# Learning Adaptive Sampling and Reconstruction for Volume Visualization

Sebastian Weiss (iD), Mustafa Işık (iD), Justus Thies (iD), and Rüdiger Westermann (iD)
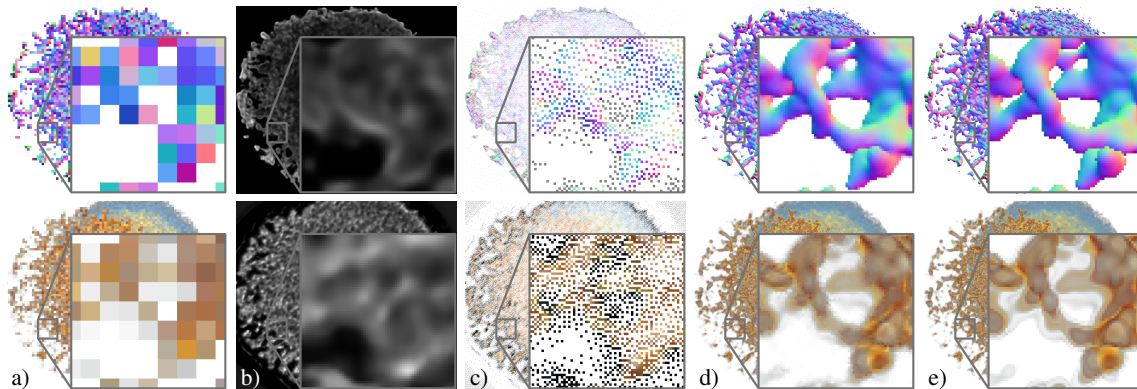
Fig. 1: An importance network, together with a differentiable sampler and a reconstruction network, takes a low resolution visualization (a) and infers an importance map (b) from it. From this map, an adaptive sampling pattern with adjustable number of samples (5% for iso, top; 10% for dvr, bottom) is derived, and a volume ray-caster samples the data according to these samples (c). The reconstruction network completes the visual representation from the sparse set of samples (d). The ground truth visualizations are shown in (e). The proposed network pipeline works on images of iso-surfaces (top) and direct volume renderings (bottom).

**Abstract**—A central challenge in data visualization is to understand which data samples are required to generate an image of a data set in which the relevant information is encoded. In this work, we make a first step towards answering the question of whether an artificial neural network can predict where to sample the data with higher or lower density, by learning of correspondences between the data, the sampling patterns and the generated images. We introduce a novel neural rendering pipeline, which is trained end-to-end to generate a sparse adaptive sampling structure from a given low-resolution input image, and reconstructs a high-resolution image from the sparse set of samples. For the first time, to the best of our knowledge, we demonstrate that the selection of structures that are relevant for the final visual representation can be jointly learned together with the reconstruction of this representation from these structures. Therefore, we introduce differentiable sampling and reconstruction stages, which can leverage back-propagation based on supervised losses solely on the final image. We shed light on the adaptive sampling patterns generated by the network pipeline and analyze its use for volume visualization including isosurface and direct volume rendering.

**Index Terms**—Volume visualization, adaptive sampling, deep learning.

✦

## 1 INTRODUCTION

W HICH are the data samples that are needed to generate an image of a data set that conveys the relevant information encoded in this data? This question is fundamental to data visualization, since it asks for the importance of data samples from a perceptual point of view, rather than a signal processing standpoint that argues in terms of numerical accuracy.

Recent works in visualization have shown that artificial neural networks can perform an accurate reconstruction from a reduced set of data samples, by learning the relationships between a sparse, yet regular input sampling and the high-resolution output. Learned representations are then applied in the reconstruction process

to infer missing data samples. This type of reconstruction has been performed in the visualization image domain to infer high-resolution images from given low-resolution images of isosurfaces [60], in the spatial domain to infer a higher resolution of a 3D data set from a low-resolution version [64], and in the temporal domain to infer a temporally dense volume sequence from a sparse temporal sequence [21].

Others have even proposed neural networks that are trained end-to-end to learn directly the visual data representations instead of the data itself. Berger *et al.* [3] propose a deep image synthesis approach to assist transfer function design, by letting an artificial neural network synthesize new volume rendered images from only a selected viewpoint and a transfer function. He *et al.* [23] demonstrate that artificial neural networks can even be used to bridge the data entirely, by learning the relationships between

---

• All authors are with Technical University of Munich, Germany.
E-mail: {sebastian13.weiss,m.isik, justus.thies, westermann}@tum.de.

the input parameters of a simulation and visualizations of the simulation results. Both approaches do not make any explicit assumptions about the relevance of certain structures in the data, yet the learned relationships between parameters and visual representations are considered in the image generation process.

## 1.1 Contribution

Our goal is to make a further step towards learning visual representations, by investigating whether a neural network can a) learn the relevance of structures for generating such representations, b) use this knowledge to adaptively sample a visual representation of a volumetric object, and c) reconstruct an accurate image from the sparse set of samples. Notably, even we can demonstrate for very large volumes and image sizes that adaptive sampling can save rendering time, performance improvement is not our main objective. It is even fair to say that an optimized GPU volume ray-caster can hardly been beaten performance-wise. Our main objective is to gain an improved understanding of the learning skills of neural networks for generating visual representations in an unsupervised manner, by letting networks learn the relevance of certain structures for obtaining such representations. It can eventually become possible to generate data representations that compactly encode relevant structures in a way they can be used by a neural network to visualize the data. Such insights can further facilitate the use of transfer learning to construct synthetic data sets that contain the structures that are important for successful learning tasks on real data. For viewpoint selection, a network might learn to recommend views showing many important structures, and for training this information can be used to acquire more data from similar views.

