viewed as a graph, and transform it from a multi-partite graph to a *complete graph* of input, hidden, and output nodes, equipped with high-dimensional features. We perform message passing on this graph and employ weight-sharing via self-attention among all the nodes. NeoMLP has a built-in mechanism for conditioning through the hidden and output nodes, which function as a set of latent codes, and as such, NeoMLP can be used straightforwardly as a conditional neural field. We demonstrate the effectiveness of our method by fitting high-resolution signals, including multi-modal audio-visual data. Furthermore, we fit datasets of neural representations, by learning instance-specific sets of latent codes using a single backbone architecture, and then use them for downstream tasks, outperforming recent state-of-the-art methods.

# 1 Introduction

2

4

5

6

10

11

12

13

14

15

16

17

18

19

The omnipresence of neural networks in the last decade has recently given rise to neural fields (NeFs) (cf. Xie et al. [46]) as a powerful and scalable method to encode continuous signals of various 21 modalities. These range from shapes [29], scenes [24], and images, [40], to physical fields [17], CT 22 scans [27, 8], and partial differential equations [48, 16]. Consequently, the popularity of NeFs has 23 spurred interest in neural representations, i.e. using NeFs as representations for downstream tasks. 24 Existing neural representations, however, suffer from notable drawbacks. Representations based on 25 unconditional neural fields, i.e. independent multi-layer perceptrons (MLPs) fitted on each signal, 26 are subject to parameter symmetries [11], which lead to extremely poor performance in downstream 27 tasks if left unattended [25]. Many recent works [25, 55, 18, 21, 28] have proposed architectures that 28 respect the underlying symmetries; the performance, however, leaves much to be desired. Another 29 line of works [29, 9] has proposed conditional neural fields with a single latent code per signal that 30 modulates the activations of a shared MLP through concatenation, FiLM [30], or hypernetworks [10], 31 while, recently, other works [35, 45] have proposed set-latent conditional neural fields—conditional 32 neural fields with a set of latent codes—that condition the signal through attention [44]. Whilst 33 the study of Rebain et al. [33] showed that set-latent neural fields outperform single latent code methods as conditioning mechanisms, existing set-latent neural fields are based on cross-attention, which limits their scalability and expressivity: coordinates are only used as queries in attention, and 36 cross-attention is limited to a single layer.

We argue that many of these drawbacks stem from the lack of a unified native architecture that integrates the necessary properties of neural representations and eliminates the shortcomings of current approaches. To address these concerns, we draw inspiration from *connectionism* and the long history of MLPs to design a new architecture that functions as a standard machine learning model—akin to an MLP—as well as a conditional neural field. The paradigm of neural networks, from the early days of Perceptron [23], to MLPs with hidden neurons trained with backpropagation [34], to modern transformers [44], shares the connectionist principle: cognitive processes can be described by interconnected networks of simple and often uniform units.

This principle is lacking from current conditional neural field architectures, since conditioning is 46 added to the network as an ad-hoc mechanism. In contrast, motivated by this principle, we take a 47 closer look at MLPs; more specifically, we look at MLPs as a graph—similar to a few recent works 48 [18, 21, 26]— and design a novel architecture that operates on this graph using message passing. First, 49 we convert the graph from a multi-partite graph to a fully-connected graph with self-edges. Instead of 50 using edge-specific weights, we employ weight-sharing via self-attention among all the nodes. We 51 initialize the hidden and output nodes with noise and optimize their values with backpropagation. Finally, we use high-dimensional features for all nodes to make self-attention and the network as a 53 whole more scalable. 54

We make the following contributions. First, we propose a new architecture, which we term *Neo*MLP, 55 by viewing MLPs as a graph, and convert this graph to a complete graph of input, hidden, and 56 output nodes with high-dimensional features. We employ message passing on that graph through 57 self-attention among the input, hidden, and output nodes. The hidden and output nodes can be used as a learnable set of latent codes, and thus, our method can function as a conditional neural field. 59 We introduce new neural representations that use sets of latent codes for each signal, which we 60 term  $\nu$ -reps, as well as datasets of neural representations, which we term  $\nu$ -sets. We fit datasets of 61 signals using a single backbone architecture, and then use the latent codes for downstream tasks, 62 outperforming recent state-of-the-art methods. We also demonstrate the effectiveness of our method 63 by fitting high-resolution audio and video signals, as well as multi-modal audio-visual data. 64

# 1.1 Background on Neural Fields

Neural fields (NeFs), often referred to as Implicit Neural Representations (INRs), are a class of neural networks that parameterize fields using neural networks (cf. Xie et al. [46]). In their simplest form, they are MLPs that take as input a single coordinate (e.g. an x-y coordinate) and output the field value for that coordinate (e.g. an RGB value). By feeding batches of coordinates to the network, and training to reconstruct the target values with backpropagation, the neural field learns to encode the target signal, without being bound to a specific resolution.

Conditional neural fields introduce a conditioning mechanism to neural fields through latent variables, often referred to as *latent codes*. This conditioning mechanism can be used to encode instance-specific information (*e.g.* encode a single image) and disentangle it from the backbone architecture, which now carries dataset-wide information.

# 76 2 NeoMLP

65

77

# 2.1 From MLP to NeoMLP

78 We begin the exposition of our method with MLPs, since our architecture is influenced by MLPs and builds on them. Without loss of generality, a multi-layer perceptron takes as input a set of scalar 79 variables  $\{x_i\}_{i=1}^I, x_i \in \mathbb{R}$ , coalesced into a single high-dimensional array  $\mathbf{x} \in \mathbb{R}^I$ . Through a series 80 of non-linear transformations, the input array is progressively transformed into intermediate (hidden) 81 representations, with the final transformation leading to the output array  $\mathbf{y} \in \mathbb{R}^{O}$ . 82 Akin to other recent works [18, 22, 26], we look at an MLP as a graph; an MLP is an L+1-partite graph, where L is the number of layers. The nodes represent the input, hidden, and output neurons, 84 and have scalar features that correspond to individual inputs, the hidden features at each layer, and 85 the individual outputs, respectively. We perform message passing on that graph, after making it 86 more amenable for learning. First, we convert the connectivity graph from an L+1-partite graph 87 to a fully-connected graph with self-edges. Since the forward pass now includes message passing from all nodes to all nodes at each step, we create learnable parameters for the initial values of the

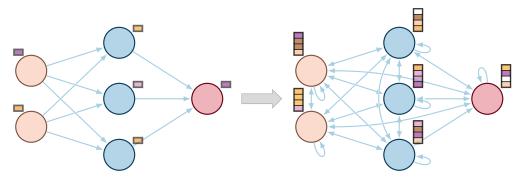


Figure 1: The connectivity graphs of MLP and *Neo*MLP. *Neo*MLP performs message passing on the MLP graph. Going from MLP to *Neo*MLP, we use a fully connected graph and high-dimensional node features. In *Neo*MLP, the traditional notion of layers of neurons, as well as the asynchronous layer-wise propagation, cease to exist. Instead, we use synchronous message passing with weight-sharing via self-attention among all the nodes. *Neo*MLP has three types of nodes: input, hidden, and output nodes. The input is fed to *Neo*MLP through the input nodes, while the output nodes capture the output of the network.

hidden and output node features. We initialize them with Gaussian noise, and optimize their values with backpropagation, simultaneously with the network parameters. Next, we observe that having dedicated edge-specific weights for all node pairs would result in an intractable spatial complexity. As such, in order to reduce the memory footprint, we follow the standard practice of graph neural networks and Transformers [44], and employ weight-sharing between the nodes, specifically via self-attention. In other words, the weights for each node pair are computed as a function of the incoming and outgoing node features, in conjunction with weights that are shared across nodes. As a by-product of the self-attention mechanism, which is permutation invariant, we use node-specific embeddings that allow us to differentiate between different nodes. Finally, instead of having scalar node features, we increase the dimensionality of node features, which makes self-attention more scalable and expressive.

We show the connectivity graph of *Neo*MLP and its conversion from a standard MLP in Figure 1. We also show the equations of the forward pass for a single layer of an MLP and a simplified version of *Neo*MLP (without softmax normalization, scaling, or multi-head attention) in Equation (1).

$$\begin{aligned} \text{MLP:} \quad \mathbf{h}_i^{(l)} &= \sum_j & \mathbf{W}_{ij}^{(l)} & \mathbf{h}_j^{(l-1)} \\ \text{NeoMLP:} \quad \mathbf{h}_i^{(l)} &= \sum_j & \left(\mathbf{W}_Q^{(l)} \mathbf{h}_i^{(l-1)}\right)^\top \mathbf{W}_K^{(l)} \mathbf{h}_j^{(l-1)} \mathbf{W}_V^{(l)} & \mathbf{h}_j^{(l-1)} \end{aligned}$$

We note that throughout this work, we retain the nomenclature of input, hidden, and output nodes, but repurpose them for *Neo*MLP. More specifically, these nodes refer to the connectivity graph of *Neo*MLP, *i.e.* the graph on which we perform message passing, shown in Figure 1, and not its computational graph, which would include layers of all the nodes. The input is fed to *Neo*MLP through the input nodes before any information propagation, while the output nodes are the ones that will capture the output of the network, after a number of message passing layers. Every other node that is not used for input or output is a hidden node. The number of hidden nodes in *Neo*MLP does not need to correspond one-to-one to the MLP hidden nodes.

