# Learned Initializations for Optimizing Coordinate-Based Neural Representations

Matthew Tancik[*1]     Ben Mildenhall[*1]     Terrance Wang[1]     Divi Schmidt[1]
Pratul P. Srinivasan[2]     Jonathan T. Barron[2]     Ren Ng[1]
[1]UC Berkeley     [2]Google Research

## Abstract

*Coordinate-based neural representations have shown significant promise as an alternative to discrete, array-based representations for complex low dimensional signals. However, optimizing a coordinate-based network from randomly initialized weights for each new signal is inefficient. We propose applying standard meta-learning algorithms to learn the initial weight parameters for these fully-connected networks based on the underlying class of signals being represented (e.g., images of faces or 3D models of chairs). Despite requiring only a minor change in implementation, using these learned initial weights enables faster convergence during optimization and can serve as a strong prior over the signal class being modeled, resulting in better generalization when only partial observations of a given signal are available. We explore these benefits across a variety of tasks, including representing 2D images, reconstructing CT scans, and recovering 3D shapes and scenes from 2D image observations.*

## 1. Introduction

Recent work has demonstrated the potential of representing complex low-dimensional signals using deep fully-connected neural networks (typically referred to as multi-layer perceptrons, or MLPs). A coordinate-based neural representation $f_\theta$ for a given signal is an MLP (with weights $\theta$) that is optimized to map from an input coordinate $\mathbf{x}$ to the signal's value at that coordinate. For example, $f_\theta$ could map from 2D pixel coordinates to RGB color values to encode an image. Unlike a signal stored as a discretely sampled array of values, a coordinate-based neural representation is *continuous* and is not constrained to have a fixed spatial resolution. This fact has recently been exploited to design representations for 3D shapes (which typically occupy a small 2D subset of 3D space) that do not require cubic storage complexity, in contrast to 3D voxel grids [23, 25, 28, 35].
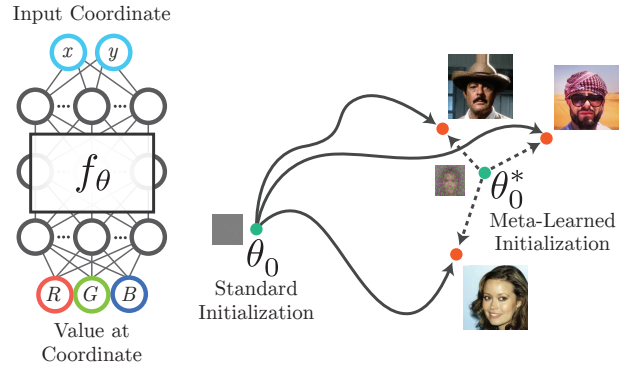


Figure 1. A coordinate-based MLP, illustrated on the left, takes a coordinate as input and outputs a value at that location. For example, the network could take in a pixel coordinate $(x, y)$ and emit the $(R, G, B)$ color at that pixel as output, thereby representing a 2D image. The network weights $\theta$ are typically optimized via gradient descent to produce the desired image, as depicted on the right. However, finding good parameters can be computationally expensive, and the full optimization process must be repeated for each new target. We propose using meta-learning to find initial network weights $\theta_0^*$ that allow for faster convergence and better generalization.

However, one limitation of these neural representations is that computing network weights $\theta$ that reproduce a given signal typically requires solving an optimization problem by running many steps of gradient descent. This can take between seconds (when encoding a small image) and hours (when solving an inverse problem to recover a high resolution radiance field, as in NeRF [25]). Common approaches to address this issue include concatenating a latent vector to the input coordinate and supervising a single neural network to represent an entire class of signals [23, 28], or training a hypernetwork to map from signal observations (or a latent code) to MLP weights [34, 35]. However, each of these strategies is restricted to representing only signals within its learned latent space, potentially limiting its ability to express previously unseen target signals.

Recent work [33] has shown that optimization-based meta-learning can dramatically reduce the number of gra-

---

dient descent steps required to optimize a neural representation to encode a new signal in the case of signed distance fields of 2D and 3D shapes. In this work, we propose learning the weight initialization for neural representations across a wide variety of underlying signal types, such as images, volumetric data, and 3D scenes. We show that compared to a standard random initialization, using fixed, learned values for the initial network weights acts as a strong prior that enables both faster convergence during optimization and better generalization when only partial observations of the target signal are available. In the context of using neural representations for 3D reconstruction from images, a learned initialization specialized to a particular ShapeNet [3] class allows the network to recover 3D shape from a single image over the course of optimization, whereas a standard randomly initialized network fails unless provided with multiple input views. Given a meta-training set consisting of observations of different signals sampled from a fixed underlying class, our setup applies an optimization-based meta-learning algorithm (MAML [7] or Reptile [26]) in order to produce initial weights better suited for representing that specific signal class (e.g., face images from CelebA [21] or 3D chairs from ShapeNet [3]).

The biggest advantage of our approach is its simplicity. Given an existing framework for test-time optimization of a neural representation, implementing an outer loop with MAML or Reptile update steps only requires a few extra lines of code and a dataset of training examples. Once the meta-learning phase is complete, the learned initial weights can be stored and later reloaded in place of a standard network initialization whenever a new signal needs to be encoded. This minor implementation change can significantly alter the behavior of the network during optimization.

## 2. Related Work

**Neural Representations**   Neural representations have recently risen to prominence as compact representations for 3D shapes. These methods represent shapes as implicit surfaces defined as a level set of an MLP network and enable full object reconstruction from incomplete 3D point cloud data or depth scans [4, 5, 10, 11, 15, 23, 24, 28]. Later work combined this idea with various formulations of differentiable rendering to recover neural representations of 3D shape using only 2D image observations [19, 20, 25, 27, 35, 38].

Coordinate-based neural networks have also been used to represent other low-dimensional signals, such as 2D images, where such networks (when trained via genetic algorithms) have been referred to as compositional pattern–producing networks [36]. Recent works have shown that standard ReLU MLPs fail to adequately represent fine details in these complex low-dimensional signals due to a spectral bias [29] and address this issue by either replacing

the ReLU activations with sine functions [34] or by lifting the input coordinates into a Fourier feature space [37]. Our work makes use of these observations and presents a technique that enables a coordinate-based MLP to learn from the process of fitting many signals within a category so that it can quickly optimize to fit any new signal using fewer steps and fewer observations.

**Meta-learning**   Meta-learning typically addresses the problem of few-shot learning, where some examples of a given task (including training and test data) are used to learn an algorithm that achieves better performance on new, previously unseen instances of the same task. A prototypical example from computer vision is few-shot image classification, where a network must learn to differentiate between new classes at test time based on only a small number of labeled instances of each class.

