

Adaptively Placed Multi-Grid Scene Representation Networks for Large-Scale Data Visualization

Category: Research

Paper Type: algorithm/technique

Abstract—Scene representation networks (SRNs) have been recently proposed for compression and visualization of scientific data. However, state-of-the-art SRNs do not adapt the allocation of available network parameters to the complex features found in scientific data, leading to a loss in reconstruction quality. We address this shortcoming with an adaptively placed multi-grid SRN (APMGSRN) and propose a domain decomposition training and inference technique for accelerated parallel training on multi-GPU systems. We also release an open-source neural volume rendering application that allows plug-and-play rendering with any PyTorch-based SRN. Our proposed APMGSRN architecture uses multiple spatially adaptive feature grids that learn where to be placed within the domain to dynamically allocate more neural network resources where error is high in the volume, improving state-of-the-art reconstruction accuracy of SRNs for scientific data without requiring expensive octree refining, pruning, and traversal like previous adaptive models. In our domain decomposition approach for representing large-scale data, we train a set of APMGSRNs in parallel on separate bricks of the volume to reduce training time while avoiding overhead necessary for an out-of-core solution for volumes too large to fit in GPU memory. After training, the lightweight SRNs are used for realtime neural volume rendering in our open-source renderer, where arbitrary view angles and transfer functions can be explored. A copy of this paper, all code, all models used in our experiments, and all supplemental materials and videos are available at <https://github.com/anonymousvisuser1036/APMGSRN>.

Index Terms—Scene representation network, deep learning, scientific visualization, volume rendering

1 INTRODUCTION

Scene representation networks (SRNs) are compact neural networks that map input coordinates to output scalar field values [15, 18, 23–25]. SRNs use a small footprint on disk (data reduction between 10 – 10000×) while being efficient to evaluate with random access. With these benefits, SRNs have been used for compression [16, 30] and volume rendering [30, 32] for scientific data up to sizes of 1TB.

Despite SRNs’ popularity, there are two shortcomings of current state-of-the-art SRNs for scientific data visualization. First, the SRN architectures used in recent approaches are not designed to adapt network resources to more important regions, and make the inherit assumption that the volume it will represent has uniform complexity, leading to inefficient use of the network parameters in homogeneous regions and limiting model performance. Recently proposed SRN models have incorporated adaptive parameter allocation with spatially adaptive quadtrees and octrees that refine on the scene’s geometry or features, but these approaches require significant extra storage and computation to store and search the tree structure during training [15, 17, 25, 33], taking hours or days to fit on images with fewer degrees of freedom than most 3D scientific data. In addition to this, the tree structure itself is not directly learnable, so the trees are updated every set number of iterations with an ad-hoc method while training. The tree structure also scales poorly with dimensionality, increasing storage and computation requirements exponentially as the number of dimensions or tree depth increases. Second, existing training routines for large-scale data are proposed for a single GPU with out-of-core sample streaming [32], which introduces I/O overhead and does not take advantage of the multiple GPUs often available on a single compute node from the supercomputers generating the large-scale data.

In this work, we address the two limitations mentioned above with a novel SRN model and domain decomposition training routine for fitting large-scale data. First, we address the lack of adaptivity in state-of-the-art SRNs for scientific data with a novel SRN architecture called adaptively placed multi-grid SRN (APMGSRN). Instead of using trees as the adaptive data model in our network, APMGSRN uses a set of spatially adaptive feature grids, shown in Figure 1, whose extents are defined by learned transformation matrices which transform global space into local feature space, where the grid is defined as the unit cube $[-1, 1]^3$. To guide the grids to cover spatial regions where the model has relatively higher error, we develop a custom *feature density loss*, that calculates the relative entropy from the current feature grid

density to a target feature density, where the assumption is higher feature density improves reconstruction (shown by other feature grid networks [3, 25, 30, 33]). This alone is not enough to train the feature grids though, as the density of a grid is a step function that has a gradient of 0 everywhere. Therefore, we approximate the feature density with a differentiable function that closely matches the step function called a *flat-top gaussian*. We also develop a specific training routine for the feature grid’s transformation matrices with delayed start and early stopping to improve accuracy and converge quicker. Our adaptive grid architecture does not require expensive tree searching or pruning like other adaptive models, and dynamically allocates more neural network resources to regions of higher error for any volume, improving state-of-the-art SRN performance by using network parameters efficiently.

Using our APMGSRN architecture as a building block, we develop a *domain decomposition* training strategy that fits a large-scale volume in a model-parallel fashion. Our domain decomposition training approach divides data in the volume on a grid of bricks, and trains one SRN per block of data. The blocks of data fit in GPU memory, so there are no inefficiencies from an out-of-core data sampling method. We also assume that the large-scale data being fit are often generated by machines equipped with multiple GPUs per node, so we use all available GPUs to train multiple networks in parallel, reducing total training time. Inference in this set of models is more complicated now, since a search is necessary to find which model was trained on the spatial domain for each point being queried. To accelerate inference in a domain decomposition model, we use a hash function that maps spatial coordinates to the hashtable entries for the correct model to use for inference in parallel. Not only do domain decomposition models allow fitting a 450 GB volume in under 7 minutes on 8 GPUs, but they also increase reconstruction accuracy over a single model at the same number of total model parameters.

In summary, our contributions are twofold:

- A novel SRN architecture called adaptively placed multi-grid SRN (APMGSRN) with adaptive feature grids that localize trainable network parameters on regions of the volume with high error during training
- A domain decomposition training and inference strategy that trains multiple SRNs in parallel for fitting large-scale data quicker on capable machines

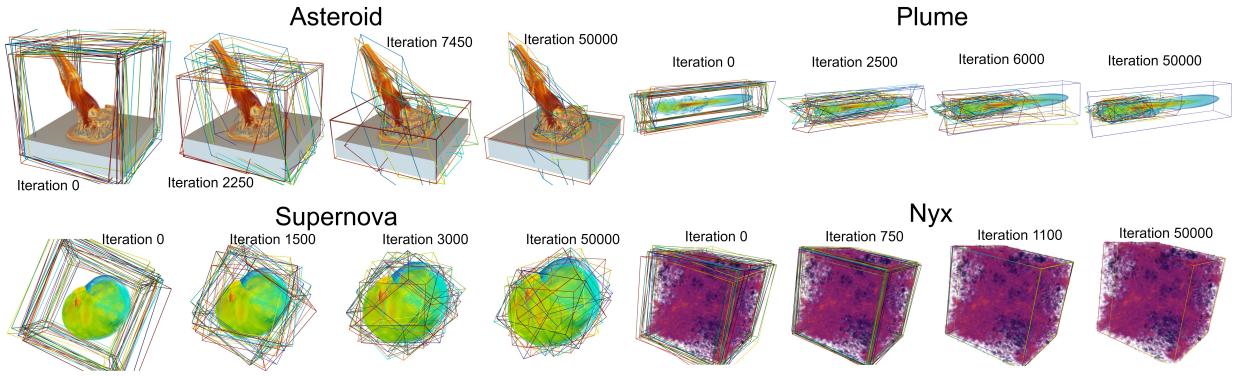


Fig. 1: Examples of our adaptively placed feature grids fitting to volumes during training. The feature grids find volume-specific regions to cover to maximize reconstruction accuracy. Video of evolution of grids during training available in supplemental materials.