# 2.2 NeoMLP Architecture

After establishing the connection with MLPs, we now discuss the architecture of our method in detail. The inputs comprise a set of scalar variables  $\{x_i\}_{i=1}^{I}, x_i \in \mathbb{R}$ . We employ random Fourier features [42] as a non-learnable method to project each scalar input (each dimension separately) to a high-dimensional space  $\mathbb{R}^{D_{\text{RFF}}}$ . This is followed by a linear layer that projects it to  $\mathbb{R}^{D}$ . We then add learnable positional embeddings to the inputs. These embeddings are required for the model to differentiate between input variables, since self-attention is a permutation invariant operation. We use similar learnable embeddings for each scalar output dimension (referred to as output embeddings), as well as H learnable embeddings for each hidden node (referred to as hidden embeddings), where H

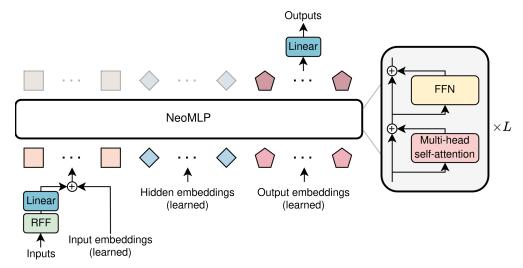


Figure 2: The architecture of NeoMLP. We pass each input dimension through an RFF layer followed by a linear layer, and then add individual input embeddings to each input. The transformed inputs, alongside the embeddings for the hidden and output nodes, comprise the inputs to NeoMLP. NeoMLP has L layers of residual self-attention and non-linear transformations. We capture the output that corresponds to the output nodes and pass it through a linear layer to get the final output of the network.

is chosen as a hyperparameter. We concatenate the transformed inputs with the hidden and output embeddings along the node (token) dimension, before feeding them to NeoMLP. We denote the concatenated tokens as  $\mathbf{T}^{(0)} \in \mathbb{R}^{(I+H+O)\times D}$ , where O is the number of output dimensions. The input, hidden, and output embeddings are initialized with Gaussian noise. We use a variance  $\sigma_i^2$  for the input embeddings and  $\sigma_o^2$  for the hidden and output embeddings; both are chosen as hyperparameters. Each NeoMLP layer comprises a multi-head self-attention layer among the tokens, and a feed-forward network that non-linearly transforms each token independently. The output of each layer consists

Each *Neo*MLP layer comprises a multi-head self-attention layer among the tokens, and a feed-forward network that non-linearly transforms each token independently. The output of each layer consists of the transformed tokens  $\mathbf{T}^{(l)} \in \mathbb{R}^{(I+H+O)\times D}$ . We use pre-LN transformer blocks [47], but omit LayerNorm [1], since we observed it does not lead to better performance or faster convergence. This also makes our method conceptually simpler. Thus, a *Neo*MLP layer is defined as follows:

$$\widetilde{\mathbf{T}}^{(l)} = \mathbf{T}^{(l-1)} + \text{SelfAttention}\left(\mathbf{T}^{(l-1)}\right)$$
 (2)

$$\mathbf{T}^{(l)} = \widetilde{\mathbf{T}}^{(l)} + \text{FeedForwardNetwork}\left(\widetilde{\mathbf{T}}^{(l)}\right)$$
(3)

We explore different variants of self-attention and find that linear attention [14, 39] performs slightly better and results in a faster model, while simultaneously requiring fewer parameters. Specifically, we use the version of Shen et al. [39] from a publicly available implementation of linear attention<sup>1</sup>.

After L NeoMLP layers, we only keep the final tokens that correspond to the output embeddings, and pass them through a linear layer that projects them back to scalars. We then concatenate all outputs together, which gives us the final output array  $y \in \mathbb{R}^O$ . The full pipeline of our method is shown in Figure 2, while the forward pass is mathematically described as follows:

$$\mathbf{i}_i = \text{Linear}(\text{RFF}(x_i)) + \text{InputEmbedding}(i), \quad i \in \{1, \dots, I\}, \, \mathbf{i}_i \in \mathbb{R}^D$$
 (4)

$$\mathbf{h}_{i} = \text{HiddenEmbedding}(j), \qquad j \in \{1, \dots, H\}, \, \mathbf{h}_{i} \in \mathbb{R}^{D}$$
 (5)

$$\mathbf{o}_k = \text{OutputEmbedding}(k), \qquad k \in \{1, \dots, O\}, \mathbf{o}_k \in \mathbb{R}^{O \times D}$$
 (6)

$$\mathbf{T}^{(0)} = \left[ \left\{ \mathbf{i}_{i} \right\}_{i=1}^{I}, \left\{ \mathbf{h}_{j} \right\}_{j=1}^{H}, \left\{ \mathbf{o}_{k} \right\}_{k=1}^{O} \right], \qquad \mathbf{T}^{(0)} \in \mathbb{R}^{(I+H+O) \times D}$$
 (7)

$$\mathbf{T}^{(l)} = \text{NeoMLPLayer}\left(\mathbf{T}^{(l-1)}\right), \qquad l \in \{1, \dots, L\}, \mathbf{T}^{(l)} \in \mathbb{R}^{(I+H+O) \times D} \quad (8)$$

$$\mathbf{y} = \operatorname{Linear}\left(\mathbf{T}_{I+H:I+H+O}^{(L)}\right), \qquad \qquad \mathbf{y} \in \mathbb{R}^{O \times 1}$$
 (9)

 $<sup>^{1}</sup> https://github.com/lucidrains/linear-attention-transformer \\$ 

#### 2.3 NeoMLP as an auto-decoding conditional neural field

One of the advantages of our method is its adaptability, since it has a built-in mechanism for conditioning, through the hidden and output embeddings. In the context of neural fields, this mechanism enables our method to function as an auto-decoding conditional neural field [29], while the embeddings can be used as neural representations for downstream tasks, shown schematically in Figure 3. We refer to these representations as  $\nu$ -reps (nu-reps), and similarly, we refer to the datasets of neural representations obtained with our method as  $\nu$ -sets (nu-sets).

As a conditional neural field, the *Neo*MLP backbone encodes the neural field parameters, while the latent variables, *i.e.* the hidden and output embeddings, encode instance-specific information. Each instance (*e.g.* each image in an image dataset) is represented with its own set of latent codes  $\mathbf{Z}_n = \left[\left\{\mathbf{h}_j^n\right\}_{j=1}^H, \left\{\mathbf{o}_k^n\right\}_{k=1}^O\right]$ . We optimize the latent codes for a particular signal by feeding them to the network as inputs alongside a coordinate  $\mathbf{x}_p^{(n)}$ , compute the field value  $\hat{\mathbf{y}}_p^{(n)}$  and the reconstruction loss, and backpropagate the loss to  $\mathbf{Z}_n$  to take one optimization step.

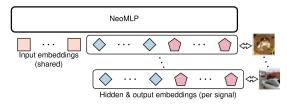


Figure 3: The hidden and output embeddings constitute a set of latent codes for each signal, and can be used as neural representations for downstream tasks. We term these neural representations as  $\nu$ -reps, and the datasets of neural representations as  $\nu$ -sets.

Our method operates in two distinct stages: fitting and finetuning. During fitting, our goal is to optimize the backbone architecture, *i.e.* the parameters of the model. We sample latent codes for all the signals of a fitting dataset and optimize them simultaneously with the backbone architecture. When the fitting stage is complete, after a predetermined set of epochs, we freeze the parameters of the backbone architecture and discard the latent codes. Then, during finetuning, given a new signal, we sample new latent codes for it and optimize them to minimize the reconstruction error for a number of epochs. We finetune the training, validation, and test sets of the downstream task from scratch, even if we used the training set to fit the model, in order to make the distance of representations between splits as small as possible.

In both the fitting and the finetuning stage, we sample completely random points from random signals. This ensures i.i.d. samples, and speeds up the training of our method. During the fitting stage, we also sample points with replacement, as we observed a spiky behaviour in the training loss otherwise. We provide the detailed algorithms of the fitting and the finetuning stage in Algorithms 1 and 2 in Appendix A, respectively. We provide further implementation details in Appendix D.

# 2.4 Using $\nu$ -reps for downstream tasks

After finetuning neural representations, our goal is to use them in downstream tasks, e.g. to train a downstream model for classification or segmentation. Our  $\nu$ -reps comprise a set of latent codes for each signal, corresponding to the finetuned hidden and output embeddings. While the space of  $\nu$ -reps is subject to permutation symmetries, which we discuss in Appendix B, we use a simple downstream model that first concatenates and flattens the hidden and output embeddings in a single vector, and then process it with an MLP. We leave more elaborate methods that exploit the inductive biases present in  $\nu$ -reps for future work.

# 3 Experiments

We gauge the effectiveness of our approach by fitting individual high-resolution signals, as well as datasets of signals. We also evaluate our method on downstream tasks on the fitted datasets. We refer to the appendix for more details. The code is included in the supplementary material and will be open-sourced to facilitate reproduction of the results.