Most relevant to this work are optimization-based meta-learning algorithms such as Model-Agnostic Meta Learning (MAML) [7] and Reptile [26], as well as various extensions [1, 6, 8, 18, 30]. Given a network architecture for performing a task, these methods use an outer loop of gradient-based learning to find a weight initialization that allows the network to more efficiently optimize for new instances of the underlying task at test time. These methods assume the use of a standard gradient-based optimization method such as stochastic gradient descent or Adam [17] at test time, making them easy to layer on top of existing implementations, as opposed to more complex methods such as Ravi *et al*. [31], which trains a "meta-learner" LSTM network to perform gradient updates for the underlying task. An exhaustive review of meta-learning algorithms is provided in the survey paper by Hospedales *et al*. [13].

MetaSDF [33] specifically applies this idea of learning a weight initialization to the task of fitting neural representations to represent signed distance fields, and shows that this strategy achieves much more rapid convergence than standard approaches such as DeepSDF [28]. Our work applies meta-learning to neural representations for a wider variety of underlying signal types and further explores the power of using initial weight settings as a prior.

## 3. Overview

We define a finite signal $T$ as a function mapping from a bounded set $C \in \mathbb{R}^d$ to $\mathbb{R}^n$, where we refer to elements $\mathbf{x} \in C$ as $d$-dimensional coordinates. Examples include images (mapping from 2D pixel coordinates to 3D color values) or volumetric representations for 3D shapes (mapping from 3D locations to 4D tuples of color and density). A coordinate-based neural representation $f_\theta$ for $T$ is a fully connected neural network with $d$ input and $n$ output channels whose weights $\theta$ are optimized such that $f_\theta$ matches $T$ as closely as possible for all coordinates in $\mathbf{x} \in C$.

If direct pointwise observations $\{(\mathbf{x}_i, T(\mathbf{x}_i)\}_i$ of the signal $T$ are available, $f_\theta$ can be supervised by gradient descent using a simple L2 loss:

$$L(\theta) = \sum_i \|f_\theta(\mathbf{x}_i) - T(\mathbf{x}_i)\|_2^2. \qquad (1)$$

Let $\theta_0$ denote the initial network weights before any gradient steps are taken, and let $\theta_i$ denote the weights after $i$ steps of optimization. Basic gradient descent applies the rule:

$$\theta_{i+1} = \theta_i - \alpha \nabla_\theta L(\theta)|_{\theta=\theta_i}, \qquad (2)$$

with a learning rate parameter $\alpha$, whereas more sophisticated optimizers such as Adam [17] keep track of gradient moments over time to redirect the optimization trajectory. Given a fixed budget of $m$ optimization steps, different initial weight values $\theta_0$ will result in different final weights $\theta_m$ and signal approximation error $L(\theta_m)$. When emphasizing the functional dependence of $\theta_m$ on the initial weights and a particular signal, we will write $\theta_m(\theta_0, T)$.

It is often the case that only indirect observations of $T$ are available, taken through some forward measurement model $M(T, \mathbf{p})$. For example, if $T$ is a 3D object, $M(T, \mathbf{p})$ could be a 2D image captured of the object from camera pose $\mathbf{p}$. In this case, recovering a neural representation for $T$ from observations $\{\mathbf{p}_i, M(T, \mathbf{p}_i)\}_i$ requires solving an inverse problem by taking gradient steps on a loss that incorporates the forward model $M$:

$$L_M(\theta) = \sum_i \|M(f_\theta, \mathbf{p}_i) - M(T, \mathbf{p}_i)\|_2^2. \qquad (3)$$

If $M$ discards too much information about $T$ or the set of provided observations is too small, the resulting network $f_\theta$ may not match $T$ closely. For example, accurately recovering a 3D object from a single 2D view may not be possible without strong a priori knowledge of the object's shape.

## 3.1. Optimizing initial weights

We assume that we are given a dataset of observations of signals $T$ from a particular distribution $\mathcal{T}$ (e.g., 2D face images or 3D chairs) and our goal is to find initial weights $\theta_0^*$ that will result in the lowest possible final loss $L(\theta_m)$ when optimizing a network $f_\theta$ to represent a new, previously unseen signal from the same distribution:

$$\theta_0^* = \arg\min_{\theta_0} E_{T\sim\mathcal{T}}[L(\theta_m(\theta_0, T))] \qquad (4)$$

This problem of trying to learn the initial weights of a network to serve as a good starting point for gradient descent across a distribution of tasks is addressed by a variety of optimization-based meta-learning algorithms, such as MAML [7] and Reptile [26].

**MAML [7]** Given a task $T$, calculating the weight values $\theta_m(\theta_0, T)$ requires taking $m$ optimization steps, which are collectively referred to as the *inner loop*. MAML wraps an *outer loop* of meta-learning around this inner loop in order to learn the initial weights $\theta_0$. Each outer loop samples a signal $T_j$ from $\mathcal{T}$ and applies the update rule:

$$\theta_0^{j+1} = \theta_0^j - \beta \nabla_\theta L(\theta_m(\theta, T_j))|_{\theta=\theta_0^j} \qquad (5)$$

with meta-learning step size $\beta$. This update rule applies gradient descent to the loss on the weights $\theta_m(\theta_0^j, T_j)$ resulting from the inner loop optimization.

**Reptile [26]** Reptile uses the same meta-learning setup as MAML but applies a simpler update rule that does not require calculating second-order gradients:

$$\theta_0^{j+1} = \theta_0^j - \beta(\theta_m(\theta_0^j, T_j) - \theta_0^j). \qquad (6)$$

This rule moves the previous weight initialization $\theta_0^j$ in the direction of the task-optimized weights $\theta_m(\theta_0^j, T_j)$.

## 3.2. Experimental setup

The meta-learning algorithms described previously are conceptually simple, requiring no changes to the architecture or optimization procedure of a coordinate-based neural representation when given a new signal to encode at "test time" (after meta-learning is complete). These algorithms produce only a set of initial network weights $\theta_0^*$ that are then used as a starting point for gradient descent. Test-time optimization on new signals is not limited to the same number of steps $m$ as were used in the inner loop during meta-learning; indeed, at test time we often observe benefits from optimizing for significantly more iterations than were used during the inner loop of the meta-learning algorithm.

MAML is typically able to produce a better initialization than Reptile given a fixed number of inner loop steps $m$, but Reptile can be unrolled for more inner loop steps because it is less memory-intensive than MAML. For some tasks, MAML's limited number of inner loop steps means that it can only observe a small percentage of the observations of a target signal. In these cases, we use Reptile to maximize the number of different observations seen over the course of the inner loop. Experimentally we find it beneficial to unroll more steps for more complex tasks.

Each of our experiments involves two phases:
1. *Meta-learning*, where we use MAML or Reptile in combination with a training dataset of example tasks (observations of different signal instances) to optimize initial network weights for that class of signals, and
2. *Test-time optimization*, where we use standard gradient-based optimization to fit the weights of a network to observations of a previously unseen signal from the same class.

We aim to answer the following question: how do different initial network weight settings influence the ability of a neural representation to fit to a new signal during test-time optimization?

