

low-pass filter

BACON: Band-limited Coordinate Networks for Multiscale Scene Representation

David B. Lindell Dave Van Veen Jeong Joon Park Gordon Wetzstein
Stanford University

<http://computationalimaging.org/publications/bacon>

Abstract

Coordinate-based networks have emerged as a powerful tool for 3D representation and scene reconstruction. These networks are trained to map continuous input coordinates to the value of a signal at each point. Still, current architectures are black boxes: their spectral characteristics cannot be easily analyzed, and their behavior at unsupervised points is difficult to predict. Moreover, these networks are typically trained to represent a signal at *a single scale*, and so naive downsampling or upsampling results in artifacts. We introduce band-limited coordinate networks (BACON), a network architecture with an analytical Fourier spectrum. BACON has predictable behavior at unsupervised points, can be designed based on the *spectral characteristics of the represented signal*, and can represent signals at multiple scales without explicit supervision. We demonstrate BACON for multiscale neural representation of images, radiance fields, and 3D scenes using *signed distance functions* and show that it outperforms conventional single-scale coordinate networks in terms of interpretability and quality.

1. Introduction

Coordinate networks are an emerging class of neural networks that can be used to represent or optimize a broad format of signals including images, video, 3D models, audio waveforms, and more [38, 41, 48, 60, 62]. As opposed to storing discrete samples of signals in conventional array- or grid-based formats, neural representations approximate signals using a continuous function that is embedded in the learned weights of a fully-connected neural network. *Given an input coordinate, these networks are trained to output the value of a signal at that point.* Since even complex or high-dimensional signals can be flexibly optimized using a coordinate network, they have become popular for applications including view synthesis [41], image processing [56], 3D reconstruction [48], and neural rendering [66].

Yet, current coordinate networks are black box models that are designed to represent signals at a single scale. As a

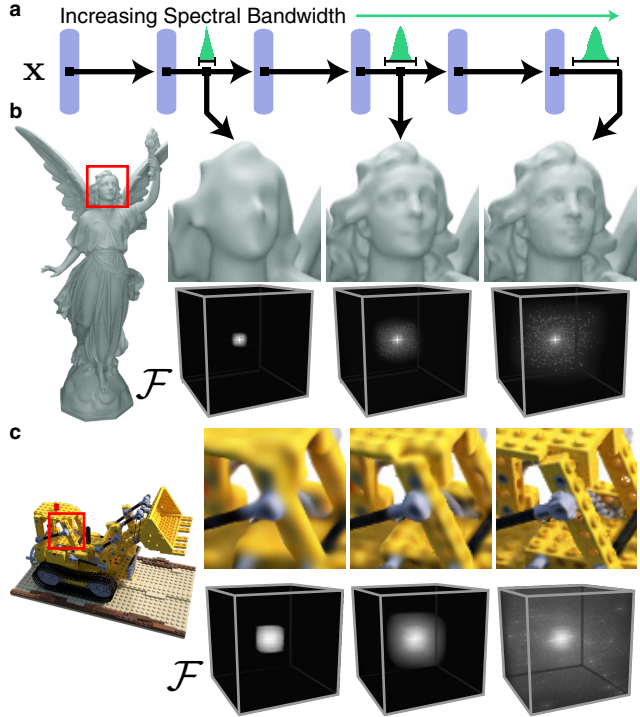


Figure 1. Overview of band-limited coordinate networks (BACON). (a) The proposed architecture produces intermediate outputs with an analytical spectral bandwidth that can be specified at initialization. When supervised on a high-resolution signal, the network learns a multi-resolution decomposition of the output, as shown for fitting 3D shapes via a signed distance function (b) and radiance fields (c). The network is characterized entirely by its Fourier spectrum (see insets) and so its behavior is predictable, even at unsupervised locations.

result, the behavior of the network at unsupervised coordinates is difficult to predict, with complex dependencies on hyperparameters such as hidden layer size, network depth, or input coordinate encoding. The black box nature of the architecture similarly inhibits multiscale signal representation, since we cannot readily filter or anti-alias these models, and the frequency spectrum of a coordinate network is

current coordinate network is only working well
for spatial spectrum, but not for frequency spectrum

difficult to analyze. Thus naive downsampling or upsampling by querying the network on a coarser or finer grid of coordinates leads to aliasing or undesired high-frequency artifacts. Ultimately, these characteristics stem from the fact that coordinate networks are not amenable to Fourier analysis and are not designed to be scale aware.

Still, being able to represent and optimize signals at multiple resolutions is an important requirement for many applications. For example in image processing, many techniques rely on image pyramids [58] (e.g., optical flow estimation, compression, filtering, etc.). Representing 3D objects or scenes at multiple levels of detail is useful for speeding up rendering and reducing memory requirements (e.g., mipmapping).

In this work, we introduce band-limited coordinate networks (BACON). The key properties of this architecture are that (1) the maximum frequency at each layer can be manipulated analytically, and (2) the behavior of a trained network is entirely characterized by its Fourier spectrum. BACON is suited to multiscale signal representation because band-limited output layers can be designed with an inductive bias towards a particular resolution or scale.

In addition to introducing BACON, we demonstrate a variety of applications including multiscale representation of images, neural radiance fields, and 3D scenes. Our work takes important steps towards making coordinate-based networks scale aware, and provides a new representation with interpretable behavior. Specifically, we make the following contributions:

- We introduce band-limited coordinate-based networks for representing and optimizing signals.
- We develop methods for spectral analysis of the architecture, and propose a principled, band-limited initialization scheme.
- We demonstrate that our architecture outperforms conventional single-scale coordinate networks for multiscale image fitting, neural rendering, and 3D scene representation.