Our APMGSRN architecture and domain decomposition modelling technique are evaluated on several scientific datasets ranging from volumes of size $128^2 \times 512$ (32MB to store) up to $10240 \times 1536 \times 7680$ (450GB to store). We compare the reconstruction quality of our proposed APMGSRN architecture (single model, not decomposed) with other state-of-the-art models, and demonstrate that our adaptive feature grids improve reconstruction quality over state-of-the-art at similar model sizes. We further evaluate the performance of our domain decomposition strategy by fitting two datasets that are 250 GB and 450 GB. Lastly, we evaluate the rendering performance of our proposed model compared with other state-of-the-art SRNs in our open-source renderer, which we release with our code to offer an easy to use, plug-and-play, neural volume renderer for INRs trained on scientific data. We provide all code for our model, training, and neural volume renderer on GitHub: <https://github.com/anonymousvisuser1036/APMGSRN>, along with installation instructions and supplemental videos.

2 RELATED WORKS

As our method is primarily related to scene representation networks, we review related literature covering SRN architectures and training, examining applications in the computer vision and sci-vis domains. Our work is part of a larger domain of research called DL4SciVis, for which Wang and Han provide a comprehensive survey [29].

2.1 Scene representation networks

Scene representation networks, also called implicit neural representations (INRs) or coordinate networks, are networks that encode a given scene with the weights of the neural network, such that input coordinates are mapped to output values. In the context of computer vision, SRNs can learn to represent images [17, 24], signed distance functions [17, 24], or radiance fields [18, 23, 33]. In the context of scientific data and visualization, SRNs have been used to model 3D scalar fields and time-varying scalar fields [16, 30, 32]. We broadly classify the architecture of SRN models into two categories: fully connected and grid-based encoding.

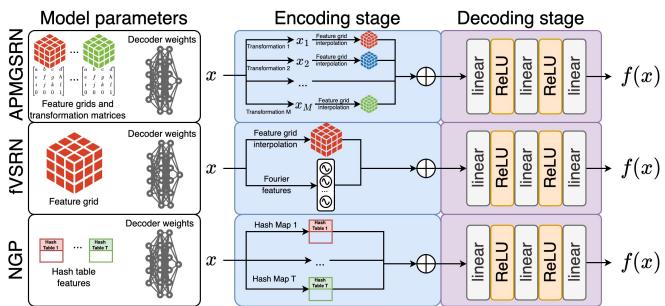


Fig. 2: Comparison of our APMGSRN with other state-of-the-art models fVSRN [30] and NGP [20]. The \oplus operator represents concatenation.

Fully connected SRNs. Fully connected SRNs use only linear layers with activations between them to map an input coordinate to output value, making them slow to train and perform inference on. SIREN [24] is an example of a fully connected SRN which uses sinusoidal activations to fit data more accurately with useful higher-order derivatives. Fourier features by Tancik et al. [27] shows that encoding input coordinates to a higher dimension using multiple frequencies before going through a fully connected can help INR accuracy. NeRF [18] is a fully connected SRN for modelling a radiance field, and trains using a set of input images taken from different points in the real-world or 3D modelling scene. NeRF models take a coordinate and viewing direction as input and return the color and volume density for that spatial location. After training, novel viewpoints can be rendered by volume rendering using the model. AutoInt [14] requires far less model inferences during neural volume rendering compared to NeRF by exploiting the fact that the gradient of a SIREN network is another SIREN network that shares the same weights as the original network. Then the “gradient” network can be trained as a NeRF, resulting in the original “integral” network learning the integration of colors and density through rays in space. Lu et al. [16] use a SIREN-based architecture with residual connections to model 3D volumetric data in a compressed format. Höhlein et al. [7] study the capabilities of INRs for compressing meteorological ensemble data. Han and Wang [6] propose CoordNet, a single implicit model based on the residual siren architecture of Lu et al. [16] that is used for spatiotemporal super resolution and novel view synthesis. Tancik et al. [26] propose Block-NeRF, which uses multiple NeRF models representing disjoint subsections of blocks of a region of San Francisco. This is similar to our approach, but models represent geometry within a specified radius from the model’s center with large overlap with other models, while our approach divides the domain with rectangular bricks as it fits the voxel grid representation better.

Grid-based encoding SRNs. In contrast to fully connected SRNs, grid-based encoding SRNs move a large majority of the network parameters to an efficient encoding scheme that transforms the input coordinates to a high-dimensional feature space. SRNs with grid-based encoding schemes train and infer faster than their fully connected counterparts due to the efficient encoding and limited fully connected operations. Weiss et al. [30] create fVSRN, an SRN that models 3D scientific data using a feature-grid encoder which places a low-resolution feature grid over the data domain and interpolates within the feature grid to obtain a feature vector for decoding. Yu et al. [33] accelerate a pre-trained NeRF model with spherical harmonics and an octree data structure for empty-space skipping, enabling realtime rendering of NeRF models. Liu et al. [15] propose neural sparse voxel fields, which uses a coarse 3D grid of voxels that are pruned and split at checkpoints during training to reduce the training time of NeRF with improved accuracy. Genova et al. [4] learn an implicit representation of a mesh from depth images using a dense grid and many local deep implicit functions. Takikawa et al. [25] propose NGLOD, which learns an octree data structure for realtime rendering of an SRN fitting a signed-distance

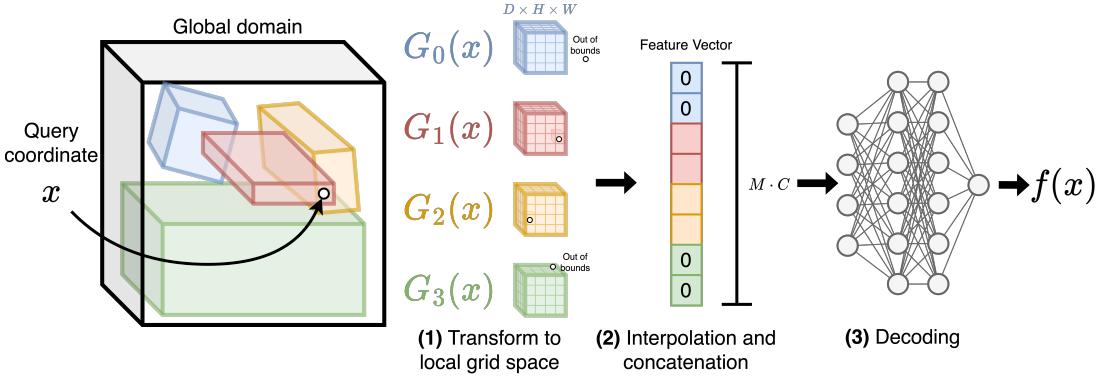


Fig. 3: An overview of the APMGSRN architecture. In (1), a query coordinate is transformed into each of the M grid’s local coordinate systems learned by the transformation matrices, where grid local extents are assumed to be $[-1, 1]^3$. (2), the local coordinates are used to trilinearly interpolate within each feature grid (of resolution $D \times H \times W$) to obtain the corresponding feature for each grid. If the coordinate is out of bounds, zeros are returned instead. The resulting $M \cdot C$ features are concatenated into a feature vector. Steps (1) and (2) are considered the encoding process. In step (3), the feature vector is decoded in a shallow MLP for the final output value $f(x)$.