To address our objectives, we introduce a novel network pipeline that is trained end-to-end to learn the relevance of certain structures in the data for generating a visual representation (Figure 1). This pipeline is comprised of two consecutive internal network stages: An importance network and a reconstruction network. Both networks work in tandem, in that the first learns to place samples along relevant structures by using the second network to give feedback on how well a visual representation of the data can be reconstructed from the sparse sampling. Our approach differs from previous adaptive sampling approaches in volume visualization [37, 31, 2] in that it does not rely on any specific saliency model to determine the image regions that need to be refined. In contrast, we propose a network-based processing pipeline that simultaneously learns where to sample and how to accurately reconstruct an image from the sparse samples, solely using losses on the reconstructed images.

For learning an importance map from a low-resolution visualization and reconstructing an image from a sparse set of pixel values, we use two modified versions of an EnhanceNet [53]. To enable network-based learning using gradient descent, two novel processing stages are introduced:

- A differentiable sampling stage that models the relationship between sample positions and visual representation.
- A differentiable image reconstruction stage using the pull-push algorithm [15, 32] to model the relationship between a sparse set of image samples and the reconstructed image.

In a number of experiments, we demonstrate that the importance network effectively selects structures that are relevant for the final visual representation. We focus on adaptive sampling in image-space, i.e., using surface samples and samples resulting from direct volume rendering. As a future direction of research, we outline adaptive sampling in object space, i.e., using data samples along view-rays. Our experiments include qualitative and quantitative evaluations, which indicate good reconstruction accuracy even from few samples. The source code of our processing pipeline is available at https://github.com/shamanDevel/AdaptiveSampling, including some of the data sets that have been used for training and validation.

## 2 RELATED WORK

In the following, we review previous works that share similarities with our approach from the fields of adaptive sampling for rendering as well as neural network-based image and volume reconstruction.

**Adaptive Sampling for Rendering** Adaptive rendering has a long tradition in computer graphics, to reduce the number of rays to trace against the scene and perform rasterization at lower image resolution. At the core of such approaches is the computation of importance values to steer the adaptive refinement, for instance, based on perceptual models [5, 42, 48], image saliency models using pixel variance [44, 49], image difference operations [40], or entropy-based measures [61], to name just a few. In the context of foveated rendering [17], where usually a static adaptive sampling pattern is used that moves with the users gaze, a luminance-contrast-aware criterion was introduced to enable feature-aware adaptivity [58]. The importance map generation process is often started from an image preview that is calculated using a low resolution render pass or a high-resolution estimate that can be created in a significantly faster way than the final image.

For volume rendering, a number of approaches have investigated adaptive sampling in object space, to reduce the number of samples along the view rays [43, 9, 38, 6]. Adaptive image-space refinement has been proposed by Levoy [37], by using the color variances between pixels at low image resolution to decide whether to refine the image resolution locally. Kratz *et al.* [31] propose to use the difference image between two coarser resolution images, and locally refine where high differences are observed. Belayev *et al.* [2] render low-resolution images of isosurfaces and refine depending on how many pixels surrounding a pixel in the low-resolution view fulfill certain requirements. Frey *et al.* [13] use a fixed random sampling structures that is applied in a hierarchical manner to progressively refine the image.

The major differences between these approaches and our proposed sampling pipeline are as follows: Firstly, the pipeline learns to adapt the sampling in an unsupervised manner. A specific feature descriptor that steers the placement of samples is not used, and importance values are learned solely using losses on the reconstructed image. Secondly, the number of samples can be prescribed, which is not easily possible with existing schemes due to their pixel-iterative nature. Thirdly, the pipeline learns simultaneously the adaptive sampling and the image reconstruction from the sparse set of samples. In all previous schemes, the final interpolation step is decoupled from the sampling process.

**Deep Learning for Upscaling and Denoising** In recent years, deep learning approaches have been used successfully for single-image and video super-resolution tasks [11, 54, 55, 56, 52, 7], i.e., the upscaling of images and videos from a lower to some higher resolution. Many previous works let the networks learn to optimize for losses between the inferred and ground-truth images based on direct vector norms [29, 27]. GANs were introduced to prevent the undesirable smoothing of direct loss formulations [53,

36], and instead use a second network that discriminates real from generated samples and guides the generator. Convolutional architectures [11] with residual blocks [22] are popular generator architectures that offer training stability as well as high-quality inference. Losses based on the feature-space differences of image classification networks, e.g., a pre-trained VGG network [26], have shown to mimic well the human's capability to assess the perceptual similarity between two images.

The approach closest to ours is by Kuznetsov *et al.* [34] for learning adaptivity in Monte-Carlo path-tracing and denoising of the final image. A first network learns to adapt the number of additional paths from an initial image at the target resolution, which is generated via one path per pixel. A second denoising network learns to model the relationship between an image with increased variance in the color samples to the ground truth rendering [47, 41]. Conceptually, our approach differs in that it works on a low-resolution input map and then learns to freely position the sample locations in image space, i.e. it learns to place zero or one sample per pixel. This requires a completely different differentiable sampling stage, as well as a differentiable image reconstruction stage that can work on a sparse set of samples. Furthermore, Kuznetsov *et al.* use finite differences between images of different sample counts for gradient estimation. Incurring noise is reduced by averaging multiple samples with different sample counts, which is not possible in our approach where at most one sample per pixel is taken. Instead, we propose a sigmoid approximation that can be differentiated analytically.