2. Related Work

Neural Scene Representation and Rendering. Emerging neural scene representations promise 3D-structure-aware, continuous, memory-efficient representations for parts [18, 19], objects [3, 6, 12, 20, 39, 48, 72], or scenes [13, 24, 51, 60, 62]. These can be supervised with 3D data, such as point clouds, and optimized as either signed distance functions [3, 20, 24, 28, 39, 48, 51, 59, 62, 64, 73] or occupancy networks [9, 38]. Using neural rendering [66], representation networks can also be trained using multiview 2D images [4, 17, 25, 31–33, 36, 41, 42, 45–47, 52, 56, 62, 63, 69, 71, 72, 77, 78]. Temporally aware extensions [44]

and multimodal variants with part-level semantic segmentation [30] have also been proposed. Recent 3D-aware GANs use related ideas but are training with 2D image collections [7, 43, 57].

Architectures for Scene Representation. Neural network architectures for scene representation networks can be roughly classified as feature-based, coordinate-based, or hybrid. Feature-based approaches represent the scene using differentiable feature primitives, such as points [14, 50, 53, 70, 75], surface patches [74], meshes [21, 55, 67, 79], multi-plane [16, 40, 80] or multi-sphere [2, 5] images, or using a voxel grid of features [34, 61]. A tradeoff with feature-based representations is that they can be quickly evaluated, but typically have a large memory footprint.

Coordinate-based representations (sometimes called implicit representations or coordinate networks), use a multilayer perceptron (MLP) to map from an input coordinate to a signal value, for example, the signed distance or occupancy of a 3D scene. These networks can represent complex, high-dimensional signals with a small memory footprint and represent signals continuously across the input domain. Ultimately, the maximum resolution of coordinate networks is tied to the capacity of the network architecture through dependencies on the number of trainable parameters and other architectural choices. There are varieties of coordinate networks that represent signals globally [38, 48, 60, 65] or locally [6, 8, 24, 54]. Hybrid architectures that combine feature-based and coordinate representations aim to combine the best of both worlds [22, 32, 35, 51].

The proposed method is also a coordinate network, but rather than using an MLP architecture, as with all coordinate networks discussed above, our method builds on recently proposed multiplicative filter networks (MFNs) [15]. We develop the theory of MFNs significantly, with new tools to describe and manipulate the Fourier spectra of these networks, and a new initialization scheme that mitigates vanishing activations in deep networks. These insights enable band-limited coordinate networks, which we demonstrate for multiscale signal representation.

Multiscale Representations. Several existing works have explored multiscale architectures in the context of scene representation networks. For example, proposed methods use an octree [64, 76] to accelerate neural rendering of radiance fields or signed distance functions, or a hierarchy of features [10] to improve 3D shape completion. Multiscale representations can be optimized directly using specialized architectures [35, 73] or progressive training strategies [23, 35]. The closest work to ours in this category is Mip-NeRF [4], which is a coordinate-based network with a scale-dependent positional encoding. After training the network with supervision at multiple scales, the resolution of the network output can be controlled by adjusting the positional encoding. Our work differs in that the bandwidth of

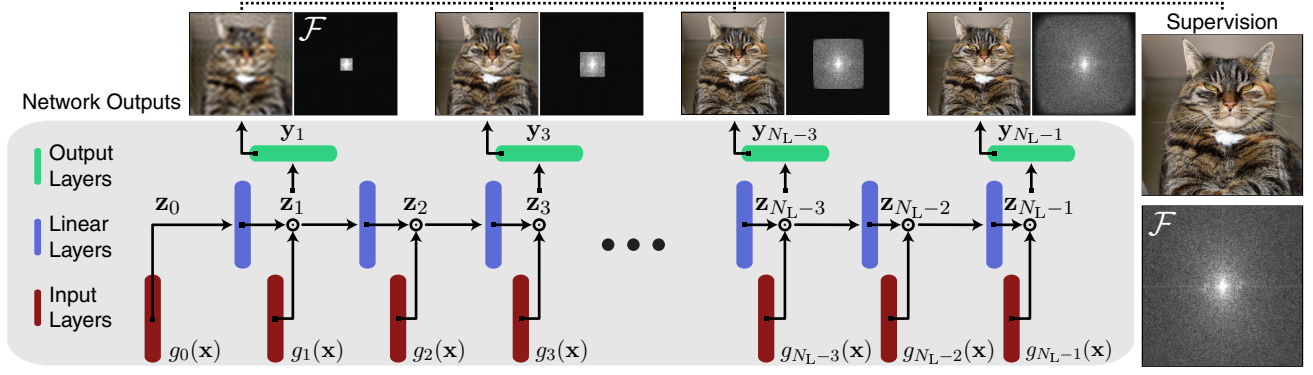


Figure 2. Overview of BACON architecture. We initialize the frequencies of the sine layers of a multiplicative filter network [15] within a limited bandwidth $[-B_i, B_i]$ (bottom row). Then, the bandwidth of each output layers is the sum of the input bandwidths up to that point (top row), allowing the network bandwidth to be explicitly specified. At training time the network can be supervised with a signal at any resolution, and the network learns to fit the signal in a band-limited fashion. Image from DIV2K dataset [1].

the network outputs are constrained by design rather than through training. Thus, our approach learns a band-limited multiscale decomposition of a signal, even without explicit training at multiple scales.