function. Martel et al. [17] create an adaptive coordinate network that uses quadtrees/octrees during training, updated by solving an integer linear programming optimization problem every set number of iterations. Chen et al. [3] use tensor (de)composition to reduce the space complexity of 3D feature grids with their VM decomposition of a tensor. Müller et al. [20] introduce hash grid encoding for neural graphics primitives (NGP), which uses random hashing to map input coordinates to features in a hash table at multiple levels of detail. Wu et al. [32] use the hash grid architecture [20] to model scientific data up to 1TB in size. Concurrent work by Wu et al. [31] uses a meta-learning approach that predicts the network weights of a set of models of the hash grid architecture [20] for novel view synthesis and neural volume rendering of temporally interpolated volumes and dynamic global shadows.

Existing adaptive models such as ACORN [17], NSVF [15], and NGLOD [25] are expensive to train, store, and require ad-hoc steps for pruning/refining the tree data structure. Additionally, the tree-based adaptivity only helps for scenes with static surfaces. Volumetric scientific data do not generally have static surfaces, since different surfaces arise from visualizing different isovalue. Thus, the space-skipping and tree structures used by other methods ignore relevant regions of the scientific data. Our approach’s adaptive feature grids do not require surfaces to refine to, and is driven only by reconstruction error. APMGSRN is conceptually an adaptive version of fVSRN that splits the single highly-parameterized feature grid into many less-parameterized feature grids, where each feature grid now has the ability to adjust its transformation within the domain to focus on high-error regions. A visual comparison of our model and two state-of-the-art models compared with in this paper, fVSRN and NGP, is provided in Figure 2.

3 OVERVIEW

Our approach is composed of two components: (1) a novel SRN architecture called APMGSRN, and (2) a domain decomposition training and inference method to train on large-scale data. In section 4, we detail the APMGSRN architecture, a SRN model that adaptively learns where multiple feature grids should be spatially located while training. In section 5, we explain the domain decomposition approach, which dissects a volume into a 3D grid of networks that are each trained on their own local region.

4 ADAPTIVELY PLACED MULTI-GRID SCENE REPRESENTATION NETWORK

Recent state-of-the-art SRNs have found that using explicitly defined feature grids within the model reduces training and inference times [3, 20, 30]. In these feature grid-based models, a feature grid is interpolated at an input spatial coordinate to retrieve a feature vector, which is fed through a shallow multi-layer perceptron (MLP) to obtain the final output value. Tree structures have been added to SRN models to support adaptive allocation of network resources for specific regions [15, 17, 33],

but the approaches are ad-hoc, require hours to days to train, and have significant training/storage overhead to manage the tree structure. Additionally, the cost of the tree structures scale exponentially as the number of dimensions or tree depth increases.

Our method improves adaptivity within SRNs with a novel adaptively placed multi-grid scene representation network (APMGSRN), depicted in Figure 3. Instead of a tree structure, APMGSRN uses a set of multiple spatially adaptive feature grids (shown in Figure 1), each described by a learnable transformation matrix, for coordinate encoding before being decoded by a shallow MLP, as described in Section 4.1. Since a basic reconstruction loss is not enough to learn the transformation matrices properly, Section 4.2 describes our method to learn grid positions with a feature density-based loss function. Finally, we discuss the network training in Section 4.3. Our model gives strong reconstruction quality even with few model parameters thanks to the ability to transform the feature grids to fit the complexities in the data. Since our model does not rely on a tree structure, inference remains quick and storage costs remain low.

4.1 APMGSRN architecture

An APMGSRN is a function $f(x)$ that maps normalized spatial coordinates $x \in [-1, 1]^3$ to the scalar value at that location in space. The model architecture is composed of an encoder e and decoder d such that $f(x) = d(e(x))$, depicted in Figure 2 and Figure 3. Our encoder, described in Section 4.1.1, is our adaptively placed multi-grid encoder, and our decoder, described in Section 4.1.2, is a small fully-connected MLP that decodes the encoded feature vector to an output value.

4.1.1 Adaptively placed multi-grid encoder

The encoder in APMGSRN contains M learnable feature grids of resolution $D \times H \times W$ with C channels, represented as a tensor F with shape $[M, C, D, H, W]$.

Transformation to local space. In our model, the spatial extents of the i -th feature grid is defined by a 4×4 transformation matrix G_i , which transforms global-space coordinates to local-space coordinates according to the following:

$$\begin{bmatrix} p_x^l \\ p_y^l \\ p_z^l \\ 1 \end{bmatrix} = \begin{bmatrix} G_{i,0,0} & G_{i,0,1} & G_{i,0,2} & G_{i,0,3} \\ G_{i,1,0} & G_{i,1,1} & G_{i,1,2} & G_{i,1,3} \\ G_{i,2,0} & G_{i,2,1} & G_{i,2,2} & G_{i,2,3} \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} p_x^g \\ p_y^g \\ p_z^g \\ 1 \end{bmatrix}, \quad (1)$$

where $G_{i,r,c}$ is the r -th row and c -th column of the transformation matrix, p_x^g, p_y^g, p_z^g is the global coordinate, and p_x^l, p_y^l, p_z^l is the resulting local coordinate. Conceptually, the first 3 rows and columns of G_i are responsible for scale, rotation, and shearing of the grid, while the 4th column is responsible for translation. The grid’s local extents are assumed to be $[-1, 1]^3$, and global coordinates for the 8 corners of a grid

can be determined by calculating the inverse of G_i and transforming the coordinates of the local extents to global space. For shorthand, we consider $G_i(x_g)$ to be a function mapping a 3D global coordinate x_g to the coordinate in grid i 's local space according to [Equation 1](#).