In visualization, Zhou *et al.* [64] presented a CNN-based solution that upscales a volumetric data set using three hidden layers designed for feature extraction, non-linear mapping, and reconstruction, respectively. Han *et al.* [20] introduced a two-stage approach for vector field reconstruction via deep learning, by refining a low-resolution vector field from a set of streamlines. Berger *et al.* [3] proposed a deep image synthesis approach to assist transfer function design using GANs, by letting a network synthesis new volume rendered images from only a selected viewpoint and a transfer function. The use of neural network-based inference of data samples in the context of in situ visualization was demonstrated by Han and Wang [21], where a network learns to infer missing time steps between 3D simulation results. He *et al.* [23] use neural networks for parameter-space exploration, by training a network to learn the dependencies between visual mappings of simulation results and the input parameters of the simulation. Guo *et al.* [18] designed a deep learning framework that produces coherent spatial super-resolution of 3D vector field data. Weiss *et al.* [60] extent image upscaling to geometry images of isosurfaces including depth and normal information. Instead of data upscaling, Tkachev *et al.* [57] predict a next time-step of a simulation and identify regions of interest by high variance between the network prediction and the ground truth. Common to all these approaches is the use of a regular sampling structure that does not consider the importance of samples in the inference step.

# 3 LEARNING TO SAMPLE

In the following, we discuss how the importance network makes use of both the adaptive sampling stage and the reconstruction network to learn where to place samples with higher density. The importance network (subsection 3.1) receives an image of the data set at low resolution. This image $L$ is of shape $C \times fH \times fW$, where $W$ and $H$ denote the screen resolution, and $f$ the downsampling factor. This factor is set to $1/8$ in all of our experiments. Each image pixel is comprised of $C$ channels, such as color, depth, and normal, representing what is seen through that pixel. The network is trained to learn an importance function $\mathcal{N}_I$ that generates a gray-scale importance map $I \in [0, \infty)^{H \times W}$ in which low and high values, respectively, indicate where less or more samples are taken.

The sampler $\mathcal{S}$ (subsection 3.2) takes the importance map and places a given number of samples, e.g., 5% of the pixel, in the full resolution image $S \in \mathbb{R}^{C \times H \times W}$ according to the importance information. Only at these samples the object is rendered. The reconstruction network learns a function $\mathcal{N}_R$ (subsection 3.3) that reconstructs the final output $O \in \mathbb{R}^{C \times H \times W}$ from the sparse set of samples. We make the sampler differentiable w.r.t. sample positions to allow gradient flow from the reconstruction network (subsection 3.3) to the importance network, so that the reconstruction network is trained simultaneously and propagates the loss information to the sampling stage. Since the entire pipeline is trained end-to-end using a loss on the reconstructed and ground truth images, the importance network and the pair of sampler and reconstruction network work together in an effort to learn the placement of samples so that high reconstruction quality is achieved.

In principle, one can refrain from using a separate importance map, by realizing the sampler as a network that directly learns the adaptive sampling. In this case, however, modelling the positional information in a network requires to represent positions explicitly, either in a graph structure or a linear field, so that less efficient graph networks or fully-connected networks need to be used. Furthermore, the sampler has to be re-trained whenever a different number of samples is used. Our approach enables to use efficient convolutional networks, and to change the number of samples at testing time.

An overview of the processing pipeline is shown in Figure 2. It works with images comprised of an arbitrary number $C$ of channels. In the first part of this work, the pipeline is introduced for isosurface rendering with $C = 5$, i.e., a binary mask (1: hit, 0: no hit), a normal vector, and a depth value. The application to direct volume rendered images is discussed in section 6.

Weiss *et al.* [60] enforce frame-to-frame coherence during animations by including a temporal loss in the training step. This loss considers the difference between the previous frame – warped by the frame-to-frame optical flow – to the current frame. In the accompanying video, this approach is used for both the importance and reconstruction network. In the following discussion, however, temporal connections are omitted and the focus is solely on single image reconstruction for clarity.

## 3.1 Importance Network

The importance network $I \leftarrow \mathcal{N}_I(L)$ determines the distribution of the samples that are required by the reconstruction network to generate the visual output according to some loss function. Deeming every pixel equally important, i.e.,

$$\mathcal{N}_{I,\text{constant}}(L)_{ij} = \mathbf{1}, \tag{1}$$

leads to a uniform distribution of the samples [64, 21, 60]. Alternatively, and in the spirit of classical edge detection filters, the screen space gradients of the individual channels can be used, i.e.,

$$\mathcal{N}_{I,\text{gradient}}(L)_{ij} = \sum_c w_c ||\nabla L_{ij,c}||_2^2, \tag{2}$$

where $\nabla L_{ij,c}$ is the screen space gradient of channel $c$ at location $ij$. The contributions of the individual channels are weighted by
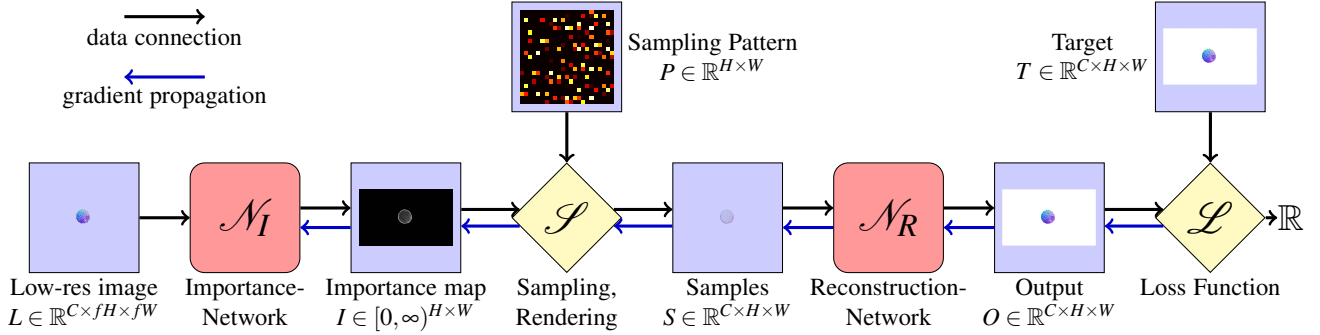
Fig. 2: Overview of network-based adaptive sampling. From a low-resolution image $L$, the importance network infers the importance map $I$. The sampler $\mathscr{S}$ uses this map together with a sampling pattern $P$ to adaptively place samples in the high-resolution image $S$. Ray-casting the object at these samples generates a sparse image. The reconstruction network recovers the dense output $O$.