3. Method

This section, provides an overview of MFNs and the BACON multiscale architecture, describes the Fourier spectra of these networks, and proposes an initialization scheme for deep networks.

3.1. Band-limited Coordinate Networks

Our approach builds on a recently introduced coordinate-based architecture called **Multiplicative Filter Networks** (MFNs) [15], which differ from conventional MLPs in that they employ a Hadamard product between linear layers and sine activation functions. While BACON uses an MFN backbone, we significantly extend the theoretical understanding and practicality of these networks by (1) proposing architectural changes to achieve multiscale, band-limited outputs, (2) deriving formulas to quantify the **expected frequencies** in the representation, and (3) deriving a principled initialization scheme that prevents vanishing activations in deep networks.

In a forward pass through the network, an input coordinate $\mathbf{x} \in \mathbb{R}^{d_{in}}$ is first passed through several layers of the form $g_i : \mathbb{R}^{d_{in}} \mapsto \mathbb{R}^{d_h}$, with $g_i(\mathbf{x}) = \sin(\omega_i \mathbf{x} + \phi_i)$, $i = 0, \dots, N_L - 1$, and N_L the number of layers in the network. We refer to the intermediate activations as $\mathbf{z}_i \in \mathbb{R}^{d_h}$, and we allow intermediate outputs of the network $\mathbf{y}_i \in \mathbb{R}^{d_{out}}$ at the i th layer, defined as follows (see also Fig. 2).

$$\begin{aligned} \mathbf{z}_0 &= g_0(\mathbf{x}) \\ \mathbf{z}_i &= g_i(\mathbf{x}) \circ (\mathbf{W}_i \mathbf{z}_{i-1} + \mathbf{b}_i), \quad 0 \leq i < N_L \\ \mathbf{y}_i &= \mathbf{W}_i^{\text{out}} \mathbf{z}_i + \mathbf{b}_i^{\text{out}}, \end{aligned} \quad (1)$$

$$\begin{aligned} \mathbf{z}_1 &= g_1(\mathbf{x}) \circ (\mathbf{W}_1 \mathbf{z}_0 + \mathbf{b}_1) \\ \mathbf{y}_1 &= \mathbf{W}_1^{\text{out}} \mathbf{z}_1 + \mathbf{b}_1^{\text{out}} \end{aligned}$$

where \circ indicates the Hadamard product. The parameters of the network are $\theta = \{\omega_i \in \mathbb{R}^{d_h \times d_{in}}, \mathbf{b}_i, \phi_i \in \mathbb{R}^{d_h}, \mathbf{W}_i \in \mathbb{R}^{d_h \times d_h}, \mathbf{W}_i^{\text{out}} \in \mathbb{R}^{d_h \times d_{out}}, \mathbf{b}_i^{\text{out}} \in \mathbb{R}^{d_{out}}\}$.

A useful property of this formulation is that the network output can be expressed equivalently as a sum of sines with varying amplitude, frequency, and phase [15].

$$\mathbf{y}_i = \sum_{j=0}^{N_{\text{sine}}^{(i)}-1} \bar{\alpha}_j \sin(\bar{\omega}_j \mathbf{x} + \bar{\phi}_j), \quad (2)$$

where $\bar{\alpha}_i$, $\bar{\omega}_i$, and $\bar{\phi}_i$ depend on the parameters of the MFN, and the number of terms in the sum for an N_L layer network is given as (see supplemental for derivation)

$$N_{\text{sine}}^{(N_L)} = \sum_{i=0}^{N_L-1} 2^i d_h^{i+1}. \quad (3)$$

This property stems from the repeated Hadamard product of sines and the trigonometric identity that

$$\sin(a) \sin(b) = \frac{1}{2} (\sin(a + b - \pi/2) + \sin(a - b + \pi/2)). \quad (4)$$

By applying this identity through the layers of the network, the output can be reduced to a single sum of sines.

3.2. Frequency Spectrum

We exploit the property that MFNs can be expressed as a sum of sines to create band-limited networks. This is achieved by designing the architecture so that the frequency of all represented sines never exceeds a desired threshold.

To this end, we freeze (i.e., do not optimize) the frequencies, or entries of ω_i , and set them to a bandwidth in $[-B_i, B_i]$ using random uniform initialization. Then, since **products of sines result in summed frequencies** (Eq. 4), the total bandwidth of an output at layer i of the network is less

than or equal to $\sum_{j=0}^i B_j$ and the maximum bandwidth is $B = \sum_{i=0}^{N_L-1} B_i$ (see Fig. 2).

When representing signals across a finite input domain, e.g., with input coordinates $\mathbf{x} \in [-0.5, 0.5]^{d_{\text{in}}}$, it is not necessary to represent all frequencies continuously. Instead, we can assume that the represented signal is periodic, so we are only required to represent discrete frequency values whose spacing is $1/T$, where T is the periodicity or extent of the signal in the primal domain. Moreover, using discrete frequencies allows complete characterization of the network spectrum by applying a fast Fourier transform to a uniformly sampled network output (shown in Fig. 2 for image fitting).