Encoding. To encode an input global coordinate, the coordinate is transformed into each grid's local coordinate system via the defining transformation matrices. Then, each 3D local coordinate x_l is used to perform trilinear interpolation within each feature grid:

$$e_i(x_l) = \begin{cases} \text{interp}(F_i, x_l) & \text{if } -1 \leq x_l \leq 1 \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

where F_i is the feature grid for the i -th grid and $\text{interp}(F_i, x_l)$ performs trilinear interpolation within F_i at point x_l assuming the grid has vertices corner-aligned at the extents of $[-1, 1]^3$. When the input coordinate is not within the feature grid, the feature returned is 0 for each of the expected C channels. The encoding step for a single grid as described in [Equation 2](#) is performed for each grid, with results concatenated together. From global space, the encoding is calculated as:

$$e(x_g) = e_0(G_0(x_g)) \oplus e_1(G_1(x_g)) \oplus \dots \oplus e_{M-1}(G_{M-1}(x_g)) \quad (3)$$

where \oplus is the concatenation operator. The result is a feature vector $y = e(x_g)$ with $C \cdot M$ features, where C is the number of channels and M is the number of grids.

4.1.2 Decoding

Our decoder d is a shallow 2-layer MLP m with no bias terms used and 64 neurons per layer, which is lightweight and efficient. After each layer besides the last layer, ReLU activation is used as a nonlinearity. We use the Tiny-CUDA-NN package [19] for tensor-core accelerated decoding, which reduces training time by about 10% compared to a pure PyTorch decoder.

Scaling. Unlike other data formats that SRNs represent such as images, radiance fields, and distance fields, scientific data values are generally not bounded, posing a challenge for the decoder to accurately represent the unknown data range. We tackle this with a preprocessing step before training that identifies the minimum and maximum values for the volume the network will represent and saving those along with the network. Then, the decoding is calculated as $d(y) = m(y) \cdot (\max - \min) + \min$ where m is the MLP in the decoder, \min is the saved minimum value, and \max is the saved maximum value. This formulation allows us to avoid numerical issues by scaling the network output, which results in the MLP learning to output values between 0 and 1. More specific scaling for specific data types, such as z-score-, log-, and exponential-scaling, is not tested in this paper and is left to future work.

4.2 Learning feature grid transformation matrices

The adaptivity of our model stems from the learnable transformations for each feature grid in the encoder. The goal during training is to localize feature grids on/around regions that have high error. We focus more network parameters to those spatial regions with high error in order to improve reconstruction quality in that region.

While the idea is intuitive, naively attempting to learn the transformation matrices defining each grid is not possible using only a reconstruction loss (such as L1 or MSE). This is a direct consequence of our encoding ([Equation 2](#)), which returns a 0 feature if a query point is outside of a grid. This means the only gradients that will update a transformation matrix are from points that reside within that grid, so high error regions outside of a feature grid will not “pull” the feature grid toward it as desired.

To properly train the feature grid's transformation matrices, we introduce a feature density loss, which is a quantifiable and differentiable metric for the difference between the current feature density ρ , and a derived *target feature density* ρ^* . We use the term *feature density* to describe the average number of features (from the feature grids) that exist per unit volume at some point in space. In this section, we will discuss how we calculate the current feature density in a differentiable manner, how we derive a *target feature density* that warps the current feature

density to increase where the error is relatively higher and decrease where error is relatively lower, and a loss function to quantify the error between the current and target feature density which is ultimately used to update the transformation matrices during training.

Feature density. The simplest way to represent feature density is by the number of grids that overlap a spatial position x . The feature density ρ at a global location x using this formulation can be described by the following equation:

$$\rho(x) = \sum_{i=0}^{M-1} \begin{cases} 1 & \text{if } -1 \leq G_i(x) \leq 1 \\ 0, & \text{otherwise} \end{cases}, \quad (4)$$

where $G_i(x)$ is the function that transforms global coordinate x into feature grid i 's local coordinate frame. The if statement is true when it holds for each of the three dimensions of $G_i(x)$. However, this approach does not accurately represent the *density* of features at that point. Each feature grid has the same resolution, but a feature grid with small extents will have more dense features than a feature grid with near global extents. Instead, a better feature grid formulation is:

$$\rho(x) = \sum_{i=0}^{M-1} \begin{cases} \det(G_{i,0:3,0:3}) & \text{if } -1 \leq G_i(x) \leq 1 \\ 0, & \text{otherwise} \end{cases}, \quad (5)$$

where $G_{i,0:3,0:3}$ is the top-left 3×3 matrix of G_i representing scale, shear, and rotation for the grid, and $\det(\cdot)$ is the determinant of a matrix. As the grid becomes more dense (smaller in the global domain), the determinant above increases, representing the true feature density of the grid.

Though an accurate description of the feature density of the encoder, this formulation suffers from a large drawback in that there is no gradient to change the transformation matrices G when a point is outside of a grid, as the gradient of ρ will be 0 everywhere. This is a critical component for learning the transformation matrices in our method, and without it, the grid's scales, translations, etc., cannot update based on the error of points outside of their local domain.

Therefore, we use a gaussian approximation of the feature density, which is differentiable everywhere. Specifically, we use a class of gaussian functions called *flat-top gaussians* or *super-gaussians*. A flat-top gaussian is the same as the normal gaussian equation with a p term in the exponent:

$$g(x, p) = A \exp\left(-\left(\frac{(x-\mu)^{2p}}{2\sigma^2}\right)\right), \quad (6)$$

where A is a normalizing coefficient, x is the location, μ is the center of the distribution, σ is the standard deviation of the distribution, and p is the strength of the flat-top. As p increases, the gaussian shape become more box-shaped, depicted in [Figure 4](#).

Using the flat-top approximation for feature density box function, we give the final form of $\rho(x)$ as:

$$\rho(x) = \sum_{i=0}^{M-1} \left(\det(G_{i,0:3,0:3}) \cdot \exp\left(-\left(\sum_{d=0}^2 (G_i(x)_d)^{2p}\right)\right) \right), \quad (7)$$

where $G_i(x)_d$ is the transformed coordinate in the x-, y-, and z-axis for $d = 0, 1, 2$. Since the center of the local coordinate space is 0 and the standard deviation is 1, the gaussian equation does not require the μ or σ terms in the transformed space. This formulation is differentiable everywhere, facilitating proper training. Experimentally, we find $p = 10$ to be strong enough to match the box shape of the feature grid without being too strong so as to have exploding gradients on the ramps to and from the flat-top.

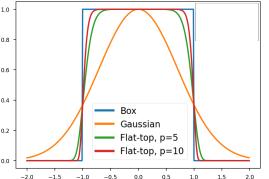


Fig. 4: Visualizing how flat-top gaussians (defined in [Equation 6](#)) with increasing flat-top power p generates a curve that is closer to a box, useful for approximating feature density in a differentiable manner.