# 4. Results

We present results on 2D image regression, 2D computed tomography (CT) reconstruction, 3D object reconstruction, and 3D scene reconstruction. For each task, we demonstrate the benefits of using meta-learned initial weights optimized to reconstruct a specific class of signals.

For 2D image regression, a meta-learned weight initialization leads to faster convergence and better performance during test-time optimization. For CT reconstruction, it allows for better reconstruction quality from fewer supervision views during test-time optimization. For 3D shape reconstruction from images, it allows for faster convergence at test time and makes single view reconstruction possible. For Phototourism landmark reconstruction, it can be optimized at test time to transfer the appearance of a single input image onto the whole landmark, which can then be rendered from novel camera views.

## 4.1. Tasks

Here we provide the basic setup for each task. Please see the supplement for full implementation details.

**Image regression** A prototypical example of a coordinate-based neural representation is an MLP optimized to represent a 2D image [34, 37] by taking in 2D pixel coordinates and outputting RGB color values. We consider four different distributions $\mathcal{T}$: images of faces (*CelebA* [21]), natural images (*Imagenette* [14]), images of text (*Text*), and 2D signed distance fields of simple curves (*SDF*). Each category contains around ten thousand examples. Given a sampled image $T \sim \mathcal{T}$, we provide all $178 \times 178$ pixels as observations for optimizing the network weights $\theta$ in the inner loop. Since this task is not memory constrained, we use MAML to meta-learn the weights over 2 unrolled gradient steps (separately for each category $\mathcal{T}$). In each of these inner loop steps, the entire image is reconstructed and used to calculate the loss. For the MLP $f_\theta$, we use 5 layers with 256 channels each and sine function nonlinearities, as in SIREN [35].

**CT reconstruction** Computed tomography (CT) is a widely used medical imaging technique that captures projective measurements of the volumetric density of a target object. Tancik *et al.* [37] use a coordinate-based neural representation to reconstruct a 2D signal from 1D integral projections; the underlying MLP $f_\theta$ takes in a 2D coordinate and outputs a scalar volume density at that location. Here
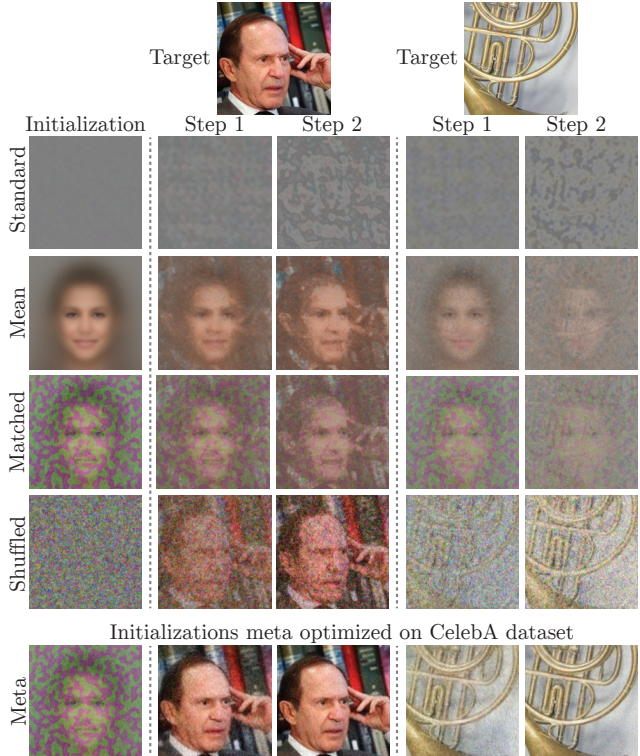


Figure 2. **Faster convergence:** Examples of optimizing a network to represent a 2D image from different initial weight settings. The meta-learned initialization (*Meta*) is specialized for the class of human face images but still helps speed up convergence on other natural images (right). Non-meta-initialized networks take 10 to 20 times as many iterations to reach the same quality as the meta-initialized network does after only 2 gradient steps (see Table 1).

$\mathcal{T}$ is a dataset of 2048 randomly generated $256 \times 256$ pixel Shepp-Logan phantoms [32], where we provide 2D integral projections of a bundle of 256 parallel rays from a random angle as the measurement for each sampled signal $T$ during meta-learning. We use Reptile to meta-learn the initial weights over 12 unrolled gradient steps. We found this to outperform MAML, which was limited to 3 unrolled steps due to memory constraints. For the MLP $f_\theta$, we use 5 layers with 256 channels each and ReLU nonlinearities, and we apply random Fourier features to the input coordinates [37].

**View synthesis for ShapeNet [3] objects** The goal of view synthesis is to generate a novel view of a scene from a set of reference images. Recently, neural radiance fields (NeRF) [25] proposed a method to accomplish this task by using a neural representation that predicts a color and density for any input 3D location and 2D viewing direction within the scene, along with a differentiable volumetric rendering model to generate new views from that representation. This network is optimized to minimize the resid-

| Init. Method | 2 Step PSNR ↑ | # of iters to match ↓ |
|---|---|---|
| Standard | 10.88 | $37.92 \pm 6.31$ |
| Mean | 14.48 | $25.59 \pm 4.57$ |
| Matched | 13.73 | $26.32 \pm 4.17$ |
| Shuffled | 16.29 | $25.80 \pm 4.02$ |
| Meta | **30.37** | - |

Table 1. Comparison of different initialization methods on an image regression task using the CelebA dataset. We report reconstruction PSNR after two steps of test-time optimization. The meta-learned initialization (*Meta*) significantly outperforms all other initializations. We also report the average number of iterations necessary to match the accuracy of *Meta* after two steps.

|  |  | Task | | | |
|---|---|---|---|---|---|
|  |  | CelebA | Imagenette | Text | SDF |
| Init. | CelebA | **30.37** | 26.44 | 21.53 | 36.45 |
|  | Imagenette | 28.51 | **27.07** | 22.63 | 34.80 |
|  | Text | 14.65 | 15.83 | **27.85** | 23.14 |
|  | SDF | 19.80 | 20.05 | 17.23 | **51.73** |

Table 2. PSNR comparison of four different learned initializations for image regression. Each row corresponds to an initialization meta-learned over a different underlying image dataset. The columns indicate which dataset images are sampled from during testing. The best initialization for each task (bolded) is the one specifically optimized on training images drawn from the same dataset. We observe that initializations transfer better between more similar datasets (*CelebA* and *Imagenette*, both natural images) and poorly between less similar datasets (the frequency spectrum of *Text* images is unlike that of the other categories).

ual of re-rendering each of the input reference images from their respective camera poses. In our view synthesis experiments, we use a simplified NeRF model (simple-NeRF) that maintains the same image supervision and volume rendering context. Unlike the original NeRF model, we do not feed in the viewing direction and we use a single model instead of the two "coarse" and "fine" models used by NeRF.