We also analyze the distribution of sine frequencies in the network. Briefly, sines in the network can be associated with one of the N_L terms in the summation of Eq. 3. Then, considering the probability of sines originating from each term results in a compound random variable that gives the overall distribution of frequencies. We provide an extended derivation in the supplemental, showing that the distribution is approximately zero-mean Gaussian with variance

$$\text{Var}(\omega_i) \cdot \sum_{m=0}^{N_L-1} m \cdot \frac{2^{N_L-1-m} d_h^{N_L-m}}{\sum_{i=0}^{N_L-1} 2^i d_h^{i+1}}. \quad (5)$$

In practice, this Gaussian distribution of frequencies results in a greater parameterization of low frequencies in the network; this may be a useful inductive bias since low-frequency Fourier coefficients typically have a much higher amplitude than high-frequency coefficients in natural signals [68].

To facilitate representing signals at multiple resolutions, we introduce linear layers at intermediate stages throughout the network to extract band-limited outputs (see Fig. 2). By supervising the outputs of these layers, we can train BACON to fit a signal at multiple scales simultaneously. Interestingly, because the outputs are band-limited, BACON can be trained in a semi-supervised fashion where the bandwidth of the supervisory signal need not match the desired bandwidth of the output of the network, demonstrated in Fig. 2 for image fitting.

3.3. Initialization Scheme

Finally, we derive a principled initialization scheme that ensures the distribution of activation functions at the output of each layer is distributed uniformly at the beginning of training. We compare our initialization scheme to the initialization proposed by Fathony et al. [15] in Fig. 3. Our proposed scheme resolves a problem with vanishingly small activations for deep networks and results in standard normal distributed activations after each linear layer.

In the supplemental, we provide an extended derivation, which we summarize as follows. Assume the input to the

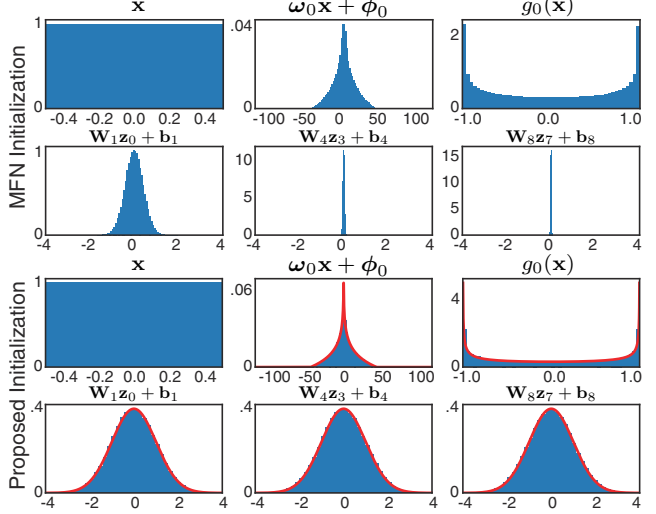


Figure 3. Comparison of distribution of activations at initialization. The initialization scheme proposed for MFNs [15] results in vanishingly small activations for deep networks (shown for layers 0, 1, 4, and 8). The proposed initialization scheme maintains a standard normal distribution after each linear layer (all distributions shown for a network with $d_h = 1024$), and activations at intermediate outputs closely match our analytical derivations (red lines, see supplemental for details).

network is uniformly distributed $\mathbf{x} \sim \mathcal{U}(-0.5, 0.5)$, with $\omega_i \sim \mathcal{U}(-B_i, B_i)$ and $\phi_i \sim \mathcal{U}(-\pi, \pi)$, where we describe the distribution of each element of the matrix or vector. Then, $\omega_i \mathbf{x} + \phi_i$ is distributed as

$$\begin{cases} 1/B_i \log(B_i / \min(|2x|, B_i)), & -B/2 \leq x \leq B/2 \\ 0 & \text{else} \end{cases}$$

and $g_i(\mathbf{x}) = \sin(\omega_i \mathbf{x} + \phi_i)$ is approximately arcsine distributed with variance 0.5 (see supplemental, red plots in Fig. 3). Now, let $\mathbf{W}_i \sim \mathcal{U}[-\sqrt{6/d_h}, \sqrt{6/d_h}]$. Then we have that $\mathbf{W}_1 g_0(\mathbf{x}) + \mathbf{b}_1$ converges to the standard normal distribution with increasing d_h (see supplemental). Finally, the Hadamard product $g_1(\mathbf{x}) \circ (\mathbf{W}_1 \mathbf{z}_0 + \mathbf{b}_1)$ is the product of arcsine distributed and standard normal random variables which again has a variance of 0.5. Applying the next linear layer results in another standard normal distribution, which is also the case after all subsequent linear layers (see red plots of Fig. 3).

4. Experiments

We demonstrate BACON on three separate tasks: image fitting, multiview reconstruction using neural radiance fields, and 3D shape fitting using signed distance functions.

4.1. Image Fitting

We use an image fitting task to evaluate the performance of BACON and to demonstrate its band-limited behavior.

the def of \mathbf{W}_i result in Normal distribution



Figure 4. Image fitting results. We train networks using Fourier Features [65], SIREN [60], and integrated positional encoding (PE) [4] to fit an image at 256×256 ($1 \times$) resolution. We show network outputs at $1/4$ and $4 \times$ resolution. Fourier Features and SIREN fit to a single scale and show aliasing when subsampled. Integrated PE is explicitly supervised at $1/4$ and $1 \times$ resolution and learns reasonable anti-aliasing; however, all methods except BACON show high-frequency artifacts at $4 \times$ resolution (insets). BACON is supervised at a single scale and learns band-limited outputs that closely match a low-pass filtered reference image (see left column and Fourier spectra insets).