Target feature density. With a way to calculate our current feature density at an arbitrary point, we now require a target feature density to steer toward during training. Our goal is to increase $\rho(x)$ where the error is relatively high and decrease $\rho(x)$ where error is relatively low. We formulate a target feature density for a coordinate $\rho^*(x)$ according to the following:

$$\rho^*(x) = \exp\left(\frac{\bar{h}}{h(x)+\epsilon} \log(\rho_{\text{scaled}}(x) + \epsilon)\right) \quad (8)$$

where $h(x)$ is the model’s error at location x , \bar{h} is the average model error (in practice, the average error over a batch), $\rho_{\text{scaled}}(x)$ is $\rho(x)$ divided by its sum over the set (to rescale the density between 0 and 1), and ϵ is a small number to avoid dividing by zero and other numerical issues. In language, Equation 8 will logarithmically scale the current (scaled) feature density according to the *relative error* at each position. If the error is exactly the average, then the target feature density is unchanged with respect to ρ_{scaled} . This logarithmic scaling assures that more importance is placed on increasing feature density in regions with very large relative error, and decreasing feature density in regions with very small relative error.

Feature density loss function. To quantify the difference between the current feature density and the target feature density, we view this as a distribution-matching problem. We calculate our feature density loss over a set of coordinates X as

$$\mathcal{L}_{\text{density}} = \frac{1}{|X|} \sum_{x \in X} \rho_{\text{scaled}}(x) \log\left(\frac{\rho_{\text{scaled}}(x)}{\rho^*(x)}\right), \quad (9)$$

where $|X|$ is the number of coordinates in X . The loss function measures the relative entropy from the target feature density to the current feature density, and is also known as the Kullback–Leibler divergence (KL divergence). This loss function provides gradients that can effectively update the parameters of each feature grid to match the calculated target feature density.

4.3 APMGSRN training

In this section, we cover the losses, training routine, and initialization we use for APMGSRN.

4.3.1 Loss functions.

We use two loss functions while training - a reconstruction loss \mathcal{L}_{rec} and the density loss $\mathcal{L}_{\text{density}}$ described in Equation 9. The reconstruction loss is the mean-squared error (MSE) between the model output and the ground truth data $\mathcal{L}_{\text{rec}} = \frac{1}{|X|} \sum_{x \in X} (f(x) - v(x))^2$, where X is a batch of coordinates, $|X|$ is the number of coordinates, f is the SRN, and $v(x)$ returns the ground truth data value at spatial position x via trilinear interpolation.

4.3.2 Training routine.

Each training step has two parts: (1) update the feature grid and decoder parameters using \mathcal{L}_{rec} , (2) update the transformation matrices using $\mathcal{L}_{\text{density}}$. Specifically, $\mathcal{L}_{\text{density}}$ is the only loss providing gradients to update the transformation matrices G , and \mathcal{L}_{rec} is the only loss providing gradients to update the feature grid values F and decoder weights. As long as the model is training, step (1) will happen every iteration. However, performing step (2) during each iteration would slow down training and potentially negatively impact reconstruction quality for reasons discussed in the following paragraphs. Therefore, we use a few techniques to minimize the number of step (2) updates.

Delayed start. Since $\mathcal{L}_{\text{density}}$ is dependent on the error for the current training iteration \mathcal{L}_{rec} , the update to the transformation matrices will be more useful when the model has already seen enough to know where the low- and high-error regions will be. When the model begins training, the error distribution may be random based on the network initialization, grid initialization, and data. Therefore, we delay the training for the transformation matrices until the rest of the network has been briefly trained on the data so as to remove most of the noise from initialization. We find that 500 iterations (with our learning rate and

network sizes) is enough to provide clear details about what regions of the volume will be hard to learn, so the transformation matrices are frozen for the first 500 iterations, and are updated according to $\mathcal{L}_{\text{density}}$ after that.

Early stopping. Our formulation presents a challenge for the neural network parameters in our APMGSRN in that if both training step (1) and (2) occur simultaneously, the model is training to chase a moving target. The parameters are updated from \mathcal{L}_{rec} with the assumption that the feature grids are static, but the grids are moving with each update step (2). Fixing the feature grids as early as possible is essential for the network to fine-tune the feature grid parameters to their location in space, as well as to reduce training time spent updating the transformation matrices when they may have already converged.

While training the transformation matrices, we keep a running track of $\mathcal{L}_{\text{density}}$ each iteration, and set an early stopping flag when the 1000-iteration moving average of $\mathcal{L}_{\text{density}}$ has not reduced by 0.01%. Alternatively, if this early stopping criteria is not met after 80% of the total training steps, we stop updating the transformation matrices in order to allow the other network parameters to learn given fixed feature grid positions.

A single APMGSRN model takes at most 4 minutes to train for 50k iterations with the hyperparameters experimented with in this paper. Often, the grids converge quickly, reducing training time to as short as 40 seconds.

4.3.3 Initialization

With feature grid positions being crucial to the performance of our model, the initialization of the transformation matrices is a relevant factor in training speed and model accuracy. We initialize our transformation matrices to cover a near global domain with small random shears, rotations, and translations. The diagonals of the matrices (representing the scale of the grid) are sampled from $\mathcal{N}(1, 0.05)$, while the remaining entries in the first three rows of the transformation matrix are sampled from $\mathcal{N}(0, 0.05)$. This initialization assures that each (relevant) entry in the transformation matrix is non-zero so that it may contribute to the final output, guaranteeing gradients can be used to update each entry of the matrix. Additionally, the domains created with this initialization scheme are all nearly equal to the global extents, and so while fixed at the beginning of training (see Section 4.3.2), each grid is initially helping learn the global domain, providing a better approximation of which regions will be challenging to learn.

Besides our transformation matrices, we initialize our feature grids from $\mathcal{U}(-0.0001, 0.0001)$, encouraging near-zero initial guesses with small randomness as following Müller et al. [20]. Our decoder network weights are initialized following Glorot and Bengio [5].

5 DOMAIN DECOMPOSITION SCENE REPRESENTATION

It may not always be feasible to train a single model for some large-scale data for two reasons. From one end, as training data become more complex with higher resolutions, larger neural networks are necessary to obtain adequate reconstruction accuracy, which means longer training times and a large model footprint on disk. As model complexity increases, GPU memory may become a bottleneck during training or training time becomes unacceptably long. From the other end, the large-scale data we wish to model with an SRN may not fit within GPU memory for efficient query in the training loop. The main CPU memory can be used to host data and support on-demand data transfer to the GPU for training, but is costly due to the random sampling of data points during training. If the CPU memory also cannot support hosting the data, then an out-of-core solution is the only option available. Both options have been explored by Wu et al. [32].

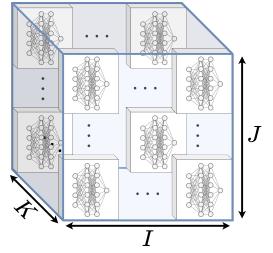


Fig. 5: A depiction of a domain decomposition of networks, which cover the domain in a grid pattern.