For view synthesis on objects from the ShapeNet [3] dataset, we consider three categories $\mathcal{T}$: *Chairs*, *Cars*, and *Lamps*. We provide 25 $128 \times 128$ pixel reference images during meta-learning for each 3D object $T$. The reference viewpoints are randomly distributed on a sphere and are oriented towards the target object, and each object is oriented in the canonical coordinate frame. The scenes are lit by a randomly selected environment map [9] and rendered using ray tracing. We use Reptile to meta-learn the initial weights (for each shape category) over 32 unrolled gradient steps. For the MLP $f_\theta$, we use 6 layers with 256 channels each and ReLU nonlinearities, and apply a positional encoding to the input coordinates [25].

**View synthesis for Phototourism [16] scenes** This dataset consists of thousands of posed tourist photographs of famous landmarks. Our objective is to use these images to create an underlying representation that can be explored and rendered from novel viewpoints with varying lighting conditions. The primary challenge is the diversity of the capture conditions: the photos are taken with different lighting conditions, camera hardware, camera viewpoint, and varying transient objects like people and cars. Each underlying dataset $\mathcal{T}$ for meta-learning $\theta_0^*$ consists of images of a single landmark (*Trevi*, *Sacre Couer*, or *Brandenburg*); the category is the overall 3D structure of the landmark itself, and the signal is its particular appearance (resulting from the time of day, lighting, weather conditions, etc) within a single photo. If a standard NeRF model is trained directly on this data, it learns a blurry representation of the scene that roughly corresponds to the mean of the environmental conditions. NeRF in the Wild [22] explores these short-

comings and proposes extensive architectural modifications to account for the variations. We find that these shortcomings can be addressed to some degree solely with a better initialization and no architectural changes.

We apply meta-learning to the same simple-NeRF model from the ShapeNet experiment. The meta-training dataset for each landmark consists of thousands of images with varying resolution and intrinsic/extrinsic camera parameters. We use Reptile to meta-learn the initial weights (for each landmark) over 64 unrolled gradient steps. At test time, we optimize the simple-NeRF (starting from the initial weights $\theta_0^*$ for that landmark) to reproduce the appearance of a new image, and then render that simple-NeRF from other viewpoints. For the underlying MLP $f_\theta$, we use 6 layers with 256 channels each and ReLU nonlinearities, and apply positional encoding to the input coordinates [25].

### 4.2. Baselines

As well as a *Standard* randomly initialized network (Glorot *et al*. [12]), we compare to various other initialization schemes in several of our experimental settings:
- *Mean:* we optimize a network from scratch such that its output matches the mean signal $E_{T \sim \mathcal{T}}[T]$ from the current class $\mathcal{T}$.
- *Matched:* we optimize a network from scratch such that its output matches the output of a network using the meta-learned initialization for the current class $\mathcal{T}$.
- *Shuffled:* we randomly permute the weights (within each network layer) of the meta-learned initialization $\theta_0^*$ for the current class $\mathcal{T}$.

Both the *Mean* and *Matched* baselines demonstrate the difference between having a good initialization in *signal* space versus *weight* space—despite *Mean* and *Matched* being initialized so that the loss against a randomly sampled signal will be low, they are a worse starting point for gradient descent than the actual meta-learned initial weights. The *Shuf-

| Init. | PSNR | | | |
|---|---|---|---|---|
| Method | 1 Views | 2 Views | 4 Views | 8 View |
| Standard | 13.63 | 14.15 | 16.31 | 21.49 |
| Mean | 14.72 | 15.39 | 17.43 | 25.19 |
| Matched | 14.07 | 15.51 | 20.25 | 24.77 |
| Shuffled | 13.64 | 14.17 | 16.69 | 22.09 |
| Meta | **15.09** | **18.70** | **22.00** | **27.34** |

Table 3. Comparison of initialization methods on a CT reconstruction task. Each "view" consists of 256 parallel rays. The data-dependent prior acquired during meta-learning improves reconstruction quality when fewer views are observed.

*fled* baseline demonstrates that matching the statistical distribution of the meta-learned initial weights is not sufficient for better convergence or generalization. We find that using the Adam [17] optimizer performs best for all of the baseline initializations, but that standard stochastic gradient descent works best for the meta-learned initializations (we choose the best optimizer and hyperparameters for each task and initialization using a held-out validation set, see supplement for details).

### 4.3. Faster convergence

**Image regression** In Figure 2, we visualize the network output for a variety of initial weight settings, showing the output images after 0, 1, and 2 gradient steps of test-time optimization. The meta-learned initial weights are optimized to represent face images (CelebA [21]). When using the learned initial weights $\theta_0^*$ (*Meta*), the target image is already clearly visible after the very first step. In contrast, the baseline initialization methods take an order of magnitude more iterations to represent the target image to the same accuracy (see Table 1). The *Mean*, *Matched*, and *Shuffled* baselines perform better than the completely random *Standard* initialization, but still take over ten times as many iterations to reach the same quality as the meta-initialized network can after 2 steps. In particular, this demonstrates that neither matching the image space output nor the statistical distribution of the meta-learned weights is sufficient for achieving a similar speedup.

**View synthesis for ShapeNet [3] objects** In Figure 5, we plot the image reconstruction accuracy for a held-out test set of objects from the *Chair* category. During test-time optimization, 25 views are observed. We find that starting from the optimized weights $\theta_0^*$ allows the network to recover the chair more quickly compared to the *Standard* weight initialization. We note that after many steps, both methods end up at a similar quality.
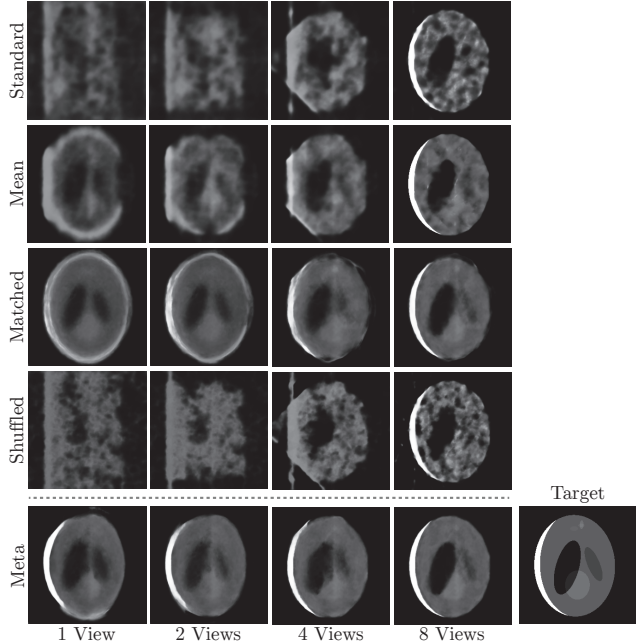


Figure 3. **Sparse Recovery:** Examples of CT reconstructions of a Shepp-Logan phantom from a sparse set of views. The meta-learned initial weights encode a data-dependent prior that improves reconstruction in the limited data regime.