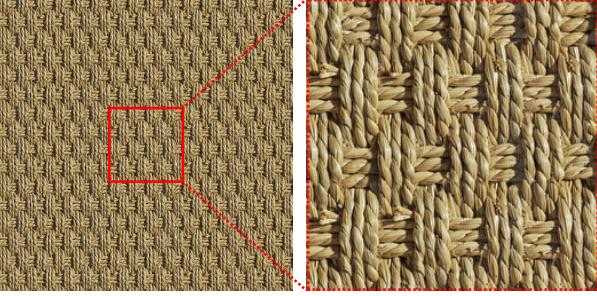


Figure 5. BACON periodic extrapolation behavior¹.

BACON is compared to three other baselines: a network with Gaussian Fourier Features positional encoding [65], SIREN [60], and the integrated positional encoding of Mip-NeRF [4], which is scale-dependent.

We initialize all networks with 4 hidden layers, 256 hidden features, and we train on the 256×256 resolution image for 5000 iterations using PyTorch [49] and Adam [29]. The batch size is set to be equal to the number of image pixels. For **Fourier Features** and **integrated positional encoding**, we use encoding scales of 6 and 10, respectively, which results in good reconstruction quality while reducing artifacts from overfitting. For **SIREN**, we initialize the frequency parameter to $\omega_0 = 30$.

Fourier Features and SIREN are trained to minimize the loss $\mathcal{L}_{\text{img}} = \|\mathbf{y} - \mathbf{y}_{\text{GT}}\|_2^2$, where \mathbf{y} is the network output and \mathbf{y}_{GT} are the image pixel values. For **BACON**, we sum this loss over all network outputs, with explicit supervision at all scales on the full-resolution image. The **integrated positional encoding network** is supervised explicitly on anti-aliased image pixels at $1/4$, $1/2$, and full resolution, following Barron et al. [4]. Finally, we initialize BACON to have a maximum bandwidth B of 0.5 cycles/pixel, which is the Nyquist limit for the image. The frequencies ω_i of BACON

are initialized so that the outputs \mathbf{y}_1 , \mathbf{y}_2 , and \mathbf{y}_4 are constrained to quarter, half, or full bandwidth. That is, we set $B_0 = B_1 = B/8$, and $B_2 = B_3 = B_4 = B/4$ such that $\sum_i B_i = B$.

Results of image fitting on a test scene from the Kodak dataset [11] are shown in Fig. 4. Since **Fourier Features** and **SIREN** only represent the signal at the trained resolution, sampling the network at $1/4$ resolution results in aliasing. We show the interpolation performance of these networks by evaluating a $4 \times$ upsampled grid of 1024×1024 pixels. The upsampled reference image shows image interpolation given by zero padding in the Fourier domain. All methods except BACON have non-zero high-frequency spectra and exhibit artifacts in the reconstruction. We show additional image fitting experiments in the supplemental.

Periodic Extrapolation. Since BACON uses discrete frequencies at each sine layer $g_i(\mathbf{x})$, the representation is periodic. We demonstrate this by fitting a seamless texture using coordinates $\mathbf{x} \in [-0.5, 0.5]$ (red square of Fig. 5) and querying the network output for $\mathbf{x} \in [-2, 2]$.

4.2. Neural Radiance Fields

Neural radiance fields (NeRF) [41] have become a popular method for multiview 3D reconstruction and neural rendering. The method operates on a dataset of multiview images with known camera positions, where each image pixel is associated with a ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ that extends from the camera center of projection \mathbf{o} in the direction \mathbf{d} passing through the pixel. A pixel color $\mathbf{C}(\mathbf{r})$ is predicted using the **volume rendering equation** to integrate predicted intermediate values of color \mathbf{c} and opacity σ along the ray [4]. In practice, a neural network is queried to evaluate samples of \mathbf{c} and σ along each ray $\mathbf{r}(t)$, and the volume rendering integral is evaluated using quadrature as [37, 40]

¹Image: <https://www.sketchuptextureclub.com/>

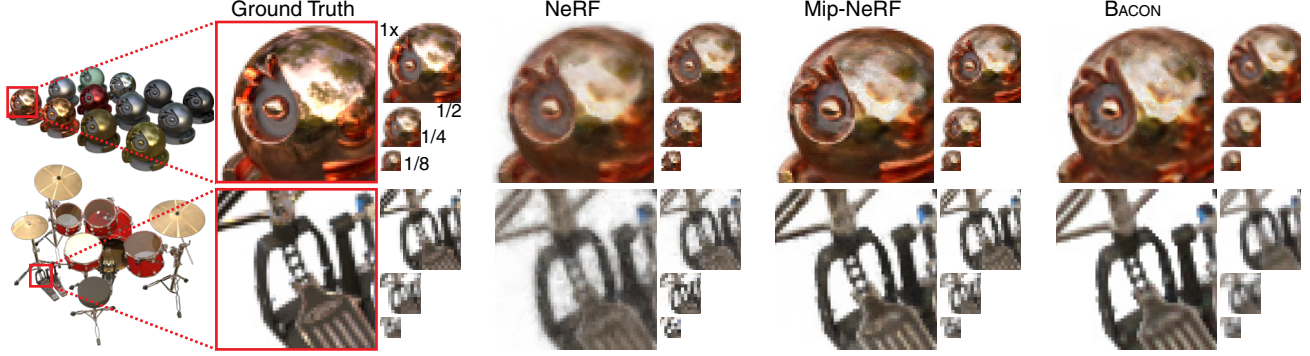


Figure 6. Neural rendering results. We compare NeRF [41], Mip-NeRF [4], and BACON supervised on a multiscale synthetic dataset [4]. BACON captures higher frequency details better than NeRF while requiring fewer parameters to render at 1/2, 1/4, and 1/8 resolution.