$w \in \mathbb{R}^C$. Other known importance measures consider screen space curvature via the variation of surface normals [46], or color contrast via the variation of luminance [58].

Alternatively, we introduce a fully convolutional neural network $\mathscr{N}_{I,\text{net}}$ (subsection 4.2) that predicts a high-resolution greyscale importance map $I$ from a low-resolution rendering $L$. Notably, this network is not trained w.r.t. specific characteristics that are derived from the image like gradients or luminance information, since this requires to heuristically decide on the importance of pixels. Instead, it is trained end-to-end with losses only on the reconstructed color information, by gradient descend all along the processing pipeline. In section 5, network-based inference of the importance map is compared to alternative approaches, showing superior prediction of regions that are important for the final image.

## 3.2 Differentiable Sampling

Given the target number of samples in the final image, e.g. $\mu = 5\%$ of all pixels, the sampler uses the importance map $I$ to determine where to place these samples. To generate the given number of samples, two main classes of algorithm are commonly used in rendering:

- Stippling starts with a given number of points at random locations and iteratively optimize these locations so that the point density matches the density of the importance map [10, 16].
- Importance sampling treats the importance map as a density function and place samples via rejection sampling or the inverse cumulative distribution function [35, 1].

These algorithms, however, are not easily differentiable w.r.t. changes in the importance map, since they use discrete optimizations or random processes, and often are too slow for real-time applications. To make the sampling process differentiable and fast, we propose a sampling strategy that computes for every pixel independently the chance of being sampled. This is achieved by a smooth approximation of rejection sampling, which is differentiable and allows for gradient propagation through the network pipeline. Since every pixel can be processed independently, this scheme can effectively leverage parallel execution on the GPU. On the other hand, it does not allow for an exact match of the prescribed number of samples, yet produces a number of samples that slightly varies around the target number.

In a first step, the importance map $I$ is normalized to have a prescribed mean $\mu$ and minimal value $l \leq \mu$. Let $\mu_I$ be the mean of $I$ over all pixels, then the image

$$I'_{ij} := \min\left\{1, l + I_{ij}\frac{\mu - l}{\mu_I + \varepsilon}\right\} \tag{3}$$

has the desired properties. A small constant $\varepsilon = 10^{-7}$ is used to avoid division by zero. The minimal value $l$ is required to maintain a lower bound on the sample distribution in empty areas, which is important to allow for an accurate reconstruction in such areas. We use $l = 0.002$ in all of our experiments. Clamping to a maximal value of 1 is required by the following sampling step, which is realized as an independent Bernoulli process via rejection sampling, i.e., a sample at location $ij$ is taken if the probability $I'_{ij}$ is larger than a uniform random value $x \in [0, 1]$.

To make the sampling deterministic and parallelizable on the GPU, a sampling pattern $P \in [0, 1]^{H \times W}$ – uniformly distributed in $[0, 1]$ – is first generated by using a permutation of the numbers $\frac{1}{HW}\{0, ..., HW - 1\}$. We analyze four different strategies for generating the permutations: Random sampling, regular sampling, Halton sampling [19], and plastic sampling [50]. Plastic sampling has been selected, since it produced slightly superior results in all of our experiments. section B provides a detailed evaluation of the different strategies.

Ray-casting is then used to compute what is seen through the pixels at the determined sample locations. This information is stored in the high resolution image $S \in \mathbb{R}^{C \times H \times W}$. Since during training the same view is rendered many times using different sampling patterns, pre-computed high-resolution target images $T \in \mathbb{R}^{C \times H \times W}$ are provided with the low-resolution inputs. Then, the sampling process simply becomes a selection of pixels from $T$:

$$S_{ij} = \mathbb{1}_{I'_{ij} - P_{ij}} T_{ij}, \tag{4}$$

where $\mathbb{1}_x$ is 1 if $x > 0$ else 0. Since the sampling function in Equation 4 is a step function with zero gradients almost everywhere, it is not differentiable w.r.t. the importance map $I$, from which $I'$ is derived. Correspondingly, gradients in the loss function w.r.t. the weights and biases of the importance network will also be zero. Therefore, Equation 4 is approximated with a smooth sigmoid function to make it differentiable, so that gradients of the loss function can be back-propagated through all network stages to change the importance map accordingly. Then, the sampling function becomes

$$S_{ij} = \text{sig}\left(\alpha(I'_{ij} - P_{ij})\right)T_{ij}, \quad \text{sig}(x) := \frac{1}{1 + e^{-x}}, \tag{5}$$

where $\alpha > 0$ determines the steepness of the function. The differentiable approximation is used only in the training phase, while in the validation phase the ray-caster renders the discrete samples obtained via rejection sampling. For $\alpha \to \infty$, Equation 5 converges to Equation 4. A large value of $\alpha$ leads to samples that are either very close to 0 or 1, but leads to exploding gradients in the backward pass. A low value leads to samples that smoothly cover the entire interval between 0 and 1. In this case, however, the mismatch between the "fractional" samples that are used only during training and the discrete "binary" samples that are used for testing and validation leads to a significant reduction of the reconstruction quality. In our experiments, a value around $\alpha = 50$ always lead to the best results. Going beyond 50 quickly introduces floating-point precision issues and exploding gradients thereof. An evaluation of the dependency between the value of $\alpha$ and the reconstruction quality is provided in subsection 5.2.