### 4.4. Generalizing from partial observations

**Image regression within a category** We perform meta-learning experiments across multiple datasets to determine the extent that the optimized weight initialization acts as a class-specific prior. We compare initializations trained on four different image datasets (*CelebA*, *Imagenette*, *Text*, and *SDF*). Table 2 presents a confusion matrix demonstrating that optimizing the network initialization does in fact induce a dataset-dependent prior, with each learned initialization generalizing best to the same dataset distribution it was trained on.

**CT reconstruction from sparse views** We report the reconstruction quality over a test set of phantoms given varying numbers of views at test time in Table 3 and visualize one test example in Figure 3. We observe poor reconstructions from the *Standard* initialization when few views are provided. The meta-learned initializations are consistently able to match the PSNR of *Standard* with half as many views. The *Mean* initialization is generated by training a network to reconstruct the mean of the training phantoms. It is better able to preserve the structure of the phantom compared to *Standard* but still performs worse than the meta-learned initializations.
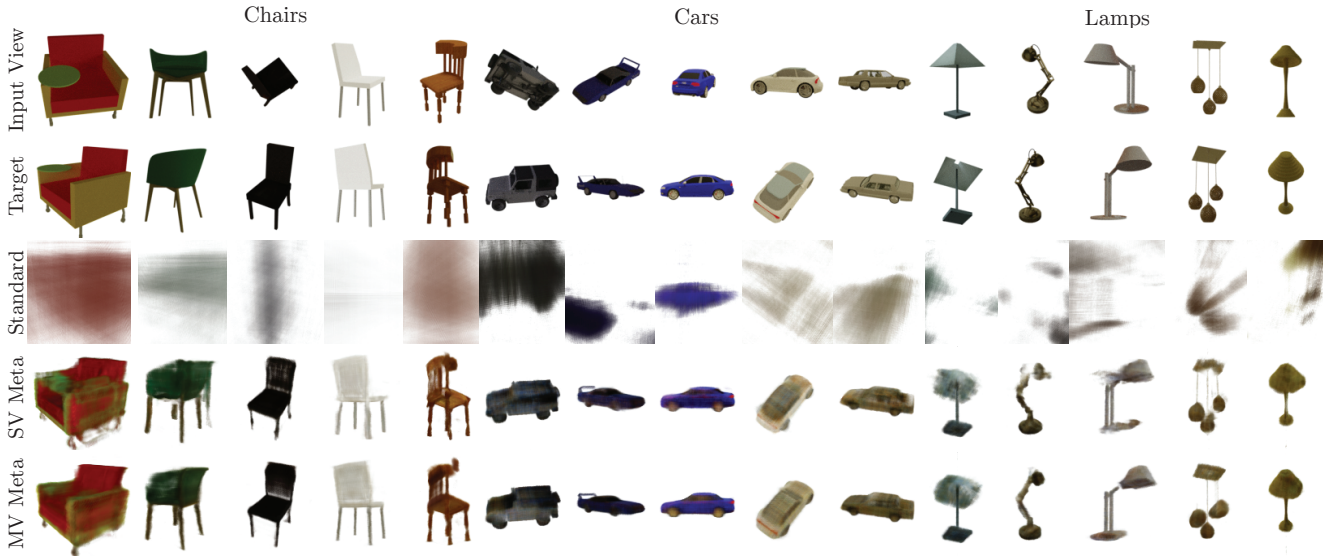
Figure 4. **Single view reconstructions of ShapeNet [3] objects.** The simple-NeRF formulation relies on multi-view consistency for supervision and therefore fails if naively applied to the task of single view reconstruction, as seen in the *Standard* column. However, if the model is trained starting from meta-learned initial weights, it is able to recover 3D geometry. The *MV Meta* initialization has access to multiple views per object during meta-learning, whereas the *SV Meta* initialization only has access to a single view per object during meta-learning. All methods only receive a single input view during test-time optimization.



|  | PSNR | | |
|---|---|---|---|
|  | Chairs | Cars | Lamps |
| Standard | 12.49 | 11.45 | 15.47 |
| MV Matched | 16.40 | 22.39 | 20.79 |
| MV Shuffled | 10.76 | 11.30 | 13.88 |
| MV Meta | **18.85** | **22.80** | **22.35** |
| SV Meta | 16.54 | 22.10 | 20.95 |

Table 4. Metrics for single image ShapeNet reconstructions using a simple-NeRF model. See Figure 4 for image examples and §4.4 for experimental details.
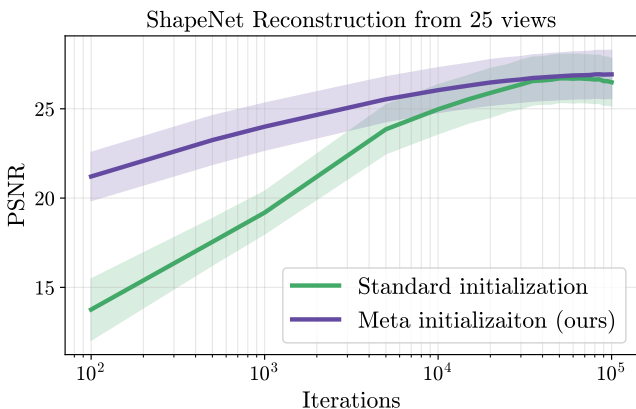
Figure 5. Reconstruction quality over the course of training for models optimized to reconstruct ShapeNet chairs from a set of 25 reference images. The model starting from the meta-learned initial weights outperforms the network using a standard random initialization throughout training.

**Single image view synthesis for ShapeNet [3]** A simple-NeRF model with a *Standard* random initialization relies on multi-view consistency to reconstruct the appearance of a 3D object. With only a single view, this naïve model is unable to recover any meaningful shape. We find that a learned initialization "bakes in" a class-specific shape prior that enables the recovery of 3D geometry (Figure 4, Table 4). We can meta-learn an effective weight initialization for single-view reconstruction by optimizing over a dataset with 25 training views of each object (*MV Meta*). We find that this prior persists even if the meta-training dataset only contains a single reference image per scene (*SV Meta*), meaning that the meta-learning phase has no access to multiview information for any particular object.