# 3.1 Fitting high-resolution signals

First, we evaluate our method at fitting high-resolution signals. We compare our method against Siren [40], an MLP with sinusoidal activations, RFFNet [42], an MLP with random Fourier features

and ReLU activations, and SPDER [38], an MLP with sublinear damping activations combined with sinusoids. Our goal is to assess the effectiveness of our method in signals of various modalities, and especially in multimodal signals, which have been underexplored in the context of neural fields. Hence, we choose signals that belong to two different modalities, namely an audio clip and a video clip, as well as a multi-modal signal, namely video with audio.

For audio, we follow Siren [40] and use the first 7 seconds from Bach's cello suite No. 1 in G Major: Prelude. The audio clip is sampled at 44.1 kHz, resulting in 308,700 points. For video, we use the "bikes" video from the scikit-video Python library, available online<sup>2</sup>. This video clip lasts for 10 seconds and is sampled at 25 fps, with a spatial resolution of  $272 \times 640$ , resulting in 43,520,000 points. Finally, we explore multimodality using the "Big Buck Bunny" video from scikit-video. This clip lasts for 5.3 seconds. The audio is sampled at 48 kHz and has 6 channels. The original spatial resolution is  $1280 \times 720$  at 25 fps. We subsample the spatial resolution by 2, which results in a resolution of  $640 \times 360$ . Overall, this results in 30,667,776 points (254,976 from audio and 30,412,800 from video). 

**Training details** For audio, we follow Siren [40] and scale the time domain to  $t \in [-100, 100]$  instead of [-1, 1], to account for the high sampling rate of the signal. For the audio-visual data, we model the signal as  $f : \mathbb{R}^3 \to \mathbb{R}^9$ , *i.e.* we have 3 input dimensions (x, y, t), and 9 output dimensions: 3 from video (RGB) and 6 from audio (6 audio channels). Similar to the audio clip, we also scale the time domain, which is now used as the time coordinate for both the audio and the video points. For the points corresponding to audio, we fill their xy coordinates with zeros. Furthermore, since all points come from either the video or the audio modality, we fill the output dimensions that correspond to the other modality with zeros. Finally, during training, we mask these placeholder output dimensions, *i.e.* we compute the loss for the video coordinates using only the RGB outputs, and the loss for the audio coordinates using only the 6-channel audio outputs.

To ensure fairness, for every signal, NeoMLP has approximately the same number of parameters as the baselines. We describe the architecture details for each experiment in Appendix E. We show the results in Table 1, measuring the reconstruction PSNR. We observe that NeoMLP comfortably outperforms the baselines in all three signals. Interestingly, the performance gap is increased in the more difficult setup of multimodal data, which suggests the suitability of our method for multimodal signals. We hypothesize that this can be attributed to our method's ability to learn faster from minibatches with i.i.d. elements, which is something we observed empirically during training and hyperparameter tuning. We visualize example frames for the video clips in Figure 4, and in Figure 6 in Appendix G. We provide further qualitative results in Appendix G and include reconstructions of all signals in the supplementary material.



Figure 4: Examples frames from fitting the "bikes" video clip. The first row shows the groundtruth, while the second and the third row show the reconstructions obtained using *Neo*MLP and Siren, respectively. We observe that *Neo*MLP learns to reconstruct the video with much greater fidelity.

# 3.2 Fitting $\nu$ -sets & Downstream tasks on $\nu$ -sets

Next, we evaluate our method on fitting  $\nu$ -sets, *i.e.* fitting datasets of neural representations of signals with *Neo*MLP, as well as performing downstream tasks on  $\nu$ -sets. We compare our method against Functa [9], DWSNet [25], Neural Graphs [18], and Fit-a-NeF [28]. Functa is a conditional neural field that uses an MLP backbone and conditioning by bias modulation. DWSNet, Neural Graphs,

<sup>2</sup>https://www.scikit-video.org/stable/datasets.html

Table 1: Performance on fitting high resolution signals. We report the PSNR (higher is better).

Method		Dataset				
	Bach	Bach Bikes		k Bunny		
			Audio	Video		
RFFNet [42]	54.62	27.00	32.88	24.59		
Siren [40]	51.65	37.02	31.55	24.82		
SPDER [38]	48.06	33.82	28.45	20.90		
NeoMLP (ours)	54.71	39.06	39.00	34.17		

Table 2: Performance on fitting neural datasets and downstream classification for neural datasets. Experiments on MNIST, CIFAR10, and ShapeNet10. Results from methods marked with † were taken from Fit-a-NeF [28]. The | symbols that appear above and below a number denote that this number is shared for these three methods. For classification, we run the experiments for 3 random seeds and report the mean and standard deviation.

Method	MNIST		CIFAR10		ShapeNet	
	PSNR (†)	Accuracy (%)	PSNR (†)	Accuracy (%)	IoU (†)	Accuracy (%)
Functa [9]	33.07	98.73±0.05	31.90	68.30±0.00	0.434	95.23±0.13
DWSNet [25] †		$85.70 \pm 0.60$		$44.01{\scriptstyle\pm0.48}$		$91.06 \scriptstyle{\pm 0.25}$
Neural Graphs [18] †	14.66	$92.40_{\pm 0.30}$	20.45	$44.11_{\pm 0.20}$	0.559	$90.31_{\pm 0.15}$
Fit-a-NeF [28] †		$96.40_{\pm0.11}$		$39.83_{\pm 1.70}$		$82.96 \scriptstyle{\pm 0.02}$
NeoMLP (ours)	33.98	$98.78 \scriptstyle{\pm 0.04}$	33.16	$73.40{\scriptstyle\pm0.12}$	0.934	$95.30 \scriptstyle{\pm 0.08}$

and Fit-a-NeF, on the other hand, are equivariant downstream models for processing datasets of unconditional neural fields. For these three methods, the process of creating datasets of neural representations corresponds to fitting separate MLPs for each signal in a dataset, a process that is independent of the downstream models themselves. Since these methods have the step of generating the neural datasets in common, we use shared datasets for these methods, provided by Fit-a-NeF.

We consider three datasets, namely MNIST [20], CIFAR10 [19], and ShapeNet10 [4]. We evaluate reconstruction quality for MNIST and CIFAR10 with PSNR, and for ShapeNet with IoU. For CIFAR10, we follow the setup of Functa [9], and use 50 augmentations for all training and validation images during finetuning. For all datasets, we only use the training set as a fitting set, since this closely mimics the real-world conditions for auto-decoding neural fields, namely that test set data can appear after the backbone is frozen, and should be finetuned without changing the backbone.

After fitting the neural datasets, we optimize the downstream model for the downstream tasks, which corresponds to classification for MNIST, CIFAR10, and ShapeNet10. We perform a hyperparameter search for *Neo*MLP to find the best downstream model. Specifically, we use Bayesian hyperparameter search from Wandb [2] to find the best performing hyperparameters for CIFAR10, and reuse these hyperparameters for all datasets.

While neural datasets can easily reach excellent reconstruction quality, it is often at the expense of representation power. This was shown in the case of unconditional neural fields by Papa et al. [28], where optimal downstream performance was often achieved with medium quality reconstructions. Since our goal in this experiment is to optimize the performance of neural representations in downstream tasks, we report the reconstruction quality of the models that achieved the best downstream performance.

We report the results in Table 2. We observe that *Neo*MLP comfortably outperforms DWSNet [25], Neural Graphs [18] and Fit-a-NeF [28], *i.e.* all methods that process unconditional neural fields, both in terms of representation quality and downstream performance. Further, these two quantities seem to be positively correlated for *Neo*MLP, in contrast to the findings of Papa et al. [28] for unconditional neural fields. Our method also outperforms Functa [9] on all three datasets regarding the classification accuracy, while maintaining an excellent reconstruction quality.

#### 3.3 Ablation studies

**Importance of hyperparameters** We perform a large ablation study to assess the importance of the latent codes, and the impact of the duration of fitting and finetuning to the quality of reconstruction and representation power. Specifically, we run two studies on CIFAR10; the first study monitors the number and the dimensionality of the latent codes, as well as the number of finetuning epochs. The second study monitors the number and the dimensionality of the latent codes, as well as the number of fitting epochs. In both studies, all other hyperparameters are fixed. We report the fitting PSNR, the test PSNR and the downstream accuracy. We summarize our findings in Tables 3 and 4.

In both studies, we observe that increasing the number of latents and their dimensionality also increases the reconstruction quality. However, the higher number of latents seems to lead to decreased downstream performance. Furthermore, we notice that increasing the number of finetuning epochs also increases the test PSNR and accuracy. Finally, somewhat surprisingly, while fitting for more epochs leads to noticeably better fitting PSNR, this translates to negligible gain in the test PSNR and accuracy, and even degrades performance in some cases.

Table 3: Ablation study on the importance of the number of latents, the dimensionality of the latents, and the number of finetuning epochs. The backbone is fitted for 50 epochs. Experiment on CIFAR10; no augmentations are used in this study.