	PSNR \uparrow					# Params.			
	1 \times	1/2	1/4	1/8	Avg.	1 \times	1/2	1/4	1/8
NeRF	26.734	28.941	29.297	26.464	27.859	511K			
Mip-NeRF	29.874	31.307	32.093	32.832	31.526	511K			
BACON	27.430	28.066	28.520	28.475	<u>28.123</u>	531K	398K	266K	133K

Table 1. Performance of NeRF, Mip-NeRF, and BACON averaged across the multiscale Blender dataset. BACON achieves better average performance than NeRF while requiring fewer parameters to render the lower resolution images.

$$\mathbf{C}(\mathbf{r}, \mathbf{t}) = \sum_j T_j (1 - \exp(-\sigma_j(t_{j+1} - t_j))) \mathbf{c}_j, \quad (6)$$

$$\text{with } T_j = \exp\left(-\sum_{i' < i} \sigma_{i'}(t_{i'+1} - t_{i'})\right),$$

where T_j represents the transmittance or visibility of a point on the ray, and the values $w_j = T_j(1 - \exp(-\sigma_j(t_{j+1} - t_j)))$ can be interpreted as alpha compositing weights applied to the predicted colors \mathbf{c}_j . After training, novel views can be rendered by simply evaluating the corresponding rays.

We evaluate BACON for this task and compare to NeRF and Mip-NeRF baselines trained on a multiscale Blender dataset [4] with images at full (512 \times 512), 1/2, 1/4, and 1/8 resolution. For the baselines, we use the implementations of Barron et al.² [4]. All networks are trained according to the procedure of Mip-NeRF; we use the Adam optimizer with a batch size of 4096 rays and 1e6 training iterations. The learning rate is annealed logarithmically from 1e-3 to 5e-6 for BACON and 5e-4 to 5e-6 for the baselines. All networks are composed of 8 hidden layers with 256 hidden features.

For BACON, we adapt the training procedure and architecture as follows. Rays within the multiscale Blender dataset fall within an 8 by 8 unit volume ($\mathbf{r}(t) \in [-4, 4]^3$), and we find that setting the maximum bandwidth B to 64 cycles per unit interval allows fitting high frequency image details. To simplify the training procedure and improve the interpretability of the learned volume, we evaluate all methods without the viewing direction input originally used for NeRF. Thus the input to all networks is a

3D coordinate corresponding to the position along the ray $\mathbf{r}(t)$. BACON produces four outputs, one for each scale of the dataset: $\mathbf{y}_i, i \in [2, 4, 6, 8]$. The B_i constrain each output to 1/8, 1/4, 1/2, and full resolution, with $\sum_{i=0}^2 B_i = B/8$, $\sum_{i=0}^4 B_i = B/4$, and so on. We also adapt the hierarchical sampling procedure of NeRF [41], wherein the alpha compositing weights w_j from an initial forward pass are used to resample the ray in regions of non-zero opacity. To improve efficiency, we use the lowest-resolution output of the network for this initial forward pass and apply the following loss function on pixels rendered using the resampled rays with 256 samples.

$$\mathcal{L}_{\text{BACON}} = \sum_{i,j,k} \|(\mathbf{C}_k(\mathbf{r}_i, \mathbf{t}_j) - \mathbf{C}_{\text{GT},k}(\mathbf{r}_i))\|_2^2, \quad (7)$$

where i, j , and k index rays, ray positions, and dataset scales, respectively. Note that we explicitly supervise each output of BACON on the multiscale ground truth images to facilitate quantitative evaluation. Finally, we adopt the regularization strategy of Hedman et al. [22] to penalize non-zero off-surface opacity values (see supplemental).

Qualitative and quantitative evaluations of BACON for neural rendering, are shown in Fig. 6 and Table 1. BACON generally achieves better performance than NeRF trained on the multiscale dataset with the most dramatic improvement at 1/8 and 1 \times resolution. We report PSNR at each scale, averaged over all scenes in the multiscale Blender dataset in Table 1. In Fig. 6, we observe that BACON recovers higher frequency details compared to NeRF on the *Materials* and *Drums* scenes. Mip-NeRF incorporates an additional mechanism which changes the positional encoding along each ray to account for the expansion of the viewing frustum, and achieves the best performance. Still, we find that BACON produces high-quality novel view synthesis while using only a fraction of the parameters at low resolution.