**View synthesis with appearance transfer for Photo-tourism [16]** As described in §4.1, these images have different camera poses and visual appearance (lighting, sky, etc.) as they are taken by tourists at different times. Our goal at test time is to explore the landmark from varying camera viewpoints but rendered with the same appearance as in a target photograph. In every step of the meta-learning outer loop, we supervise the simple-NeRF model to match the appearance of a random photo of the landmark (with varying pose and appearance). We find that performing test-

Figure 6. Reconstructions of the Trevi Fountain and Sacre Coeur landmarks from the Phototourism dataset [16]. The meta-learning algorithm is run over tourist images taken at different locations and times. During the test-time optimization, the neural representation is trained to recover the input view on the left. The strong prior from the initialization captures the underlying geometry, allowing us to render views from the camera positions of the images in the top row while retaining the appearance of the input view.

|  | PSNR | | |
|---|---|---|---|
|  | Trevi | Sacre Coeur | Brandenburg |
| Basic NeRF | 17.14 | 17.59 | 17.77 |
| Meta | **19.35** | **19.33** | **19.11** |

Table 5. Reconstruction results on Phototourism data. Multi-view data with consistent appearance is not available in this dataset, so we optimize on one half of an image and report image metrics on the other half. We compare our Reptile setup (*Meta*) with a standard NeRF network trained on all images of the landmark and then test-time optimized to fit each held-out target image. This is equivalent to training Reptile with one inner loop gradient step.

time optimization using a single new photograph allows us to render convincing unobserved viewpoints of the scene with the same environmental conditions.

In Figure 6, we show results for two landmarks. We test-time optimize the meta-learned weights for five target images (shown on the left side of the grid), taking 150 gradient steps for each image. We then render each of the resulting simple-NeRF networks from the five different viewpoints (shown in the row above the grid). The result is an image from the camera position of the corresponding top row image and matching the appearance of the left column image.

Quantitative evaluation on the Phototourism dataset is difficult as multiple views with the same environmental conditions do not exist. To overcome this, for Table 5 we optimize and evaluate on the same image, by optimizing to match the appearance of the left half of the image and subsequently evaluating metrics on the right half. For comparison, we train a simple-NeRF model with a standard random initialization from scratch on each landmark, then test-time optimize it to match the left half of each new view before evaluating it on the right half. This is algorithmically equivalent to Reptile with one inner optimization step. We find that unrolling Reptile for 64 inner steps performs better, producing significantly clearer renderings of the landmark.

## 5. Conclusion

Our results show that simply modifying a coordinate-based neural representation's initial weight values can guide the network along a significantly better optimization trajectory, without changing the underlying architecture or test-time optimization procedure. These meta-learned initial weights can result in faster convergence or act as a strong prior for representing signals from a given distribution. This partially ameliorates a major shortcoming of neural representations (separately optimizing a network for each new signal) without limiting their representational power.

There are many additional directions to explore, such as applying more sophisticated meta-learning algorithms or more precisely characterizing the geometry of weight space for these networks. One limitation of our current approach is that it requires a sizable dataset of example signals from a target distribution in order to derive beneficial initial weights. Another shortcoming is that our method still requires some amount of test-time optimization.

As the number of use cases for neural representations continues to rapidly expand, we believe this work takes an important step toward understanding the importance of their initial weights and optimization behavior.

# 6. Acknowledgements

# References

[1] Antreas Antoniou, Harrison Edwards, and Amos Storkey. How to train your maml. *ICLR*, 2018. 2

[2] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, and Skye Wanderman-Milne. JAX: composable transformations of Python+NumPy programs, 2018. 10

[3] Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. Shapenet: An information-rich 3d model repository. Technical report, 2015. 2, 4, 5, 6, 7, 11

[4] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. *CVPR*, 2019. 2

[5] Boyang Deng, JP Lewis, Timothy Jeruzalski, Gerard Pons-Moll, Geoffrey Hinton, Mohammad Norouzi, and Andrea Tagliasacchi. Neural articulated shape approximation. *ECCV*, 2020. 2

[6] Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. On the convergence theory of gradient-based model-agnostic meta-learning algorithms. In *AISTATS*, 2020. 2

[7] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. *ICML*, 2017. 2, 3, 10

[8] Sebastian Flennerhag, Andrei A Rusu, Razvan Pascanu, Francesco Visin, Hujun Yin, and Raia Hadsell. Meta-learning with warped gradient descent. *ICLR*, 2020. 2

[9] Marc-André Gardner, Kalyan Sunkavalli, Ersin Yumer, Xiaohui Shen, Emiliano Gambaretto, Christian Gagné, and Jean-François Lalonde. Learning to predict indoor illumination from a single image. *arXiv preprint arXiv:1704.00090*, 2017. 5

[10] Kyle Genova, Forrester Cole, Aaron Sarna Daniel Vlasic, William T. Freeman, and Thomas Funkhouser. Learning shape templates with structured implicit functions. *ICCV*, 2019. 2

[11] Kyle Genova, Forrester Cole, Avneesh Sud, Aaron Sarna, and Thomas Funkhouser. Local deep implicit functions for 3D shape. *CVPR*, 2020. 2

[12] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. *AISTATS*, 2010. 5

[13] Timothy Hospedales, Antreas Antoniou, Paul Micaelli, and Amos Storkey. Meta-learning in neural networks: A survey. *arXiv preprint arXiv:2004.05439*, 2020. 2

[14] Jeremy Howard. imagenette. 4

[15] Chiyu Jiang, Avneesh Sud, Ameesh Makadia, Jingwei Huang, Matthias Nießner, and Thomas Funkhouser. Local implicit grid representations for 3D scenes. *CVPR*, 2020. 2

[16] Yuhe Jin, Dmytro Mishkin, Anastasiia Mishchuk, Jiri Matas, Pascal Fua, Kwang Moo Yi, and Eduard Trulls. Image matching across wide baselines: From paper to practice. *International Journal of Computer Vision*, pages 1–31, 2020. 5, 7, 8, 11, 13

[17] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *ICLR*, 2015. 2, 3, 6, 10

[18] Zhenguo Li, Fengwei Zhou, Fei Chen, and Hang Li. Meta-sgd: Learning to learn quickly for few-shot learning. *arXiv preprint arXiv:1707.09835*, 2017. 2

[19] Shichen Liu, Shunsuke Saito, Weikai Chen, and Hao Li. Learning to infer implicit surfaces without 3D supervision. *NeurIPS*, 2019. 2

[20] Shaohui Liu, Yinda Zhang, Songyou Peng, Boxin Shi, Marc Pollefeys, and Zhaopeng Cui. Dist: Rendering deep implicit signed distance function with differentiable sphere tracing. *CVPR*, 2020. 2

[21] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. *ICCV*, 2015. 2, 4, 6

[22] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. *CVPR*, 2021. 5

[23] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. *CVPR*, 2019. 1, 2

[24] Mateusz Michalkiewicz, Jhony K Pontes, Dominic Jack, Mahsa Baktashmotlagh, and Anders Eriksson. Implicit surface representations as layers in neural networks. *ICCV*, 2019. 2

[25] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. NeRF: Representing scenes as neural radiance fields for view synthesis. *ECCV*, 2020. 1, 2, 4, 5, 11

[26] Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms. *arXiv preprint arXiv:1803.02999*, 2018. 2, 3, 10

[27] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. *CVPR*, 2020. 2

[28] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. *CVPR*, 2019. 1, 2

[29] Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. *ICML*, 2019. 2

[30] Aravind Rajeswaran, Chelsea Finn, Sham M Kakade, and Sergey Levine. Meta-learning with implicit gradients. *NeurIPS*, 2019. 2