Num. latents	Latent dim.	Fit PSNR (↑)	Finetune fo	Finetune for 5 epochs		Finetune for 10 epochs		
			Test PSNR (↑)	Accuracy (%)	Test PSNR (↑)	Accuracy (%)		
6	64	27.04	24.67	51.23	26.00	50.86		
	128	30.01	26.46	53.30	28.41	53.25		
	256	33.10	28.17	53.76	30.82	54.52		
	512	37.49	30.89	54.66	34.98	56.23		
14	64	30.58	26.28	49.36	28.58	49.69		
	128	34.59	28.34	50.74	31.52	51.28		
	256	37.65	29.63	53.35	33.70	54.06		
	512	39.30	30.77	53.26	33.99	53.65		

Table 4: Ablation study on the importance of the number of latents, the dimensionality of the latents, and the number of fitting epochs. The latents are finetuned for 5 epochs. Experiment on CIFAR10; no augmentations are used in this study.

Num. latents	Latent dim.	atent dim. Fit 20 epochs		Fit 50 epochs			
		Fit PSNR (†)	Test PSNR (↑)	Accuracy (%)	Fit PSNR (†)	Test PSNR (↑)	Accuracy (%)
6	64	25.68	24.68	51.03	27.04	24.67	51.23
	128	28.05	26.40	52.67	30.01	26.46	53.30
	256	30.04	28.17	54.56	33.10	28.17	53.76
	512	33.91	30.84	55.14	37.49	30.89	54.66
14	64	28.34	26.18	49.67	30.58	26.28	49.36
	128	31.63	28.03	52.12	34.59	28.34	50.74
	256	33.02	29.24	53.52	37.65	29.63	53.35
	512	31.94	30.54	54.42	39.30	30.77	53.26

**Importance of RFF** As shown by Rahaman et al. [31], neural networks suffer from *spectral bias*, *i.e.* they prioritize learning low frequency components, and have difficulties learning high frequency functions. We expect that these spectral biases would also be present in NeoMLP if left unattended. To that end, we employed Random Fourier Features (RFF) [42] to project our scalar inputs to higher dimensions. Compared to alternatives like sinusoidal activations [40], RFFs allow our architecture to use a standard transformer.

To examine the spectral bias hypothesis, we train *Neo*MLP without RFF, using a learnable linear layer instead. We train this new model on the "bikes" video, and on MNIST. We present the results in Table 5. The study shows that RFFs clearly help with reconstruction quality, both in reconstructing a high-resolution video signal, and on a dataset of images. Interestingly, the reconstruction quality drop from removing RFFs does not translate to downstream performance drop, where, in fact, the model without Fourier features is marginally better than the original.

Table 5: Ablation study on the importance of random Fourier features on (a) the bikes video, (b) on MNIST.

(a) "Bikes" video

(b) MNIST

Method	PSNR (↑)	Method	PSNR (†)	Accuracy (%)
NeoMLP (without RFF)	35.92	NeoMLP (without RFF)	30.33	98.81±0.03
NeoMLP	<b>39.06</b>	NeoMLP	<b>33.98</b>	98.78±0.04

# 4 Related work

**Neural representations** An increasingly large body of works [25, 55, 18, 21, 28, 43, 13] has proposed downstream methods that process datasets of unconditional neural fields, *i.e.* the parameters and the architectures of MLPs. They are all addressing the parameter symmetries present in MLPs, and while the performance of such methods is constantly increasing, it still leaves much to be desired. Closer to our work is another body of works [29, 9, 35, 6, 52, 54, 45] that proposes neural representations through conditional neural fields. Of those, Sajjadi et al. [35], Zhang et al. [52], Wessels et al. [45] have proposed set-latent conditional neural fields that condition the signal through attention [44]. Zhang et al. [52] proposed 3DShape2VecSet, an architecture that employs cross-attention and self-attention to encode shapes into sets of latent vectors and decode them. Our method differs from this method, since it does not rely on cross-attention to fully encode a coordinate in a set of latents. Instead, it employs self-attention, which allows for better information propagation and enables the model to scale to multiple layers.

MLPs as graphs A few recent works [18, 21, 22, 26, 13] have viewed neural networks as graphs and proposed methods that leverage the graph structure. Kofinas et al. [18] focus on the task of processing the parameters of neural networks and represent neural networks as computational graphs of parameters. Their method includes applications to downstream tasks on neural fields. Lim et al. [22] investigate the impact of parameter symmetries, and introduce new neural network architectures that have reduced parameter space symmetries. Nikolentzos et al. [26] show that MLPs can be formalized as GNNs with asynchronous message passing, and propose a model that employs synchronous message passing on a nearly complete graph. Similar to this work, we use a complete graph and employ a synchronous message passing scheme. In contrast to this work, we employ weight-sharing via self-attention and high-dimensional node features. Further, we focus on NeF applications instead of tabular data, and explore conditioning via the hidden and output embeddings.

# 5 Conclusion

In this work, we presented *Neo*MLP, a novel architecture inspired by the principles of connectionism and the graph perspective of MLPs. We perform message passing on the graph of MLPs, after transforming it to a complete graph of input, hidden, and output nodes equipped with high-dimensional features. We also employ weight-sharing through self-attention among all the nodes. *Neo*MLP is a transformer architecture that uses individual input and output dimensions as tokens, along with a number of hidden tokens. We also introduced new neural representations based on the hidden and output embeddings, as well as datasets of neural representations. Our method achieves state-of-the-art performance in fitting high-resolution signals, including multimodal audio-visual data, and outperforms state-of-the-art methods in downstream tasks on neural representations.

**Limitations** Our  $\nu$ -reps are subject to permutation symmetries, indicating that inductive biases can be leveraged to increase downstream performance. Namely, while the output embeddings are already ordered, as they correspond to individual outputs, the hidden embeddings are subject to permutation symmetries. Future work can explore more elaborate methods based on set neural networks, such as Deep Sets [51], that exploit the inductive biases present in  $\nu$ -reps. Further, the latent codes used in  $\nu$ -reps, namely the hidden and output embeddings, carry global information. Instilling locality in latent codes can be useful for fine-grained downstream tasks, such as segmentation. Future work can explore equivariant neural fields [45], which would localize the latent codes by augmenting them with positions or orientations.

# References

- [1] J. L. Ba, J. R. Kiros, and G. E. Hinton. Layer Normalization. arXiv preprint arXiv:1607.06450, 2016. 4
- 325 [2] L. Biewald. Experiment tracking with weights and biases, 2020. URL https://www.wandb.com/.
  326 Software available from wandb.com. 7, 18
- J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula, A. Paszke,
   J. VanderPlas, S. Wanderman-Milne, and Q. Zhang. JAX: composable transformations of Python+NumPy
   programs, 2018. URL http://github.com/jax-ml/jax. 16
- [4] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song,
   H. Su, J. Xiao, L. Yi, and F. Yu. ShapeNet: An Information-Rich 3D Model Repository. Technical Report
   arXiv:1512.03012 [cs.GR], Stanford University Princeton University Toyota Technological Institute
   at Chicago, 2015. 7, 17, 18
- [5] H. Chen, B. He, H. Wang, Y. Ren, S. N. Lim, and A. Shrivastava. Nerv: Neural representations for videos.
   In Advances in Neural Information Processing Systems 34 (NeurIPS), 2021.
- Y. Chen and X. Wang. Transformers as meta-learners for implicit neural representations. In *European Conference on Computer Vision*. Springer, 2022.
- Z. Chen, Y. Chen, J. Liu, X. Xu, V. Goel, Z. Wang, H. Shi, and X. Wang. Videoinr: Learning video implicit neural representation for continuous space-time super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2047–2057, 2022.
- [8] L. de Vries, R. L. M. Van Herten, J. W. Hoving, I. Isgum, B. Emmer, C. B. Majoie, H. Marquering, and
   E. Gavves. Accelerating physics-informed neural fields for fast CT perfusion analysis in acute ischemic
   stroke. In *Medical Imaging with Deep Learning*, 2024.
- [9] E. Dupont, H. Kim, S. Eslami, D. Rezende, and D. Rosenbaum. From data to functa: Your data point is a
   function and you can treat it like one. In *Proceedings of the 39th International Conference on Machine Learning (ICML)*, 2022. 1, 6, 7, 9, 15, 16, 18
- [10] D. Ha, A. Dai, and Q. V. Le. Hypernetworks. arXiv preprint arXiv:1609.09106, 2016.
- [11] R. Hecht-Nielsen. On the algebraic structure of feedforward network weight spaces. In *Advanced Neural Computers*, pages 129–135. Elsevier, 1990.
- 1350 [12] K. Hornik, M. Stinchcombe, and H. White. Multilayer feedforward networks are universal approximators.