Instead of a complicated and inefficient out-of-core sampling method, we take a model-parallel approach to modelling a large-scale volume. We propose that instead of a single network representing the volume, we use a domain decomposition of models, where each model has learned a separate brick of the volume. The models are arranged in a grid, and do not require communication during training or inference. We discuss the data partitioning in [Section 5.1](#), the training procedure in [Section 5.2](#), and querying a set of models during inference in [Section 5.3](#).

5.1 Data partitioning

To partition data into bricks for one network to be trained per brick, we use a grid of networks as shown in [Figure 5](#). Ahead of training, grid resolution I, J, K must be chosen, representing the number of networks for the width, height, and depth of the volume, respectively. We recommend picking an I, J, K such that the resolution of the grid assigned to each network is roughly the same aspect ratio as the feature grids in the model. For instance, if the feature grids are 32^3 , choose I, J, K such that each model is assigned a roughly cube-shaped volume.

Ghost cells. In order to mitigate boundary artifacts along adjacent network boundaries during visualization, we experiment with the use of a ghost layer of cells, which are cells that overlap between networks. During data partitioning, this means that the extents of a network will get extended by some number of ghost cells along each axis. This solution does not guarantee that seams will be removed, but tends to reduce the effect of them, as the networks should have less of a difference along boundaries if they are learning beyond the actual extent. In our experiments in [Section 7.2](#), we test with between 1 and 16 ghost cells, for example. We discuss a more in-depth solution in our future work, but do not go beyond ghost cells in this paper.

5.2 Domain decomposition training

To reduce the training time needed for the $I \cdot J \cdot K$ models, we dissect the domain and train the models in parallel across available GPUs. Since large-scale data are often generated by powerful machines, we assume a user is likely training on a compute node with multiple GPUs. Our domain decomposition training routine pre-calculates the extents of each network within the volume’s extents (including ghost cells) and generates a list of jobs that are assigned to the available GPUs. When a GPU finishes training, the GPU returns to the available GPUs list, where it waits to be issued the next model to train. When a GPU begins training, an APMGSRN model is initialized, and only the data within its assigned extents are loaded to the GPU memory. This reduces I/O overhead and is efficient with data formats such as NetCDF supporting parallel file access and arbitrary cropping from disk.

Early stopping. Just as we employ early stopping to stop learning the transformation matrices earlier during training to converge faster, we also use early stopping on models within the domain decomposition while training so the GPU is freed for the next model to train quicker. We use a plateau learning rate scheduler that detects when \mathcal{L}_{rec} has not decreased by 0.01% over 500 iterations. When this is triggered, the learning rate of the model’s parameters is reduced by a factor of 10. When this is triggered 3 times, we finish training the model. This can reduce training times dramatically in regions where there may be very sparse data (down to 20 seconds per model in our experiments).

5.3 Domain decomposition inference

After training (presumably on a remote server), the neural networks can be moved to a local workstation for visualization. Once on a local machine, all models are loaded into GPU memory. Since each network in the domain decomposition represents a certain spatial extent, each global query location needs to be mapped to the correct model for inference, so we develop a spatial hashing function for efficient inference.

Spatial hashing function. To efficiently determine what model a query location belongs to, we implement a spatial hashing function to hash input 3D global coordinates to the correct index in the array of models. We load models into an array in x-,y-,z-dimension order. With this order in mind, we can hash input global coordinates directly to array

index. Assuming a global spatial coordinate p is in the domain $[-1, 1]^3$, we calculate the grid index for p with $i = \left\lfloor \frac{P_x+1}{2} \right\rfloor, j = \left\lfloor \frac{P_y+1}{2} \right\rfloor, k = \left\lfloor \frac{P_z+1}{2} \right\rfloor$, where $\lfloor \cdot \rfloor$ is the integer floor operation. If any dimension of the coordinate has exactly 1.0 for an axis, then the result is an out of bound location, so these values need to be clamped such that $i < I, j < J, k < K$. Since the networks are stored in a list, the 3D i, j, k index for the network is flattened in C-order, giving the final hash as $i + Ij + IJk$. With the correct network determined to perform inference with, the global coordinate can be scaled to the networks local domain for proper inference.

Our implementation loops over all networks to do evaluations for each, but a significant speedup could be achieved with further engineering effort, as shown by Reiser et al. [23], who develop custom CUDA code for inference within their set of networks for radiance fields. We leave this to future work, but expect the inference time of large domain decomposition models to be much quicker when inference code is written in custom CUDA code that parallelizes network evaluation.

6 NEURAL VOLUME RENDERING

For fair comparison between the models, we develop a Python/PyTorch [21]-based neural volume renderer (included with our code on GitHub) that supports our own APMGSRN as well as state-of-the-art models fVSRN [30] and neural graphics primitives (NGP) using a hash-grid encoding [20], on top of raw data volume rendering. Any PyTorch model that can support mapping 3D coordinates to an output density can also be plugged in with minimal reconfiguration. We favor ease of use over raw performance with this renderer with a plug-and-play style coding requiring minimal changes to a user’s PyTorch model, and hope it is useful for future research on SRNs for sci-vis. Our renderer supports arbitrary transfer function and view direction, and uses progressive rendering to support 60+ fps rendering for immediate feedback while panning, rotating, and zooming. We support transfer functions exported from ParaView directly to our renderer, and interacting with the scene follows the typical 3D viewer paradigm of clicking and dragging rotating around the scene, the scroll wheel zooms, and middle-mouse clicking and dragging pans the camera. We use 3D rendering helper functions from the nerfacc Python package [10], which provides CUDA-accelerated ray marching and compositing.

For improved interactivity, we implement a progressive rendering scheme that renders the image in a checkerboard pattern, evaluating the pixels in order of an image hierarchy. After each progressive rendering pass, the in-progress image displayed on the screen is a bilinearly-upscaled version of the current finest-available image from the image hierarchy masked with the fully evaluated pixels. We also offer a trade-off between interactivity and full render time with a user choice for batch size, with large batch sizes slowing down interactivity but speeding up full-frame render times, and vice versa for small batch sizes. Our progressive rendering keeps GPU memory minimized and quickly generates good approximations of the final image while offering interactive frame rates (30-400 fps).

7 EXPERIMENTS

We perform experiments to test the training and reconstruction metrics of our proposed APMGSRN model in [Section 7.1](#). Our domain decomposition training approach is evaluated in [Section 7.2](#). Lastly, we visualize trained models with our renderer in [Section 7.3](#). Our datasets and brief descriptions are shown in [Table 1](#).