### 3.3 Differentiable Reconstruction

Given the sparse set of samples $S$, the reconstruction function $\mathcal{N}_R$ needs to estimate the undefined pixel values to produce the dense high-resolution output image $O \in \mathbb{R}^{C \times H \times W}$. By using a differentiable reconstruction function, gradients of the loss function on the reconstructed images and the ground truth image can be back-propagated through the sampling stage to the importance map.

In principle, there are different possibilities to fulfill the requirement of differentiability: Firstly, a neural inpainting network can be trained on sparse inputs and the ground truth outputs to learn the reconstruction. However, as we have verified in a number of experiments, network-based inpainting [25, 39, 62] at a sparsity level as used in our application leads to low reconstruction quality (see Figure 6b). The highly varying sample density with gaps between valid pixel values of up to 20 pixels poses a challenging problem for known network architectures. Furthermore, since during training the sampling mask in our proposed pipeline is not binary but contains continuous values, techniques like Partial Convolutions [39] are not applicable.

Secondly, classical non-network-based inpainting methods can be employed, for instance, PDE-based methods solving a constrained Laplace problem [4, 14], or patch-based methods using non-local cost functions involving correspondence functions [24, 12, 8]. These methods, however, are not easily differentiable w.r.t. the sampling mask. For example, PDE-based methods use the samples as Dirichlet boundaries and, to the best of our knowledge, there is no meaningful interpretation of a "fractional" Dirichlet boundary. Patch-based methods, on the other hand, use a discrete search over the image space to find a correspondence function, which makes the derivation of continuous gradients impossible.

Therefore, we introduce a novel reconstruction approach that combines a differentiable inpainting method with a residual neural network that learns to improve the inpainting result. In particular, we propose a variation of the pull-push algorithm [15, 32], which is differentiable with respect to the sampling mask and can cope with a mask that comprises fractional values.

The pull-push algorithm builds upon the idea of mipmap hierarchies. Firstly, the sparsely sampled high-resolution image and the mask are recursively filtered and downsampled by a factor of 2. The pixel values are averaged using the fractional values in the sampling mask as weights (average pooling), and max-pooling is used to combine the values in the mask. This has the effect of filling
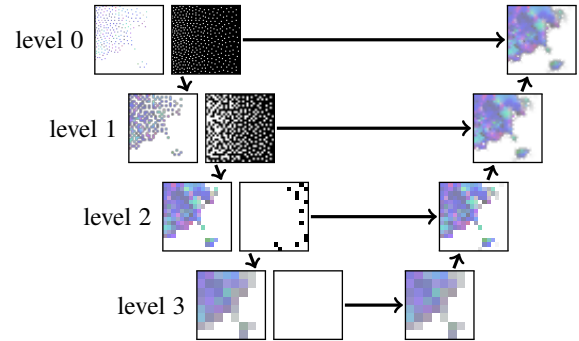


Fig. 3: Pull-push-based inpainting using a mipmap hierarchy of image samples and masks. The image is downsampled until all pixels are filled, and then upsampled by combining interpolated values from lower levels with the pixels at the current level. Masks are propagated through the hierarchy to obtain proper interpolation weights.

the undefined pixels with values that are averaged from a gradually increasing surrounding. Upon reaching a termination criterion, either a maximal number of steps or complete restoration of the undefined pixels, the images are bilinearly upscaled again. During upscaling, the pixel values from the coarse levels are weighted by the values in the mask at this level, and they are then blended with the value at the fine level based on the sampling values at that level. This allows to smoothly transition from filled pixels at the fine level that are kept in the output towards interpolated values for lower values in the sampling mask. A schematic illustration of the process is shown in Figure 3. Since the algorithm makes use exclusively of continuous pooling and interpolation operations, it is fully differentiable with respect to changes in the pixel data and the sampling mask. The forward code and a manually derived backward code are given in section D. The algorithm has been implemented via custom CUDA operations in PyTorch [45].

After inpainting the sparse samples via the pull-push algorithm, a fully convolutional network is used to improve the reconstruction by modeling the relationship between the inpainting result and the ground truth. The network sharpens the results and resolves blurred silhouettes created by the inpainting algorithm. We use the EnhanceNet [53] as base architecture for this learning task, which is discussed in detail in subsection 4.2. In particular, we use the EnhanceNet as a residual network that starts with the inpainting result and learns to infer the changes to the reconstructed samples. A quantitative comparison of different learning approaches is provided in subsection 5.2.

## 4 TRAINING METHODOLOGY

In this chapter, we provide a detailed discussion of the used network architectures, as well as the training and inference steps. We also shed light on the dependency of the reconstruction quality on the used loss functions.