  Neural networks, 2(5):359–366, 1989.
- [13] I. Kalogeropoulos, G. Bouritsas, and Y. Panagakis. Scale equivariant graph metanetworks. In Advances in
   Neural Information Processing Systems 37 (NeurIPS), 2024.
- 134 [14] A. Katharopoulos, A. Vyas, N. Pappas, and F. Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International conference on machine learning*, pages 5156–5165. PMLR, 2020. 4
- [15] D. P. Kingma and J. Ba. Adam: A Method for Stochastic Optimization. In 3rd International Conference
   on Learning Representations (ICLR), 2015. 13, 14, 19
- 159 [16] D. M. Knigge, D. R. Wessels, R. Valperga, S. Papa, J.-J. Sonke, E. Gavves, and E. J. Bekkers. Space-time continuous pde forecasting using equivariant neural fields. *arXiv preprint arXiv:2406.06660*, 2024. 1
- [17] M. Kofinas, E. J. Bekkers, N. S. Nagaraja, and E. Gavves. Latent Field Discovery in Interacting Dynamical
   Systems with Neural Fields. In Advances in Neural Information Processing Systems 36 (NeurIPS), 2023.
- [18] M. Kofinas, B. Knyazev, Y. Zhang, Y. Chen, G. J. Burghouts, E. Gavves, C. G. M. Snoek, and D. W.
   Zhang. Graph Neural Networks for Learning Equivariant Representations of Neural Networks. In 12th International Conference on Learning Representations (ICLR), 2024. 1, 2, 6, 7, 9
- 366 [19] A. Krizhevsky, G. Hinton, et al. Learning Multiple Layers of Features from Tiny Images, 2009. 7, 18
- [20] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition.
   Proceedings of the IEEE, 86(11):2278–2324, 1998. 7, 18
- 369 [21] D. Lim, H. Maron, M. T. Law, J. Lorraine, and J. Lucas. Graph metanetworks for processing diverse neural architectures. In 12th International Conference on Learning Representations (ICLR), 2024. 1, 2, 9

- 371 [22] D. Lim, M. Putterman, R. Walters, H. Maron, and S. Jegelka. The empirical impact of neural parameter symmetries, or lack thereof. *arXiv preprint arXiv:2405.20231*, 2024. 2, 9
- [23] W. S. McCulloch and W. Pitts. A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*, 5:115–133, 1943.
- B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. NeRF: Representing
   Scenes as Neural Radiance Fields for View Synthesis. In *European Conference on Computer Vision* (ECCV), 2020. 1
- [25] A. Navon, A. Shamsian, I. Achituve, E. Fetaya, G. Chechik, and H. Maron. Equivariant Architectures
   for Learning in Deep Weight Spaces. In *Proceedings of the 40th International Conference on Machine Learning (ICML)*, 2023. 1, 6, 7, 9
- [26] G. Nikolentzos, S. Wang, J. Lutzeyer, and M. Vazirgiannis. Graph neural machine: A new model for
   learning with tabular data. arXiv preprint arXiv:2402.02862, 2024. 2, 9
- [27] S. Papa, D. M. Knigge, R. Valperga, N. Moriakov, M. Kofinas, J.-J. Sonke, and E. Gavves. Neural
   modulation fields for conditional cone beam neural tomography. In *ICML Workshop on the Synergy of Scientific and Machine Learning Modeling*, 2023.
- [28] S. Papa, R. Valperga, D. M. Knigge, M. Kofinas, P. Lippe, J.-j. Sonke, and E. Gavves. How to Train
   Neural Field Representations: A Comprehensive Study and Benchmark. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 1, 6, 7, 9
- [29] J. J. Park, P. Florence, J. Straub, R. Newcombe, and S. 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. 1, 5, 9, 19
- [30] E. Perez, F. Strub, H. De Vries, V. Dumoulin, and A. Courville. FiLM: Visual Reasoning with a General
   Conditioning Layer. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2018.
- [31] N. Rahaman, A. Baratin, D. Arpit, F. Draxler, M. Lin, F. Hamprecht, Y. Bengio, and A. Courville. On
   the Spectral Bias of Neural Networks. In *Proceedings of the 36th International Conference on Machine Learning (ICML)*, 2019.
- [32] P. Ramachandran, B. Zoph, and Q. V. Le. Searching for Activation Functions. In 6th International
   Conference on Learning Representations, (ICLR), 2018.
- 399 [33] D. Rebain, M. J. Matthews, K. M. Yi, G. Sharma, D. Lagun, and A. Tagliasacchi. Attention beats 400 concatenation for conditioning neural fields. arXiv preprint arXiv:2209.10684, 2022.
- 401 [34] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Learning representations by back-propagating errors.

  402 nature, 323(6088):533–536, 1986. 2
- 403 [35] M. S. Sajjadi, H. Meyer, E. Pot, U. Bergmann, K. Greff, N. Radwan, S. Vora, M. Lučić, D. Duckworth,
  404 A. Dosovitskiy, et al. Scene representation transformer: Geometry-free novel view synthesis through
  405 set-latent scene representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*406 *Pattern Recognition*, pages 6229–6238, 2022. 1, 9
- 407 [36] V. Saragadam, J. Tan, G. Balakrishnan, R. G. Baraniuk, and A. Veeraraghavan. Miner: Multiscale implicit 408 neural representation. In *European Conference on Computer Vision*, 2022.
- V. Saragadam, D. LeJeune, J. Tan, G. Balakrishnan, A. Veeraraghavan, and R. G. Baraniuk. Wire: Wavelet implicit neural representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023.
- 412 [38] K. Shah and C. Sitawarin. Spder: Semiperiodic damping-enabled object representation. In *12th Interna-*413 tional Conference on Learning Representations (ICLR), 2024. 6, 7, 15
- [39] Z. Shen, M. Zhang, H. Zhao, S. Yi, and H. Li. Efficient attention: Attention with linear complexities. In
   Proceedings of the IEEE/CVF winter conference on applications of computer vision, pages 3531–3539,
   2021. 4
- [40] V. Sitzmann, J. Martel, A. Bergman, D. Lindell, and G. Wetzstein. Implicit Neural Representations with
   Periodic Activation Functions. In *Advances in Neural Information Processing Systems 33 (NeurIPS)*, 2020.
   1, 5, 6, 7, 8, 15, 17

- [41] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: a simple way to
   prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958,
   2014. 18
- [42] M. Tancik, P. Srinivasan, B. Mildenhall, S. Fridovich-Keil, N. Raghavan, U. Singhal, R. Ramamoorthi,
   J. Barron, and R. Ng. Fourier Features Let Networks Learn High Frequency Functions in Low Dimensional
   Domains. In Advances in Neural Information Processing Systems 33 (NeurIPS), 2020. 3, 5, 7, 8, 15
- 426 [43] H. V. Tran, T. N. Vo, T. H. Tran, A. T. Nguyen, and T. M. Nguyen. Monomial matrix group equivariant 427 neural functional networks. In *Advances in Neural Information Processing Systems 37 (NeurIPS)*, 2024. 9
- 428 [44] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. 429 Attention Is All You Need. In *Advances in Neural Information Processing Systems 30 (NeurIPS)*, 2017. 1, 430 2, 3, 9
- [45] D. R. Wessels, D. M. Knigge, S. Papa, R. Valperga, S. Vadgama, E. Gavves, and E. J. Bekkers. Grounding continuous representations in geometry: Equivariant neural fields. *arXiv preprint arXiv:2406.05753*, 2024.
   1, 9
- 434 [46] Y. Xie, T. Takikawa, S. Saito, O. Litany, S. Yan, N. Khan, F. Tombari, J. Tompkin, V. Sitzmann, and S. Sridhar. Neural Fields in Visual Computing and Beyond. In *Computer Graphics Forum*, 2022. 1, 2
- [47] R. Xiong, Y. Yang, D. He, K. Zheng, S. Zheng, C. Xing, H. Zhang, Y. Lan, L. Wang, and T. Liu. On layer normalization in the transformer architecture. In *International Conference on Machine Learning*, pages 10524–10533. PMLR, 2020. 4
- 439 [48] Y. Yin, M. Kirchmeyer, J.-Y. Franceschi, A. Rakotomamonjy, and P. Gallinari. Continuous pde dynamics forecasting with implicit neural representations. *arXiv* preprint arXiv:2209.14855, 2022. 1
- [49] T. You, M. Kim, J. Kim, and B. Han. Generative neural fields by mixtures of neural implicit functions. In
   Advances in Neural Information Processing Systems 36 (NeurIPS), 2023.
- [50] C. Yun, S. Bhojanapalli, A. S. Rawat, S. J. Reddi, and S. Kumar. Are transformers universal approximators of sequence-to-sequence functions? In 8th International Conference on Learning Representations (ICLR), 2020.
- 446 [51] M. Zaheer, S. Kottur, S. Ravanbakhsh, B. Poczos, R. R. Salakhutdinov, and A. J. Smola. Deep sets.
   447 Advances in neural information processing systems, 30, 2017.
- 448 [52] B. Zhang, J. Tang, M. Niessner, and P. Wonka. 3dshape2vecset: A 3d shape representation for neural fields and generative diffusion models. *ACM Transactions on Graphics (TOG)*, 42(4):1–16, 2023. 9
- 450 [53] H. Zhang. mixup: Beyond empirical risk minimization. arXiv preprint arXiv:1710.09412, 2017. 18
- 451 [54] S. Zhang, K. Liu, J. Gu, X. Cai, Z. Wang, J. Bu, and H. Wang. Attention beats linear for fast implicit 452 neural representation generation. *arXiv* preprint arXiv:2407.15355, 2024. 9
- 453 [55] A. Zhou, K. Yang, K. Burns, Y. Jiang, S. Sokota, J. Z. Kolter, and C. Finn. Permutation Equivariant Neural 454 Functionals. In *Advances in Neural Information Processing Systems 36 (NeurIPS)*, 2023. 1, 9