[31] S. Ravi and H. Larochelle. Optimization as a model for few-shot learning. *ICLR*, 2017. 2

[32] Lawrence A. Shepp and Benjamin F. Logan. The Fourier reconstruction of a head section. *IEEE Transactions on nuclear science*, 1974. 4

[33] Vincent Sitzmann, Eric R. Chan, Richard Tucker, Noah Snavely, and Gordon Wetzstein. MetaSDF: Meta-learning signed distance functions. *NeurIPS*, 2020. 1, 2

[34] Vincent Sitzmann, Julien N.P. Martel, Alexander W. Bergman, David B. Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. *NeurIPS*, 2020. 1, 2, 4, 10

[35] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. *NeurIPS*, 2019. 1, 2, 4

[36] Kenneth O Stanley. Compositional pattern producing networks: A novel abstraction of development. *Genetic programming and evolvable machines*, 8(2):131–162, 2007. 2

[37] Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T. Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *NeurIPS*, 2020. 2, 4, 10

[38] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Ronen Basri, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. *NeurIPS*, 2020. 2

## A. Implementation details

We found that modifying the weight initialization for these coordinate-based networks drastically changed their convergence behavior during test-time optimization. As a result, we tuned the optimization method and hyperparameters for each part of each experiment (using held-out validation sets) in order to provide the fairest possible comparison and to not bias the results against the non-meta-learned initializations. For example, we often found that SGD outperformed Adam when doing test-time optimization using meta-learned initializations, but that Adam was significantly better than SGD with a standard random initialization.

All experiments are implemented in JAX [2]. Each experiment is trained on either a single NVIDIA V100, 2080 Ti, or 3080 Ti. In all cases where the Adam optimizer [17] is used, we keep the standard parameter choices for $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$.

### A.1. Image regression

For this task we use a SIREN [34] architecture ($\omega_0 = 200$) with 5 layers of 256 channels each. For the randomly initialized *Standard* baseline, we use the specific initialization procedure as proposed in the SIREN paper.

MAML [7] is trained for 150K iterations. Each iteration has an outer batch size of 3 target images. The inner batch

| Init. Method | 2 Step PSNR ↑ | # of iters to match ↓ |
|---|---|---|
| Standard | 10.88 | 37.92 ± 6.31 |
| Mean | 14.48 | 25.59 ± 4.57 |
| Matched | 13.73 | 26.32 ± 4.17 |
| Shuffled | 16.29 | 25.80 ± 4.02 |
| Reptile | 25.55 | **9.86 ± 7.42** |
| MAML | **30.37** | - |

Table 6. Image reconstruction results with meta-learning results for both MAML and Reptile. MAML performs best, but Reptile also outperforms the non-meta-learned weights.

contains all pixels of the target image. The outer loop uses the Adam optimizer with learning rate of $10^{-5}$. The inner loop performs two steps of gradient descent with a learning rate of $10^{-2}$.

We additionally meta-learn another initialization using Reptile [26]. We use the same learning rates as in MAML but with an outer batch size of 10 target images. We report the Reptile reconstruction accuracy in Table 6. We note that Reptile also outperforms the non-meta-learned weights.

During test-time optimization, we use gradient descent with learning rate of $10^{-2}$ when starting from the MAML initial weights. For the baseline methods (*Standard*, *Mean*, *Matched*, *Shuffled*) we used Adam with learning rate of $10^{-4}$, which performed significantly better than than gradient descent.

### A.2. CT reconstruction

For this task we use an MLP with 5 layers of 256 channels each. The network uses a ReLU activation after each layer with the exception of the last layer, which has a sigmoid activation. Prior to inputting the coordinates into the network, we encode them using random Fourier features sampled from a normal distribution with $\sigma = 30$, as was done in Tancik *et al.* [37].

Reptile [26] is trained for 100K iterations. Each iteration has an outer batch size of 1. The inner batch contains 20 CT projections, each with 256 measurements, taken from a randomly sampled direction. The outer loop uses the Adam optimizer with learning rate of $5 \times 10^{-5}$. The inner loop performs 12 inner loop steps of gradient descent with a learning rate of $10^1$.

We perform test-time optimization experiments with different numbers of supervision views to compare reconstruction quality. We found that the models are more prone to overfitting when fewer views are provided. We tune the learning rate and number of gradient steps for each initialization method according to a held-out set of 16 validation images. We report all of the test-time optimization hyperparameters in Table 7.

### A.3. ShapeNet [3] view synthesis

We use a simplified NeRF [25] model for our view synthesis tasks. This model uses a single network rather than two networks (coarse and fine), and we do not provide view directions as input. The network is an MLP with 6 layers, each with 256 channels and ReLU activations. As in NeRF [25], we apply a positional encoding to each input coordinate with the form

$$\bigcup_{i=0}^{N} \left\{ \cos\left( 2^{fi/N} x \right), \sin\left( 2^{fi/N} x \right) \right\}, \qquad (7)$$

with $N = 20$ encodings and log-max frequency $f = 8$. We accumulate 128 samples per ray for rendering.

Reptile is trained for 100K iterations with an outer batch size of 1. The inner loop step optimizes over a batch of 128 rays. We perform 32 inner loop steps for every outer loop step. The outer loop uses the Adam optimizer with learning rate $5 \times 10^{-4}$ for the *Chairs* scenes and $5 \times 10^{-5}$ for the *Lamps* and *Cars* scenes.

The test-time optimization parameters vary depending on the scene and the number of views available during meta-learning. Each experiment uses an inner batch of 64 rays. The *Shuffled* and *Matched* initializations are computed based on the *MV Meta* weights. For the 25 view chair reconstruction, we use stochastic gradient descent with a learning rate of $10^{-1}$ for the Reptile initialization; for the standard initialization, we use Adam with a learning rate of $10^{-4}$. The test-time hyper parameters for the single view experiments are listed in Table 8.

### A.4. Phototourism [16] view synthesis

We use the same architecture as described in §A.3. Reptile is trained for 150K iterations with an outer batch size of 1. The inner loop step optimizes over a batch of 64 rays, with 128 volume rendering samples per ray. The outer loop uses the Adam optimizer with a learning rate of $5^{-4}$. We train with 64 inner loop steps using gradient descent with a learning rate of 10. We compare to *Basic NeRF* which has the same setup, but only one inner step. For *Basic NeRF* we train *Trevi* for 60K iterations, *Brandenburg* for 100K iterations, and *Sacre Coeur* for 200K iterations. To transfer the appearance of a new photo during test-time optimization, we take 150 gradient steps with a learning rate of 10.