Hyperparameters and evaluation hardware. To save page space, model hyperparameters for each size SRN (small/medium/large/decomposition) are available in GitHub in the “BatchRunSettings” folder, as there are many hyperparameters to enumerate for each architecture. For our model specifically, the number of grids we use is either 16 or 32, the grid resolution varies from 16^3 to 64^3 , and the number of features per feature grid vertex is either 1 or 2. We find that inference speed degrades above 64 grids, and accuracy does not improve significantly with more grids or more features per grid to justify the increased storage cost. Grid resolution has the largest impact on both network performance and network storage size.

gaussian can make up for this for some cases, but that may cause larger-scale grids to see unstable updates to their parameters. For best results, we recommend our default initialization of global-scale grids with small random perturbations, but we believe there may be room for improved performance with different differentiable approximations of grid density for smoother learning. We recommend readers view our supplemental material for videos of feature grids during training on the smaller datasets experimented with, as it gives a strong intuition for what the network is learning and how it is learning it. In the future, other grid initialization strategies may prove useful for better learned grid positions. Additionally, we believe that scaling up the flat-top gaussian strength (starting from 1.0) during training can help situations like our small initialization scheme, because gradients will not approach 0 rapidly for data outside the grid.

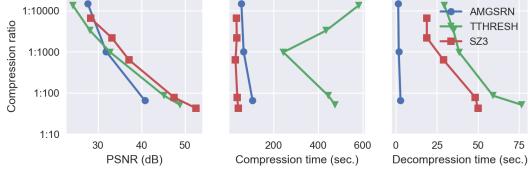


Fig. 7: Comparison of APMGSRN with TTHRESH and SZ for compressing the isotropic volume (1024^3). Our "compression" time is the time it takes to train the model, and our decompression time is the time it takes to query each point in the volume.

Compression. Though data reduction is not the only intended use case for our approach, nor is it something we design for specifically (such as with network weight quantization for further data reduction [16]), we believe it is useful to the community to compare the compression ability of state-of-the-art SRNs with state-of-the-art compressors. We compare compression results from TTHRESH [2] and SZ3 [12, 13, 34] with our approach in Figure 7. We do not compare to a rendering-focused compressor such as cudaCompress [28] or a bricked version of TTHRESH [2], as this has been done by recent work [30] showing that its decoding is significantly slower, compression rates smaller, memory use is higher, and image quality is worse than state-of-the-art SNRs.

As shown in Figure 7, our approach only provides a higher compression ratio than TTHRESH or SZ3 at a $10000\times$ data reduction rate with a reconstruction quality of under 30 dB PSNR, and is otherwise less compressive than state-of-the-art compressors. The throughput (decompression) from APMGSRN is much quicker regardless of compression level, decoding on the order of 500 million points per second, but keep in mind that APMGSRN is decoding on the GPU whereas SZ3 and TTHRESH are on CPU. The largest difference between APMGSRN and these compressors is that our approach can efficiently perform arbitrary point evaluations, whereas SZ3 and TTHRESH both require decompressing to the original data size before trilinear interpolation.

We also test compression on our two large datasets, rotstrat and channel. The tested compressors take significantly longer to compress our large-scale data, or fail to do so at all. In fact, TTHRESH and SZ3 both run out of memory on our machine with 1TB of memory when trying to compress the channel dataset (450GB), and TTHRESH also runs out of memory when trying to compress the rotstrat dataset (250GB). SZ3 successfully compresses the 250GB rotstrat data in 38 minutes, resulting in a compression ratio of $509\times$ and a reconstruction quality of 20.60 dB PSNR after another 14 minutes of decompression. Both our small and medium sized (domain decomposition) models for rotstrat (see Table 1) achieve higher PSNRs (41.32 dB and 44.29 dB) with higher compression rates ($4791\times$ and $1432\times$) while only needing to train for 34 minutes if all models are trained sequentially on 1 GPU, or 5 minutes in parallel on 8 GPUs. Additionally, our decoding takes 3 minutes for the whole 4096^3 volume.

7.2 Domain decomposition network evaluation

We evaluate our domain decomposition training approach on each dataset, and show results in Table 1. Even though the smaller datasets do not require the domain decomposition to train without memory

limitations, the domain decomposition model outperforms a single model at the same storage size for each dataset. This is expected, as each of the domain decomposition models will have their own decoder and set of feature grids that may create and advantage over the single large model. The downside of a domain decomposition model is that training may take longer if you only have 1 GPU and train serially.

For data that cannot fit in a single GPU, such as rotstrat and channel, our domain decomposition approach provides an efficient alternative to an out-of-core training routine. In Table 1, we list the *average* training time per model for the domain decomposition models. Our training machine had 8 GPUs training in parallel and the total training time was 7 minutes with data I/O for rotstrat.

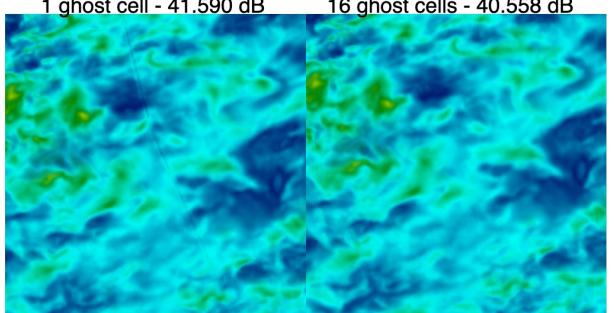


Fig. 8: Zoom in of volume render from identical $4 \times 4 \times 4$ domain decomposition networks trained with either 1 or 16 ghost cells. A seam is clear with 1 ghost cell.

Effect of ghost cells. By default, our domain decomposition has 1 ghost cell for corner alignment across models. We show the effect of increasing ghost cells within our domain decomposition models in Figure 8, which is a volume rendering of two networks - one trained with the default of 1 ghost cell, and one with 16 ghost cells. The boundary artifact is visible only in the render with only 1 ghost cell. However, a significant PSNR drop of -1 dB is noticed with this increase in ghost cells, attributed to the fact that a network in the decomposition learning a volume of size 288^3 is a 42% larger volume than the original 256^3 volume learned on in the version without ghost cells. The redundancy of learning a significant portion of the full volume in multiple networks reduces the learning capacity of the set of networks. Since the artifacts are fairly minor while reducing performance, we recommend using 1 or few ghost cells for general volume rendering, but 16 ghost cells for tasks like isosurface extraction where seams can be very distracting. We consider future work for reducing the boundary effect between networks with a network communication scheme during training.

7.3 Neural volume rendering

In this section, we evaluate the single frame performance of neural volume rendering. The raymarching and compositing are done with the same code, so the timing and memory use difference is accounted for solely by querying the SRN or raw volume data.

Our single frame rendering results are shown in Figure 9 where we compare the render results of the large sized (64 MB on disk) models for the supernova dataset. The reference render using the volume is the quickest, but assumes the volume can fit in GPU memory, which is not possible for datasets larger than 11 GB on our test machine. NGP provides the quickest render times of the networks on our 2080Ti, as well as the lowest memory use. fVSRN is slower likely due to the Fourier frequency encoding. Our models are slowest to render, but have considerably better image reconstruction shown by the higher PSNR and structural similarity index measure (SSIM) [35]. Our approach (single model) requires the most memory due to the transformation to the 64 local grid spaces, requiring $64\times$ the memory of the global coordinates for the local space coordinates.