# 455 **A** Fitting and finetuning $\nu$ -sets

```
Algorithm 1 Fit NeoMLP as a conditional neural field
Require: Randomly initialized backbone network f_{\Theta}
\textbf{Require:} \ \ \text{Fitting dataset:} \ \mathcal{D}_{\text{fit}} = \left\{ \left\{ \mathbf{x}_p^{(n)}, \mathbf{y}_p^{(n)} \right\}_{p=1}^{P_n} \right\}_{n=1}^{N_{\text{fit}}} \qquad \triangleright N_{\text{fit}} \ \text{signals, Coordinate} \ \mathbf{x}_p^{(n)} \in \mathbb{R}^I
                                                                                                                                                        \triangleright Field value \mathbf{y}_{n}^{(n)} \in \mathbb{R}^{O}
Require: Randomly initialized latents: \mathcal{Z}_{\text{fit}} = \{\mathbf{Z}_n\}_{n=1}^{N_{\text{fit}}}
Require: Initialized optimizer: O_{\text{fit}}
                                                                                                                                                                                 ⊳ Adam [15]
Require: Number of fitting epochs E
Require: Fitting minibatch size B
                                                                                                                                      ▷ Number of points per minibatch
     P \leftarrow \sum_{n=1}^{N_{\mathrm{fit}}} P_n
M \leftarrow \lfloor \frac{P}{B} \rfloor
                                                                                                                             > Total number of points in the dataset
                                                                 ▷ Number of iterations per epoch. We drop incomplete minibatches
     function FITNEOMLP
             for epoch \in \{1, \ldots, E\} do
                    for iteration \in \{1, \dots, M\} do
                           Sample point indices \mathbb{P} = \{p_b\}_{b=1}^B
Sample signal indices \mathbb{S} = \{n_b\}_{b=1}^B
\mathcal{B} \leftarrow \left\{\mathbf{x}_{p_b}^{(n_b)}, \mathbf{y}_{p_b}^{(n_b)}, \mathbf{Z}_{n_b}\right\}_{b=1}^B
                                                                                                                                  \triangleright Sample \mathbb{P} and \mathbb{S} with replacement
                           \hat{\mathbf{y}}_{p_b}^{(n_b)} \leftarrow \mathbf{f}_{\Theta} \Big( \mathbf{x}_{p_b}^{(n_b)}, \mathbf{Z}_{n_b} \Big)
                                                                                                                                              \triangleright In parallel \forall b \in \{1, \dots, B\}
                           \mathcal{L} \leftarrow \frac{1}{B} \sum_{b=1}^{B} \left\| \mathbf{y}_{p_b}^{(n_b)} - \hat{\mathbf{y}}_{p_b}^{(n_b)} \right\|_{2}^{2}\Theta \leftarrow \Theta - O_{\text{fit}}(\nabla_{\Theta} \mathcal{L})
                           \mathbf{Z}_{n_b} \leftarrow \mathbf{Z}_{n_b} - O_{\mathrm{fit}} \Big( \nabla_{\mathbf{Z}_{n_b}} \mathcal{L} \Big)
                                                                                                                                             \triangleright In parallel \forall b \in \{1, \dots, B\}
                    end for
             end for
             Freeze Θ
             Discard \mathcal{Z}_{\text{fit}}
             return ⊖
     end function
```

# 456 B NeoMLP symmetries

Our  $\nu$ -reps, and more specifically, the hidden embeddings, are subject to permutation symmetries. 457 Intuitively, when we permute two hidden embeddings from a randomly initialized or a trained 458 model, we expect the behaviour of the network to remain the same. In this section, we formalize 459 the permutation symmetries present in our method. NeoMLP is a function  $f: \mathbb{R}^{(I+H+O)\times D} \to$ 460  $\mathbb{R}^{(I+H+O)\times D}$  that comprises self-attention and feed-forward networks applied interchangeably 461 for a number of layers, following Equations (2) and (3). As a transformer architecture, it is a 462 permutation equivariant function. Thus, the following property holds:  $f(\mathbf{PX}) = \mathbf{P}f(\mathbf{X})$ , where  $\mathbf{P}$ 463 is a permutation matrix, and X is a set of tokens fed as input to the transformer. 464 Now consider the input to NeoMLP:  $\mathbf{T}^{(0)} = \left[\left\{\mathbf{i}_i\right\}_{i=1}^{I}, \left\{\mathbf{h}_j\right\}_{j=1}^{H}, \left\{\mathbf{o}_k\right\}_{k=1}^{O}\right], \mathbf{T}^{(0)} \in \mathbb{R}^{(I+H+O)\times D}$ 465 We look at two cases of permutations, namely permuting only the hidden neurons, and permuting only 466 the output neurons. The permutation matrix for the first case, i.e. permuting only the hidden neurons, 467 is  $\mathbf{P}_1 = \mathbf{I}_{I \times I} \oplus \mathbf{P}_{H \times H} \oplus \mathbf{I}_{O \times O}$ , where  $\mathbf{I}$  is the identity matrix,  $\mathbf{P}_{H \times H}$  is a permutation matrix, and 468 ## denotes the direct sum operator, i.e. stacking matrix blocks diagonally, with zero matrices in the 469 off-diagonal blocks. Each  $P_1$  corresponds to a permutation  $\pi_1 \in S_H$ . Applying this permutation to  $T^{(0)}$  permutes only the hidden neurons:

$$\mathbf{P}_{1}\mathbf{T}^{(0)} = \left[ \left\{ \mathbf{i}_{i} \right\}_{i=1}^{I}, \left\{ \mathbf{h}_{\pi_{1}^{-1}(j)} \right\}_{j=1}^{H}, \left\{ \mathbf{o}_{k} \right\}_{k=1}^{O} \right]$$
(10)

# Algorithm 2 Finetune NeoMLP as a conditional neural field

```
Require: Frozen backbone network f_{\Theta}
Require: Train, validation, test datasets: \mathcal{D}_{\text{train}}, \mathcal{D}_{\text{validation}}, \mathcal{D}_{\text{test}}
Require: Randomly initialized latents: \mathcal{Z}_{\text{train}}, \mathcal{Z}_{\text{validation}}, \mathcal{Z}_{\text{test}}
Require: Initialized optimizers: O_{\text{train}}, O_{\text{validation}}, O_{\text{test}}
                                                                                                                                                                                                                 ⊳ Adam [15]
Require: Number of finetuning epochs E'
Require: Finetuning minibatch size B'
     function FINETUNENEOMLP
               for split \in {train, validation, test} do
                       \begin{aligned} M_{\text{split}} &\leftarrow \lceil \frac{\sum_{n=1}^{N_{\text{split}}} P_n}{B'} \rceil \\ \text{for epoch} &\in \{1, \dots, E'\} \text{ do} \\ \text{for iteration} &\in \{1, \dots, M_{\text{split}}\} \text{ do} \end{aligned}
                                         Sample point indices \mathbb{P} = \{p_b\}_{b=1}^{B'}
                                        Sample signal indices \mathbb{S} = \{n_b\}_{b=1}^{B'}
\mathcal{B} \leftarrow \left\{\mathbf{x}_{p_b}^{(n_b)}, \mathbf{y}_{p_b}^{(n_b)}, \mathbf{Z}_{n_b}\right\}_{b=1}^{B'}
\hat{\mathbf{y}}_{p_b}^{(n_b)} \leftarrow \mathbf{f}_{\Theta}\left(\mathbf{x}_{p_b}^{(n_b)}, \mathbf{Z}_{n_b}\right)
                                                                                                                                                 \triangleright Sample \mathbb{P} and \mathbb{S} without replacement
                                                                                                                                                                     \triangleright In parallel \forall b \in \{1, \dots, B'\}
                                        \mathcal{L} \leftarrow \frac{1}{B'} \sum_{b=1}^{B'} \left\| \mathbf{y}_{p_b}^{(n_b)} - \hat{\mathbf{y}}_{p_b}^{(n_b)} \right\|_2^2
\mathbf{Z}_{n_b} \leftarrow \mathbf{Z}_{n_b} - O_{\text{split}} \left( \nabla_{\mathbf{Z}_{n_b}} \mathcal{L} \right)
                                                                                                                                                                     \triangleright In parallel \forall b \in \{1, \dots, B'\}
                       end for
               end for
               return \mathcal{Z}_{train}, \mathcal{Z}_{validation}, \mathcal{Z}_{test}
     end function
```

Next, we apply *Neo*MLP on the permuted inputs. Making use of the equivariance property, the output of the function applied to the permuted inputs is equivalent to the permutation of the output of the function applied to the original inputs.

$$f\left(\mathbf{P}_{1}\mathbf{T}^{(0)}\right) = \mathbf{P}_{1}f\left(\mathbf{T}^{(0)}\right) \tag{11}$$

Since the network is only using the output tokens in the final step as an output of the network, the overall behaviour of *Neo*MLP is invariant to the permutations of the hidden nodes.