### 4.1 Training Data

As training and validation input, 5000 images of randomly selected isosurfaces in the Ejecta data set, a particle-based supernova simulation, were generated via GPU ray-casting at a screen resolution of $512^2$. Each time step was resampled to Cartesian grids with a resolution of $256^3$ and $512^3$. The surfaces were rendered

(a) Importance Network $\mathcal{N}_I$      (b) Reconstruction Network $\mathcal{N}_R$ (EnhanceNet)
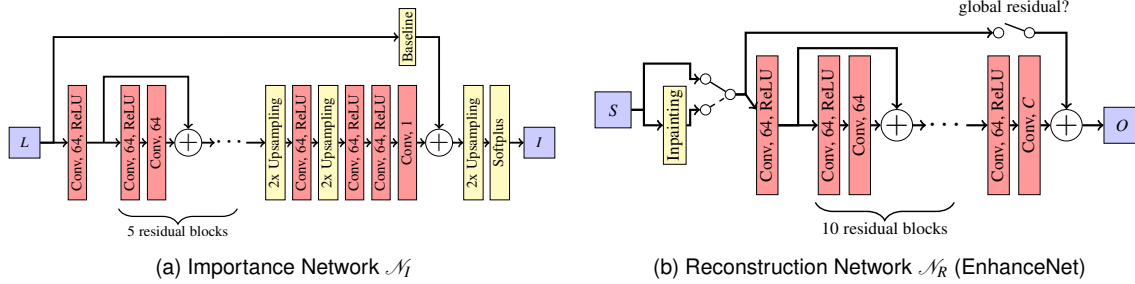
Fig. 4: Network architectures used in the proposed pipeline: To estimate the importance map, we use a smaller version of the EnhanceNet [53] with a 4x-upsampling factor, a residual connection with screen space gradient magnitude as baseline, and a 2x-upsampling network as a post-process. For the reconstruction network, we experimented with the option of passing the raw samples or interpolated samples as input and using a global residual connection or not. As network architecture, an EnhanceNet with 10 residual blocks is used.

from random camera positions, at varying distance to the object and always facing the object center. Renderings are taken from different time steps and resolution levels to let the pipeline learn features at different granularity [60]. The renderer provides the normals at the surface points, which are used in a post-process to compute colors via the Phong illumination model. From this image set, about 20.000 random crops of size $256^2$ and showing the isosurface in at least 50% of the pixels were taken, and split between training (80%) and validation (20%). For training, the mean importance value was set to $\mu = 0.1$, i.e., 10% of the samples (see Equation 3). This does not prohibit using less samples for validation and testing, yet we found it beneficial to allow the network to use more samples during training. We used the Adam [30] optimizer with a learning rate of $10^{-4}$. The networks were trained on a single GeForce GTX 1080 for 300 to 500 epochs in around 5-6 days.

### 4.2 Network Architectures

The proposed sampling pipeline comprises two trainable blocks: The importance network $\mathcal{N}_I$ and the reconstruction network $\mathcal{N}_R$. Both networks use 3x3 convolutions with zero-padding and a stride of one. The importance network is a variant of EnhanceNet [53], yet with only 5 residual blocks ( 4a). Instead of directly estimating the importance map, the network takes as input an importance map that is computed using screen space gradient magnitudes (subsection 3.1), and learns to improve this map using a residual connection. We refer to subsection 5.2 for a quantitative comparison of the network results w/ and w/o an initial gradient-based importance estimate.

The importance network performs 4x-upscaling of a low-resolution input image with $1/4$ the resolution of the final image. Thus, generating the input image requires to sample $1/4^2 = 6.25\%$ of the pixels in the target image, which already exceeds a prescribed limit of, e.g., 5% of the pixels. Therefore, an image with $1/8$ the final resolution is used as input, and the network performs 4x-upscaling to an intermediate image with $1/2$ the final resolution, followed by an additional 2x-upscaling of the inferred importance map. This allows to more aggressively reduce the number of initially required samples, i.e., only $1/8^2 \approx 1.56\%$ of the pixels in the final importance map need to be rendered.

The reconstruction network $\mathcal{N}_R$ estimates the mask, normal, and depth values at all pixels, thereby also changing the initial values that were drawn in the rendering process. A modified EnhanceNet

( 4b) shows superior reconstruction results compared to alternative architectures such as the U-Net [51]. Let us refer to section A for a more detailed analysis of both architectures. Both networks are provided in the code repository accompanying this paper.

Our experiments (subsection 5.2) show improved reconstruction quality if inpaining is performed first and the result is then passed to a network that uses a residual connection to learn the differences between this result and the ground truth. In addition to the inpainted input samples, we pass the sample mask to the network as a per-sample measure of certainty. Since the network produces output values in $\mathbb{R}$, both the mask and depth values are clamped to $[0,1]$ and the normals are scaled to unit length before shading is applied.

### 4.3 Loss Functions

We employ regular vector norms between the network prediction $O$ and the target image $T$ as primary loss functions on the individual output channels. Since the $L_2$ norm tends to smooth out the resulting images, we make use of the $L_1$ norm in this work. With the channels of the output image, i.e., the mask $M$, the normal map $N$, and the depth $D$, given as subscript, the $L_1$ loss of a selected channel $X$ is

$$\mathcal{L}_{1,X} = ||T_X - O_X||_1. \tag{6}$$

We do not employ additional perceptual losses, which were shown less effective for isosurface upsampling tasks [60].