Additionally, we can use BACON to learn semi-supervised multiscale decompositions of the neural radiance fields. In this case, we train each output scale at the full resolution, and BACON automatically learns band-limited representations at the intermediate output layers. We show

²<https://github.com/google/mipnerf>

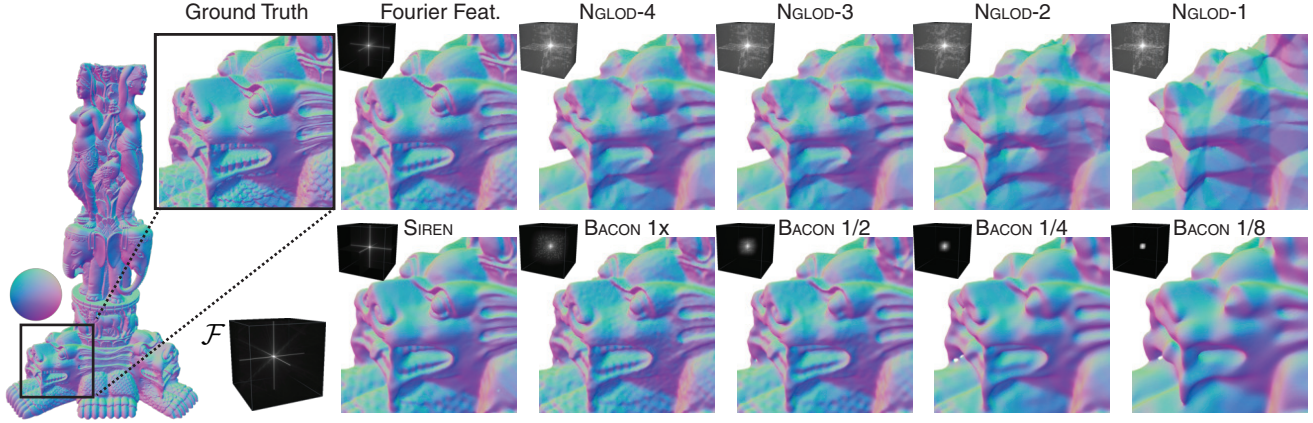


Figure 7. Shape fitting results. Results on the Thai Statue from the Stanford 3D Scanning Repository are shown for levels-of-detail 1–4 of Neural Geometric Level of Detail (NGLOD) [64], Fourier Features [65], SIREN [60], and BACON. All methods perform similarly at their highest detail output (see Table 2), but BACON learns a smooth multiscale decomposition of the shape. Insets show the spectra of the extracted signed-distance functions, revealing the band-limited output of BACON. Additional results included in the supplemental.

	FF	SIREN	NGLOD-4	NGLOD-5	BACON 1x
# Params.	527K	528K	1.35M	10.1M	531K
Chamfer↓	2.166e-6	2.780e-6	8.358e-6	2.422e-6	<u>2.198e-6</u>
IOU↑	9.841e-1	9.751e-1	9.479e-1	9.811e-1	<u>9.833e-1</u>

Table 2. Shape fitting performance of Fourier Features [65], SIREN [60], Neural Geometric Level of Detail (NGLOD) [64], and BACON averaged across 5 test scenes (detailed in main text). All methods achieve roughly comparable performance, including BACON despite simultaneously representing multiple scales. Multiple levels of detail are shown for NGLOD, which require more parameters to populate the explicit feature grids.

an example of this for the *Lego* scene in Fig. 1. Additional results for BACON in the explicitly supervised and semi-supervised cases are shown in the supplemental.

4.3. 3D Shape Representation

Neural representation networks have shown promise for representing and manipulating 3D shapes. BACON is well-suited for this task, and we evaluate its performance on a range of shapes from the Stanford 3D scanning repository³.

We compare BACON to Fourier Features [65], SIREN [60], and Neural Geometric Level of Detail (NGLOD) [64], a representation which optimizes explicit features stored on a sparse voxel octree. All networks are trained to directly fit a signed distance function (SDF) estimated from a ground truth mesh.

For BACON, Fourier Features, and SIREN, we use networks with 8 hidden layers and 256 hidden features. For Fourier Features (Gaussian encoding), we set the encoding scale to 8 and for SIREN, we set $\omega_0 = 30$. We train on locations sampled from the zero level set and add Laplacian noise; this results in an exponential decay in the number of samples off the zero level, as proposed by Davies et al. [12].

We find that the width of the Laplacian distribution has a large impact on performance. Setting the variance σ_L^2 too small results in poor **off-surface fitting**, but setting the variance too high reduces the number of samples on the zero level set, degrading the appearance of the surface. Thus, we introduce a **coarse and fine sampling procedure** wherein we produce “fine” samples using a small variance of $\sigma_L^2 = 2e-6$ and “coarse” samples with $\sigma_L^2 = 2e-2$. Samples are drawn in the domain $[-0.5, 0.5]^3$, and we initialize the frequencies of BACON similar to the NeRF experiments (additional details in supplemental). We train using a **loss function**

$$\mathcal{L}_{\text{SDF}} = \lambda_{\text{SDF}} \|\mathbf{y}^c - \mathbf{y}_{\text{GT}}^c\|_2^2 + \|\mathbf{y}^f - \mathbf{y}_{\text{GT}}^f\|_2^2, \quad (8)$$

where \mathbf{y} is the network output, \mathbf{y}_{GT} represents the ground truth SDF values, the f and c superscripts indicate fine and coarse samples, and λ_{SDF} is a weight that we set to 0.01 for all experiments. For BACON we sum the loss at all outputs.

We train Fourier Features, SIREN, and BACON on each dataset for 200,000 iterations with a batch size of 5,000 coarse and 5,000 fine SDF samples. Models are optimized using Adam [29], and we logarithmically anneal the learning rate of each method from $1e-2$ (BACON), $1e-3$ (Fourier Features), and $1e-4$ (SIREN) to a final value of $1e-4$ during the course of training. For NGLOD, we use the default training settings in the authors’ code⁴, which samples 500,000 points at each training epoch, uses a batch size of 512, and trains for 250 epochs. We train NGLOD models with a maximum of 4 or 5 levels of detail. Additional levels of detail add successive layers in their octree representation, resulting in improved performance, but also an increased number of parameters.