We can follow the same principle to show that permuting the output nodes results in different outputs. The permutation matrix in this case is  $\mathbf{P}_2 = \mathbf{I}_{I \times I} \oplus \mathbf{I}_{H \times H} \oplus \mathbf{P}_{O \times O}$ . The equivariance property still holds, namely  $f(\mathbf{P}_2\mathbf{T}^{(0)}) = \mathbf{P}_2f(\mathbf{T}^{(0)})$ . However, the output tokens are now used as the output of the network. This means that permuting the output tokens would result in permuting the output dimensions of a signal, which is clearly not equivalent to the original signal.

A corollary of the permutation symmetries is that if we start with a randomly initialized model,

482 A coronary of the permutation symmetries is that if we start with a randomly initialized model, apply a permutation on the hidden nodes to create another model, and then train the two models independently, these two trained models would be identical up to the permutation of the hidden nodes. This observation is important for downstream tasks, as it shows the existence of equivalence classes that should be taken into account by the downstream models.

# C Computational complexity

487

While *Neo*MLP comfortably outperforms Siren in the task of fitting high-resolution signals, it is also more computationally expensive. We quantitatively measure the computational complexity of our method using the fvcore library<sup>3</sup>. We evaluate on the "bikes" video signal, and use the hyperparameters described in Appendix E. We report the FLOPs for 1 input (*i.e.* 1 coordinate) in the forward pass. *Neo*MLP has 51.479 MFLOPs, out of which 17.83 MFLOPs correspond to the

<sup>3</sup>https://github.com/facebookresearch/fvcore

attention itself and 33.55 MFLOPs correspond to the FFNs. In the same setup, Siren [40] has 3.15 MFLOPs.

Despite having a higher computational complexity compared to the baselines, NeoMLP can actually fit high resolution signals faster, and does so while having a smaller memory footprint, since it can make use of small batch sizes. Figure 5 shows the runtime of NeoMLP for fitting high-resolution signals, compared to the baselines. The x-axis represents wall time in seconds and the y-axis represents the reconstruction quality (PSNR). Table 6 shows the corresponding GPU memory and batch size, along with the total runtime for fitting high resolution signals.

501

502

503

504

Finally, despite the large difference in FLOPs, the forward pass of *Neo*MLP is almost as fast as the forward pass of Siren, considering the same batch size. Namely, we ran a full evaluation on the "bikes" signal, on an Nvidia H100 GPU, using a batch size of 32,768. *Neo*MLP takes 139.74 seconds, while Siren takes 131.01 seconds. *Neo*MLP, however, cannot fit larger batch sizes in memory, while Siren can fit as big as 1,048,576. With this batch size, Siren requires 79.18 seconds for a full evaluation.

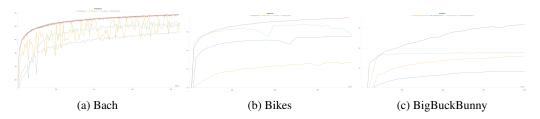


Figure 5: Runtime for fitting high-resolution signals. The *x*-axis represents wall time in seconds and the *y*-axis represents the reconstruction quality (PSNR). *Neo*MLP fits signals faster and with better reconstruction quality.

Table 6: Runtime, GPU memory, and batch size on fitting high resolution signals. For each dataset, we trained all methods for the same amount of time for fair comparison.

	(a) Bach					
Method	GPU memory (GB)	Batch size	Runtime (hours)			
RFFNet [42]	3.7	308,207	2.33			
Siren [40]	3.9	$308,\!207$				
SPDER [38]	6.0	$308,\!207$				
NeoMLP (ours)	2.2	4,096	ĺ			
Method	GPU memory (GB)	Batch size	Runtime (hours)			
RFFNet [42]	11.2	262,144	19.07			
Siren [40]	16.8	262,144				
SPDER [38]	37.3	262,144				
NeoMLP (ours)	11.1	4,096				
	(c) BigBuckB	unny				
Method	GPU memory (GB)	Batch size	Runtime (hours)			
RFFNet [42]	13.9	262,144	24.73			
Siren [40]	18.7	262,144				
<b>SPDER</b> [38]	39.2	262,144				
NeoMLP (ours)	13.2	4,096				

We also monitor the runtime of *Neo*MLP on fitting datasets of signals, and compare against Functa [9]. We report the results in Table 7. *Neo*MLP consistently exhibits lower runtimes for the fitting stage,

while Functa is much faster during the finetuning stage, which can be attributed to the meta-learning employed for finetuning, and the highly efficient JAX [3] implementation. As noted by Dupont et al. [9], however, meta-learning may come at the expense of limiting reconstruction accuracy for more complex datasets, since the latent codes lie within a few gradient steps from the initialization.

Table 7: Runtime on fitting datasets of signals. The finetuning runtime is measured on the test set only. The runtime for fitting is measured in minutes, while the runtime for finetuning is measured in seconds.

#### (a) MNIST

Method	Fitting		Finetuning		
	Num. epochs	Runtime (min.)	Num. epochs	Runtime (sec.)	
Functa [9]	192	240	3	16	
NeoMLP (ours)	20	63	10	318	
(I.) CITADIO					

#### (b) CIFAR10

Method	Fitting		Fine	tuning
	Num. epochs Runtime (min.)		Num. epochs	Runtime (sec.)
Functa [9] NeoMLP (ours)	213 50	418 305	3 10	16 646

# (c) ShapeNet

Method	Fitting		Fine	tuning
	Num. epochs Runtime (min.)		Num. epochs	Runtime (sec.)
Functa [9]	20	1002	3	250
NeoMLP (ours)	20	713	2	1680

# D Implementation details

# 514 D.1 Embedding initialization

513

517

519

520

Fitting high-resolution signals We initialize input embeddings by sampling from a normal distribution with variance  $\sigma_i^2=1$ . For hidden and output embeddings, we use a variance  $\sigma_o^2=1e-3$ .

Fitting  $\nu$ -sets During fitting, we initialize the input, hidden, and output embeddings by sampling a normal distribution with variance  $\sigma_i^2 = \sigma_o^2 = 1e - 3$ . During finetuning, we sample embeddings for new signals from a normal distribution with variance  $\sigma_o^2 = 1e - 3$ .

# D.2 Weight initialization

We initialize the bias of the final output linear layer to zeros, as we observe this leads to faster convergence and better stability at the beginning of training. Further, we initialize the weights of the linear projection following the random Fourier features by sampling from a normal distribution  $\mathcal{N}\left(0,\frac{2}{D_{\text{RFF}}}\right)$ . This results in a unit normal distribution of the inputs after the linear projection.

# E Experiment details

# 526 E.1 High-resolution signals

Below we provide the hyperparameters for *Neo*MLP.

Audio (Bach)

- Number of parameters: 182, 017
- FFN hidden dim: 256
- Token dimensionality D: 64
- Number of self-attention heads: 4
- Number of layers: 3
- RFF dimensionality  $D_{RFF}$ : 512
- 535 **–** RFF  $\sigma$ : 20
- Total number of nodes: 8
- Number of epochs: 5,000
- Batch size: 4,096
- Learning rate: 0.005
- Video (Bikes) & Video with audio (Big Buck Bunny)
- Number of parameters: 3, 189, 249
- **FFN hidden dim:** 1,024
- Token dimensionality D: 256
- Number of self-attention heads: 8
- Number of layers: 4
- RFF dimensionality  $D_{RFF}$ : 128
- **RFF**  $\sigma$ : 20
- Total number of nodes: 16
- Number of epochs: 200 (400 for BigBuckBunny)
- Batch size: 4,096
- Learning rate: 0.0005
- For the audio fitting, Siren [40] has 198,145 parameters. It is a 5-layer MLP, with a hidden dimension of 256, and it is trained with full batch training and a learning rate of 5e 5.
- 555 of 250, and it is trained with rail batch training and a featining rate of 60 of
- For the video fitting, Siren has 3,155,971 parameters, and for the audio-visual data, Siren has
- 555 3,162,121 parameters. It both settings, it is using the exact same architecture with 5 layers and a
- hidden dimension of 1024. We train it with a learning rate of 1e-4 and a batch size of 262,144.