The mask channel has a special meaning as it indicates whether or not a ray hits the isosurface. It is used in the final output to perform a hard selection between the reconstructed color values and the background. To make the mask differentiable, however, its values must be continuous, leading to a smooth blend rather than a binary decision. While this is acceptable along the silhouettes, in the interior it would noticeably distort the reconstruction. In principle, via a sigmoidal mapping it can be enforced that the mask values spread continuously between 0 and 1, yet we observed undesirable blurring when using this approach. To produce sharp masks that are either close to zero or one, we therefore constrain the reconstruction via two losses that are added to the regular $L_1$ loss on the mask. The first loss is a binary cross entropy (BCE) loss that "pulls" the values closer to either zero or one than a normal $L_1$ loss:

$$\mathcal{L}_{\text{bce}} = -\frac{1}{WH} \sum_{ij} \left( T_{M,ij} \log(O_{M,ij}) + (1 - T_{M,ij}) \log(1 - O_{M,ij}) \right).$$
$$\tag{7}$$

The BCE loss, however, requires that the output mask lies within $[0,1]$ and thus the mask is clamped beforehand. This leads to zero gradients once the mask reaches values outside of $[0,1]$. Therefore, we add the loss

$$\mathscr{L}_{\text{bounds}} = \frac{1}{WH}\sum_{ij}\left(\max(0,(2O_{M,ij}-1)^2-1)\right), \qquad (8)$$

which pushes values outside $[0,1]$ back into $[0,1]$ and leaves values within $[0,1]$ unchanged.

An additional loss term is required to account for the normalization step in Equation 3. The output of the importance map is normalized to limit the number of available samples. Hence, scaling the network output does not influence the values after normalization. Therefore, during training it can happen that the output values increase or decrease in an unbounded manner. To prevent this, a prior on the importance map is used to enforce that the mean is equal to one before the normalization step:

$$\mathscr{L}_{I,\text{prior}} = \left(1-\frac{1}{WH}\sum_{ij}I_{ij}\right)^2. \qquad (9)$$

The final loss function is a weighted sum of the individual loss terms over all channels, i.e., with $X \in \{M,N,D\}$ it becomes

$$\mathscr{L} = \sum_X \lambda_X \mathscr{L}_{1,X} + \lambda_{\text{bce}}\mathscr{L}_{\text{bce}} + \lambda_{\text{bounds}}\mathscr{L}_{\text{bounds}} + \rho\mathscr{L}_{I,\text{prior}}. \quad (10)$$

Loss weights around $\lambda_M = 5, \lambda_{\text{bce}} = 5, \lambda_{\text{bounds}} = 0.01, \lambda_N = 50, \lambda_D = 5$, and $\rho = 0.1$ lead to equally good reconstruction quality, while deviations from these values quickly worsen the reconstruction quality significantly.

## 5 RESULTS AND EVALUATION

In the following, we evaluate the proposed network pipeline. First, we introduce the quality metrics that are used to compare the results. We then analyze how our design decisions influence the reconstruction quality on the validation data (subsection 5.2). These statistics help to identify the network configurations with the best predictive skills. Next, the proposed network pipeline is compared to a fixed super-resolution network ( subsection 5.3). Finally, we shed light on the generalizability of the network pipeline to new views of Ejecta and data sets that were never seen during training (subsection 5.4).

### 5.1 Quality Metrics

The quality of network-based reconstruction is assessed using three different image quality metrics commonly used in image processing. These metrics compare the output $O$ of the network pipeline with a ground truth rendering $T$ at the target resolution.

The peak signal-to-noise ratio (PSNR) is based on the $L_2$ loss and is defined as

$$\text{PSNR}(O,T) = -10\log_{10}(||O-T||_2^2), \qquad (11)$$

where $O$ and $T$ are the network output and target image, respectively.

The Structural Similarity Index (SSIM) [59] extends on the idea of per-pixel losses by measuring the perceived quality using the mean and variance of contiguous pixel blocks in the images. It is defined as

$$\text{SSIM}(O,T) = \frac{(2\mu_O\mu_T + c_1)(2\sigma_{O,T} + c_2)}{(\mu_O^2 + \mu_T^2 + c_1)(\sigma_O^2 + \sigma_T^2 + c_2)}, \qquad (12)$$
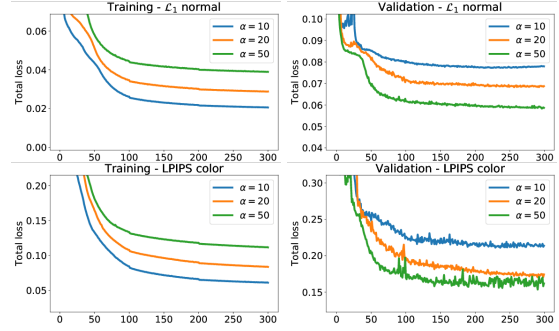


Fig. 5: Influence of the sharpness parameter $\alpha$ on the training process. A lower value leads to a lower cost during training, but increases the cost in the validation phase where a perfect step function is used.

where $\mu_O$ and $\mu_T$ are the average values of $O$ and $T$, $\sigma_O^2$ and $\sigma_T^2$ are the variances of $O$ and $T$, $\sigma_{O,T}$ is the covariance between $O$ and $T$, and $c_1$ and $c_2$ are small constants to avoid division by zero.

We also use the network-based Learned Perceptual Image Patch Similarity (LPIPS) metric [63] that predicts human perception of relative image similarities. LPIPS builds upon a network that is pretrained on an image classification task—using the AlexNet [33]—and computes a weighted average of the activations at hidden layers for a given output and target imag. Note that a lower LPIPS score is better, whereas PSNR and SSIM indicate higher quality by a higher score. Therefore, $1 - \text{LPIPS}$ is shown in our statistics for better comparison.

### 5.2 Validation of Design Decisions

Unless otherwise mentioned, all statistics presented in this section were computed on a validation data set using 2000 novel views of Ejecta at the resolution of $512^2$. The importance map was normalized to have a minimal value $l = 0.002$ and a mean value $\mu = 0.05$. "plastic" sampling was used in the sampling stage.