We fit each method to four scenes from the Stanford 3D Scanning Repository (*Armadillo*, *Dragon*, *Lucy*, and *Thai Statue*), as well as a simple sphere baseline (all ob-

³<http://graphics.stanford.edu/data/3Dscanrep/>

⁴<https://github.com/nv-tlabs/nglod>

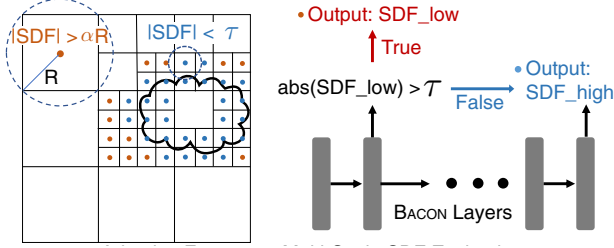


Figure 8. Adaptive-frequency multiscale SDF evaluation for fast mesh extraction. We propose a multi-scale evaluation which subdivides non-empty cells, i.e., when the SDF is smaller than the cell radius. We further accelerate the evaluation by using fewer layers (red) for regions that do not require high-frequency details. We evaluate the full network (blue) only when query locations are sufficiently close ($|\text{SDF}| < \tau$) to the surface.

	Dense Grid (512^3)	Adaptive-Frequency	Adaptive + Multiscale (Proposed)
Time (s)	17.91	5.50	0.222

Table 3. SDF evaluation time. The proposed Adaptive + Multiscale method achieves a roughly $80\times$ speedup over naive evaluation of the SDF on a dense grid (averaged over 5 test scenes).

jects are shown in the supplemental). The models are extracted at 512^3 resolution using marching cubes and evaluated using Chamfer distance and **intersection over union** (IOU), and we report these numbers averaged over the 5 scenes in Table 2. The highest resolution outputs of all methods achieve generally comparable performance, though note that NGLOD-5 requires over an order of magnitude more parameters than the other representations, and BACON achieves this despite representing all scales simultaneously.

We see similar qualitative trends in Fig. 7, with Fourier Features, SIREN, NGLOD-4, and BACON all producing detailed reconstructions of the *Thai Statue* scene. BACON produces a smooth reconstruction at multiple scales because of its band-limited output layers. This can be compared to the low-resolution outputs of NGLOD, which show fewer finer details, but also have coarse, angular artifacts from the ReLU non-linearity used in the network. This follows from the NGLOD frequency spectrum (see Fig. 7), which is non-zero for high frequencies, including at the coarsest scale.

Accelerated Marching Cubes. We observe that the band-limited, multi-output nature of our network allows efficient allocation of resources when evaluating SDFs on a dense grid for mesh extraction. The key idea is to use the lower-layer output of BACON when a cell is far away from the surface (Fig. 8). That is, we early-stop the computation within the network when $|\text{SDF}| < \tau$, where we set τ to $0.7\times$ the finest voxel size. This adaptive computation significantly reduces the mesh extraction time (Table 3).

Moreover, we propose a multiscale approach for further acceleration. As SDF values indicate the distance to the closest surface, we can consider a cell to be empty

when $|\text{SDF}| > \alpha R$, for circumsphere radius R and some $\alpha > 1$ that improves robustness to imperfect SDFs (we use $\alpha = 2$). Starting from the coarsest resolution grid, we subdivide a cell only when $|\text{SDF}| < \alpha R$, efficiently pruning empty spaces, until the finest level. Multiscale extraction approaches have been proposed for extracting occupancy fields from coordinate networks [38], but using an SDF facilitates pruning since each sample reveals a region of empty space.

We combine the two strategies by applying the adaptive-frequency evaluation on each level of the multiscale grids, leading to roughly $80\times$ faster SDF evaluation than the naive approach as shown in Table 3 (see supplemental results, all timings evaluated on an NVIDIA RTX A6000 GPU).

5. Conclusion

In this work we take important steps towards making coordinate networks interpretable and scale aware. Our approach enables analyzing and controlling the spectral bandwidth of the network at intermediate layers, allowing multiscale signal representation, even without explicit supervision. Since we can characterize the behavior of the network using Fourier analysis, its outputs are predictable, even at unsupervised locations. Moreover, BACON’s unique intermediate outputs can address the characteristic slow inference time of the coordinate-networks, via adaptive frequency evaluation. We show that BACON outperforms other single-scale coordinate networks for multiscale image fitting, neural rendering, and 3D scene representation.

Limitations. We also highlight a few limitations of BACON and promising future directions. In this work, we demonstrated BACON for representing two- and three-dimensional signals. Fitting signals in higher dimensions may require more parameters to achieve dense spectral coverage due to the curse of dimensionality. Still, it may be possible to optimize the initialization of frequencies in a way that maximizes spectral coverage and mitigates this challenge. Also, our current work is limited to single scene overfitting. However, many generative models work by increasing the frequency of the output using successive upsampling layers [27], which is similar in spirit to our method. Recent work on band-limited models for image synthesis has shown great promise [26], so applying BACON for generative modeling is an exciting area for future research.

Societal Impact. We condemn the use of scene representation networks, including BACON, for ill-intentioned purposes such as malicious deepfakes or spreading misinformation, and we emphasize the importance of research to detect and thwart such efforts. See, e.g., Tewari et al. [66] for an extended discussion of strategies for mitigating misuse of neural rendering techniques.