# 557 E.2 Fitting $\nu$ -sets

- For ShapeNet10 [4], we fit the dataset for 20 epochs. In each epoch, we stop when we have used 10% of the available points, which effectively results in 2 epochs in total. We finetune for 2 epochs,
- and use the 20% of the available points. We use a minibatch size of 32,768 points, and a learning rate of 0.005. The healthough has the following hyperpersemeters:
- of 0.005. The backbone has the following hyperparameters:
- FFN hidden dim: 512
- Token dimensionality D: 256
- Number of self-attention heads: 4
- Number of layers: 3
- RFF dimensionality  $D_{\rm RFF}$ : 512
- **•** RFF *σ*: 20
- Total number of nodes: 8
- For MNIST, we fit the dataset for 20 epochs and finetune for 10 epochs. We use a minibatch of 12,288
- points (the equivalent of 16 images), and a learning rate of 0.005. The backbone has the following
- 571 hyperparameters:
- FFN hidden dim: 512
- Token dimensionality D: 256
- Number of self-attention heads: 4

- Number of layers: 3
- RFF dimensionality  $D_{\rm RFF}$ : 512
- RFF σ: 20
- Total number of nodes: 8
- For CIFAR10, we fit the dataset for 50 epochs and finetune for 20 epochs. We use a minibatch of
- 16,384 points (the equivalent of 16 images), and a learning rate of 0.005. The backbone has the
- 581 following hyperparameters:
- FFN hidden dim: 128
- Token dimensionality D: 512
- Number of self-attention heads: 4
- Number of layers: 3
- RFF dimensionality  $D_{RFF}$ : 128
- 587 RFF σ: 20
- Total number of nodes: 8

# 589 E.3 Downstream tasks on $\nu$ -sets

- 590 We perform a hyperparameter search for *Neo*MLP to find the best downstream model. Specifically,
- we use Bayesian hyperparameter search from Wandb [2] to find the best performing hyperparameters
- for CIFAR10, and reuse these hyperparameters for all datasets. We perform our search over the choice
- of Mixup [53], batch size, learning rate, noise added to the data, data dropout, hidden dimension and
- model dropout [41].
- Our downstream model is a 3 layer MLP with SiLU activations [32], a hidden dimension of 2048,
- and dropout of 0.3. We train the model with a learning rate of 8e 3, and batch size of 256. We
- use Mixup, weight decay with  $\lambda = 0.05$ , and add noise to the data with scale 0.05. Finally, we use
- weight averaging with exponential moving average (EMA).
- For CIFAR10 [19], the model takes as input 6 embeddings (the *Neo*MLP had 8 nodes in total). We
- 600 train for 100 epochs.
- For ShapeNet10 [4], the model takes as input 13 embeddings (the *Neo*MLP had 16 nodes in total).
- We use a higher weight decay  $\lambda = 0.25$  to further prevent overfitting, and train for 500 epochs.
- For MNIST [20], the model takes as input 6 embeddings (the *Neo*MLP had 8 nodes in total). We use
- a higher weight decay  $\lambda = 0.2$  and train for 500 epochs.

# 605 F Dataset details

# 606 F.1 ShapeNet10

- We use the following 10 classes for ShapeNet10 classification: loudspeaker, bench, watercraft, lamp,
- rifle, sofa, cap, airplane, chair, table.
- The dataset comprises 35,984 shapes. We use 29,000 shapes for training, 2,000 as a validation set,
- and 4,984 as a test set.
- For CIFAR10, following Functa [9], we use 50 augmentations per training and validation image. This
- results in a total of 2,500,000 training and validation images. We use 5,000 of those for validation.

# 613 G Qualitative results

- We show example frame for the "BigBuckBunny" video clip in Figure 6.
- We show the reconstructions for the "Bach" audio clip in Figure 7, and the errors between the
- groundtruth signal and reconstructions in Figure 8.



Figure 6: Examples frames from fitting the "BigBuckBunny" video clip. The first row shows the groundtruth, while the following rows show the reconstructions obtained using *Neo*MLP, RFFNet, Siren, and SPDER, respectively. We observe that *Neo*MLP learns to reconstruct the video with much greater fidelity.

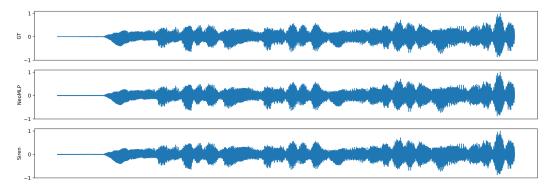


Figure 7: Predictions for the "Bach" audio clip. The first row shows the groundtruth signal, while the second and third row show the reconstructions from *Neo*MLP and Siren, respectively.

# 617 H Visualizations

# 618 I Latent space

We visualize the learned MNIST data manifold for a two-dimensional latent space (with a single embedding) in Figure 10, following Kingma and Ba [15], Park et al. [29]. We assume that the latent

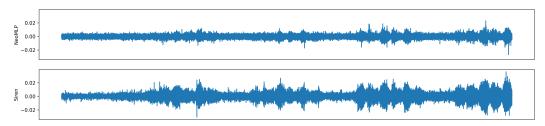


Figure 8: Errors  $\epsilon = y - \hat{y}$  between predictions y and groundtruth  $\hat{y}$ . The top row shows the error for NeoMLP, while the bottom row shows the error for Siren. Both the x-axis and the y-axis are shared in this figure, but the y-axis is different from Figure 7. We see that the errors from Siren have a much larger amplitude, and still seem to capture signal components.

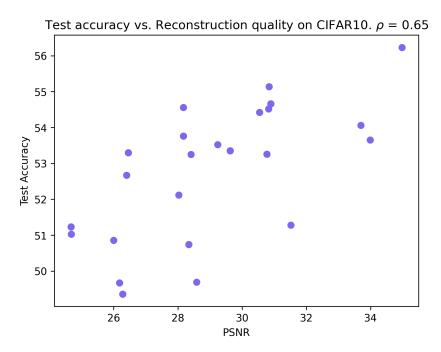


Figure 9: Test accuracy vs. reconstruction quality (PSNR). Experiments on CIFAR10, with different hyperparameters, *without* augmentations.

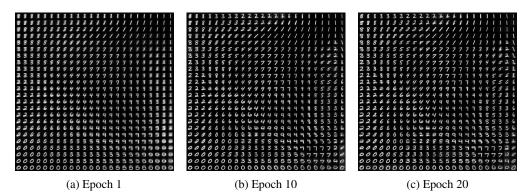


Figure 10: The data manifold of *Neo*MLP with a 2D latent space and a single embedding, *i.e.* O = 1, H = 0, D = 2. We visualize the manifold as the fitting stage progresses.

space is Gaussian, with a sample mean and variance estimated from the latents of the training set. We sample linearly spaced coordinates on the unit square and transform them through the Percent Point Function (PPF) of the Gaussian to produce the values of the latent variables.

We also visualize random samples from the latent space of *Neo*MLP across hyperparameter configurations varying in the number of embeddings and dimensionality of the latents in Figure 11, as well as the corresponding reconstruction quality and downstream performance in Table 8.

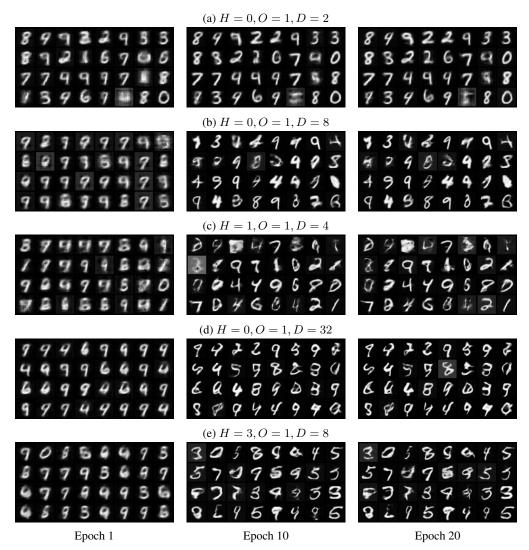


Figure 11: Random samples from the latent space of *Neo*MLP, as the fitting stage progresses. We visualize various configurations of the number of embeddings and the dimension of the latents.

Table 8: Reconstruction quality and downstream performance for the configurations corresponding to Figure 11.

Num. latents	Latent dim.	Fit PSNR (↑)	Accuracy (%)
1	2	15.61	23.7
1	8	20.47	48.2
2	4	20.19	49.3
1	32	24.75	76.2
4	8	24.44	72.5

# NeurIPS Paper Checklist

# 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: See Section 3.

# Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the
  contributions made in the paper and important assumptions and limitations. A No or
  NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals
  are not attained by the paper.

#### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: Yes

Justification: See Section 5.

# Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

# 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: We discuss the permutation symmetries of *Neo*MLP in Appendix B. Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

# 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Equations (2) and (4) mathematically describe our method. Further, we describe the algorithms for fitting and finetuning *Neo*MLP in Algorithms 1 and 2, respectively. We report details regarding the implementation in Appendix D, dataset details in Appendix F, and details about the hyperparameters used in each experiment in Appendix E.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

# 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: Our source code is included in the supplementary material. We use publicly available data and datasets, which are described in Section 3.

#### Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be
  possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not
  including code, unless this is central to the contribution (e.g., for a new open-source
  benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how
  to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

# 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We report details regarding the implementation in Appendix D, dataset details in Appendix F, and details about the hyperparameters used in each experiment in Appendix E.

# Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

#### 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: See Table 2.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

# 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824 825

826

827

828

829

830

831

832

Justification: See Appendix C.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

# 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: Neither our work nor our data involve human subjects.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

# 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

Guidelines:

• The answer NA means that there is no societal impact of the work performed.

- If the authors answer NA or No, they should explain why their work has no societal
  impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

#### 11. Safeguards

ຂຂດ

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]
Justification:

#### Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

# 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: See Section 3.

# Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- If assets are released, the license, copyright information, and terms of use in the
  package should be provided. For popular datasets, paperswithcode.com/datasets
  has curated licenses for some datasets. Their licensing guide can help determine the
  license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

#### 13. New assets

885

886

887

888

889

890

891

892

893 894

895

896

897

898 899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

935

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]
Justification:

#### Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

# 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA] Justification:

# Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

# 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]
Justification:

# Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent)
  may be required for any human subjects research. If you obtained IRB approval, you
  should clearly state this in the paper.

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

# 16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]
Justification:
Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.