## B. Weight space interpolation

We find that linearly interpolating between networks in weight space produces meaningful outputs when using meta-learned weights. Figure 7 shows interpolation between networks trained to represent images, and Figure 8 shows interpolations between networks that are trained to reconstruct a Phototourism landmark.

|  | Opt. method | | Number of steps | | | |
|---|---|---|---|---|---|---|
|  | Adam | LR | 1 View | 2 Views | 4 Views | 8 Views |
| Standard | ✓ | $5 \times 10^{-4}$ | 50 | 100 | 250 | 1000 |
| Mean | ✓ | $10^{-5}$ | 25 | 50 | 100 | 1000 |
| Matched | ✓ | $5 \times 10^{-4}$ | 50 | 100 | 500 | 1000 |
| Shuffled | ✓ | $5 \times 10^{-4}$ | 50 | 100 | 250 | 1000 |
| Meta | ✗ | $10^{1}$ | 50 | 100 | 1000 | 1000 |

Table 7. Hyper-parameters for CT test-time optimization. Each value is tuned on a held-out validation set.

|  | Chairs | | | Cars | | | Lamps | | |
|---|---|---|---|---|---|---|---|---|---|
|  | Adam | LR | Steps | Adam | LR | Steps | Adam | LR | Steps |
| Standard | ✓ | $10^{-5}$ | 1000 | ✓ | $5 \times 10^{-5}$ | 2000 | ✓ | $5 \times 10^{-5}$ | 2000 |
| Matched | ✓ | $10^{-4}$ | 2000 | ✓ | $5 \times 10^{-5}$ | 2000 | ✓ | $5 \times 10^{-5}$ | 2000 |
| Shuffled | ✓ | $10^{-4}$ | 2000 | ✓ | $5 \times 10^{-5}$ | 2000 | ✓ | $5 \times 10^{-5}$ | 2000 |
| MV Meta | ✗ | $10^{-1}$ | 1000 | ✗ | $5 \times 10^{-1}$ | 2000 | ✗ | $5 \times 10^{-1}$ | 2000 |
| SV Meta | ✗ | $5 \times 10^{-1}$ | 1000 | ✗ | $5 \times 10^{-1}$ | 2000 | ✗ | $5 \times 10^{-1}$ | 2000 |

Table 8. Hyper-parameters for ShapeNet test-time optimization from a single view. Each value is tuned on a held-out validation set.



Figure 7. Weight space interpolation for networks optimized to represent 2D images. We use test-time optimization to fit network weights for three different images (denoted with arrows), then linearly interpolate between those weight values and visualize the resulting outputs. When test-time optimization begins from a standard random initialization (*Standard*, top), weight space interpolation produces displeasing artifacts, but when it begins from a meta-learned initialization (*Meta*, bottom) the resulting outputs maintain an image-like appearance.
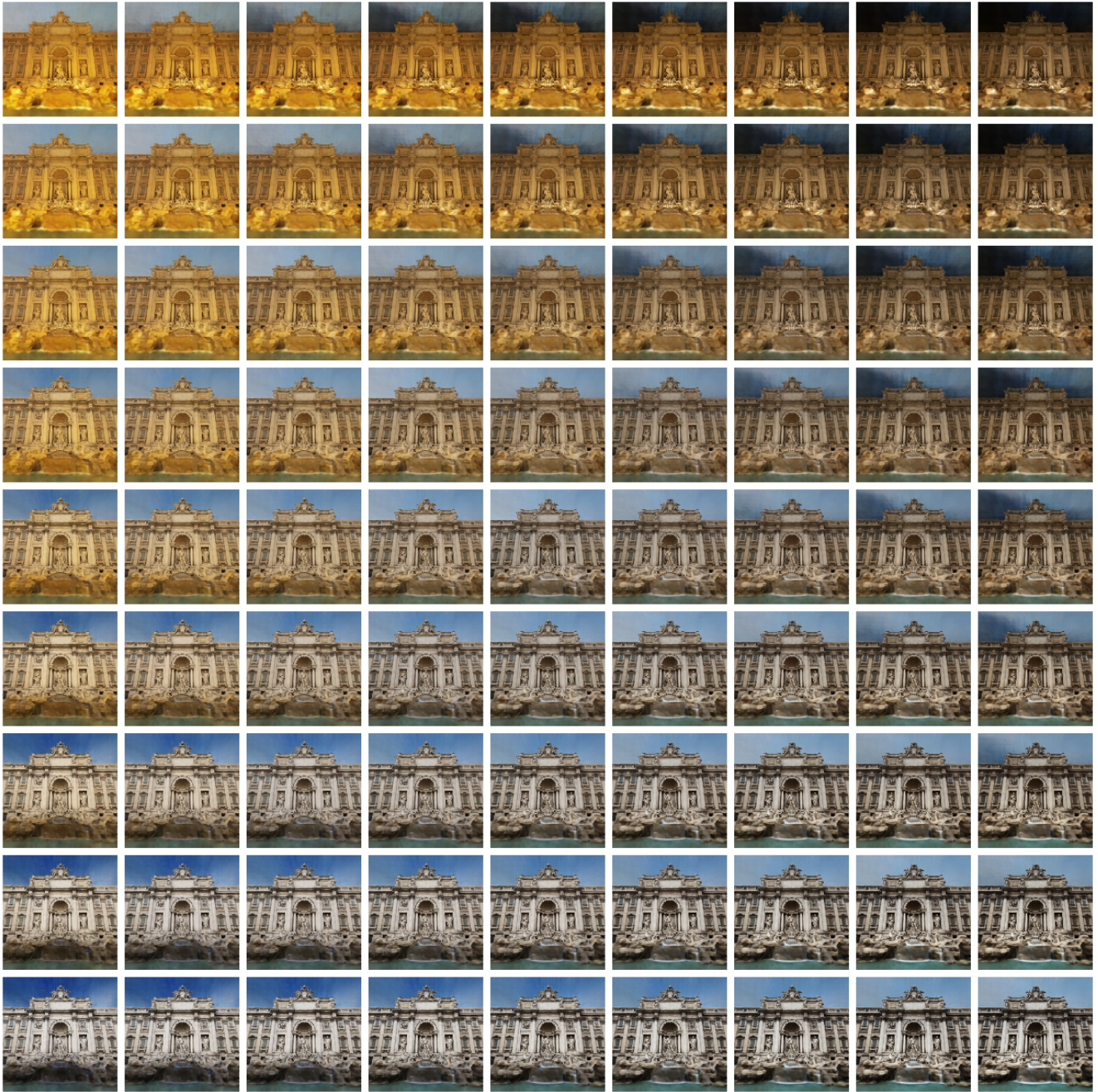
Figure 8. Appearance interpolation on the Trevi Fountain scene from the Phototourism dataset [16]. We render the scene from a single fixed camera pose. Each corner of the grid represents a NeRF network that has been test-time optimized to match the appearance of a single image (starting from the meta-learned initial weights $\theta_0^*$). We linearly interpolate between these networks in weight space to render the grid of images shown here. The image in the center is produced by directly rendering the meta-learned initial weights (with no test-time optimziation), representing an "average" appearance for the scene. Please see the supplemental video for an animated version of this figure with a moving camera path.