In Figure 10, we compare our trained “medium” sized domain decomposition model for the 250 GB rotstrat data with a downsampled version of the raw data with a similar storage footprint as the model. We subsample the raw data by 11 in each dimension to create a volume that takes 198 MB to store, which is only slightly larger than our network

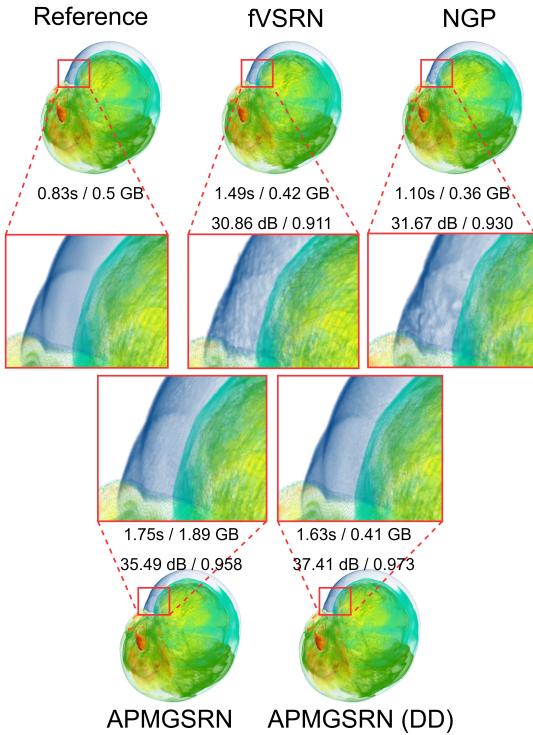


Fig. 9: 1024^2 volume renders of Supernova data (432^3) with a step size of 1 voxel. A batch size of 2^{23} is used during forward passes. The reference is the raw data (308 MB), while the others are neural representations of the data in their “large” configuration (64 MB), where (DD) means domain decomposition. The first row of metrics listed is render time / memory use, and the 2nd row is PSNR (dB) / SSIM. Metrics are each for the entire image, not the zoom in.

which takes 183 MB to store. We use our renderer to generate 8000^2 images of each representation of the data with the same visual mapping parameters. The image from our model achieves a significantly higher 30.70 dB PSNR compared with the render from the downsampled data which only achieves 22.24 dB. Additionally, the neural network captures features that are entirely missed by the downsampled data.

We include an image of our rendering application in Figure 11, but do not discuss it in detail here. Please see supplemental materials for a video using the tool in real-time, as well as additional qualitative comparisons of volume renders from each model.

8 LIMITATIONS

One limitation of our model is that memory use blows up and training speeds reduce as the number of feature grids increases. This is from the global to local transformation, which needs to allocate memory equal to the memory of the query coordinates times the number of feature grids, and then perform interpolation for each of those coordinates. Further engineering effort could reduce the overhead with custom CUDA code or by limiting the number of grids operated on at once.

Though our neural renderer is compatible with any PyTorch model coded in Python, the renderer is slower than what is reported by contemporary work [30–32]. Further optimizations to the code could improve the rendering speed while remaining flexible for new PyTorch models to be plugged in with ease.

9 CONCLUSION AND FUTURE WORK

In this work, we present a novel SRN called APMGSRN, a model which excels at representing scientific data due to the adaptivity of the model. In order to train on datasets that cannot fit in-core during training, we propose an domain decomposition approach for training and inference that divides the domain into separate chunks to be learned

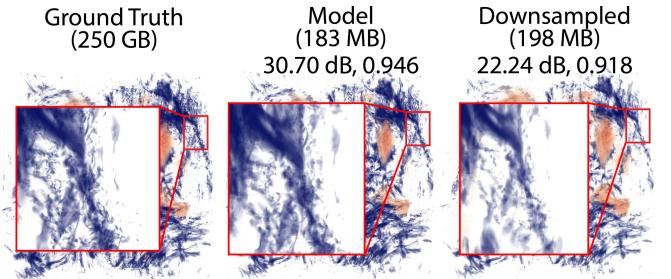


Fig. 10: Volume renders of the original rotstrat data, our “medium” size model representing the data, and a low resolution version of the original data, which has been downsampled by 11 in each direction to have a similar storage footprint as the model. Each image is rendered at 8000^2 with the same visual mapping parameters. Metrics listed are PSNR/SSIM for the entire image (not just zoomed portion).

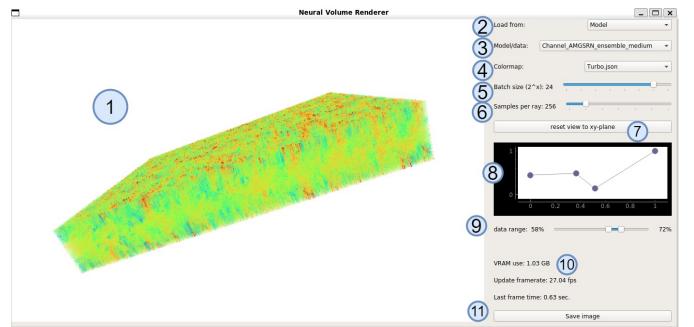


Fig. 11: The interface for our volume renderer with a volume render of an domain decomposition model trained on the 450GB channel dataset. (1) The viewing area for interaction. (2) Dropdown to choose to load a model or data (NetCDF). (3) Which model/data to load from. (4) Colormap to use. (5) Batch size when querying the network. Larger reduces total frame time, but requires more GPU memory and slows down interactivity. (6) Samples per ray. (7) Resets camera view on data. (8) Opacity transfer function editor. (9) Relative min/max for transfer function. (10) Render statistics. (11) Save current frame button.

by multiple networks. In our evaluations, we find that our proposed APMGSRN architecture outperforms the reconstruction performance of other state-of-the-art models at similar model sizes in both data space and image space. Our model cannot compress as well as state-of-the-art compressors like TTHRESH, but achieves high compression rate with efficient volume rendering in the compression domain without the need to decompress to the original data’s size.

In the future, we think our adaptive feature grids can be augmented in two ways: (1) the feature grids need not be trilinearly interpolated – another encoding scheme, such as a hash grid from NGP, can be used for query within each grid, (2) the feature grids could extend the time dimension, creating 4D learnable feature grids supporting temporal interpolation. Reducing boundary artifacts is another key direction to improve upon to reduce distracting seams in visualizations. On the engineering side, there are many places where optimizations can improve efficiency, such as the multi-grid encoder or domain decomposition model inference. Our approach could also work with minimal changes to represent non-uniform grid data types, such as unstructured, point-based, or AMR, so long as the data can return a scalar value for an arbitrary point query. With the increasing popularity and usefulness of SRNs, future work to support neural volume rendering in mature open-source software such as ParaView would greatly benefit the visualization community.