**Steepness of the Sampling Function** The parameter $\alpha$ in Equation 5 determines the steepness of the sampling function. A perfect step function, as used for testing, is obtained for $\alpha \to \infty$. Figure 5 compares the total loss on the training and validation data over the course of the optimization for different values of $\alpha$. A lower value of $\alpha$ leads to a lower cost on the training data, because smoother variations in the fractional samples can be used for reconstruction. However, this behaviour is reversed during validation, because the perfect step function corresponds to a lesser and lesser extent with an increasingly smooth sampling function. Higher values of $\alpha$, on the other hand, lead to better generalization, yet beyond 100 we observed instabilities in the training as well as numerical precision issues. Therefore, we decided to use $\alpha = 50$ in all of our experiments.

**Residual Connections for Reconstruction** In principle, there are different options to reconstruct a dense image from the sparse set of samples, including sole inpainting via the pull-push algorithm as well as inpainting in combination with network-based reconstruction w/ or w/o residual connections. In Figure 6, the reconstruction quality of all options is compared, using screen space gradient magnitudes as measure for generating the importance map.

As can be seen, the pull-push algorithm already provides a good initial guess on the reconstructed image, and reconstruction quality reduces significantly when it is not used. On the other hand,
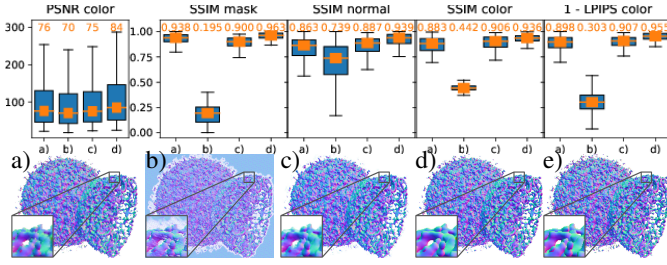
Fig. 6: Comparison of different reconstruction methods: (a) Only pull-push-based inpainting. (b) Only network-based reconstruction without residual. (c) Pull-push plus reconstruction network w/o residual. (d) Pull-push inpainting plus reconstruction network with residual. (e) Ground truth. An importance map from screen space gradient magnitudes is used in all examples, with $\mu = 5\%$ of samples.

the network-based approach fails to reliably fill the empty pixels, which is probably due to the vastly different distances between the sparse samples. When using the pull-push algorithm in combination with network-based reconstruction, but with disabled residual connections, no benefit over sole pull-push-based inpainting is gained. The best result is achieved with both pull-push-based inpainting and residual network connections. This is in line with the findings of Kim *et al.* [28], that the quality of network-based reconstruction improves if the network needs to learn only the changes to the baseline method.

**Residual Connections for Importance Mapping** On the validation data, we then analyze the reconstruction quality using different approaches for generating the importance map, i.e., constant importance, importance derived from screen space gradient magnitudes, as well as network-based importance with or without learning a residual to screen space gradient magnitudes. Figure 7 shows the results using the quality metrics described above.
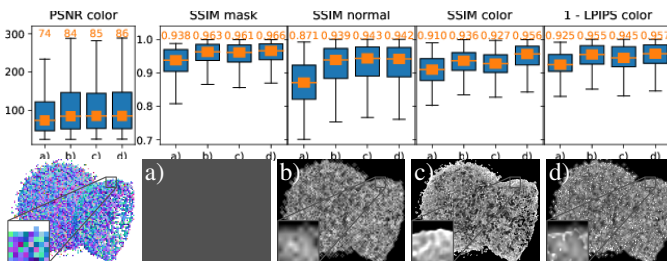


Fig. 7: Reconstruction quality using different importance maps. Bottom left: Low-resolution input. (a) Constant map, (b) based on gradient magnitudes, (c) only network-based learning, (d) network-based learning with residual on gradient magnitudes. $\mu = 5\%$ of samples were used. Top: Quality metrics for options (a) to (d).

As expected, screen space gradient magnitudes already hint to some important regions that should be sampled with higher density, significantly outperforming a constant importance map. For reconstructing the mask and normal channels, gradient magnitudes and network-based importance learning differ only marginally w.r.t. reconstruction quality. The importance network puts more emphasis on the object silhouettes and leads to an improved reconstruction of the normals over gradient magnitudes. On the other hand, it is important to note that the network learns the importance of features for an accurate screen space reconstruction
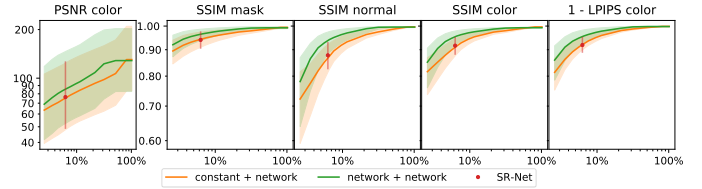


Fig. 8: Median reconstruction quality and 25% / 75% quantiles shown as confidence bands for increasing number of samples. Orange: Network-based pipeline using a constant importance map. Green: Network-based pipeline with the network-based importance map. The red dot represents the 4x-upsampling network from Weiss *et al.* [60].

without any prior information (Figure 7c). The best results are achieved by combining network-based importance learning and screen space gradient magnitudes via a residual network connection, demonstrating the feasibility of learning features that are important for an accurate reconstruction.

### 5.3 Convergence and Regular Sampling

We further analyse the convergence of the proposed sampling pipeline with increasing number of samples. The network is trained with 10% of the samples, but during inference the available number of samples is varied. The results in Figure 8 indicate, that with increasing number of samples the SSIM and LPIPS scores converge against their optima. Even though this seems logical at first, since the reconstruction network modifies the given samples, it could, in principal, converge against some other solution. Notably, already after taking 20% to 30% of samples the reconstruction is very close to the target.

We also compare the quality of adaptive sampling to fixed regular sampling using a 4x-upsampling network [60]. The 4x-upsampling network uses a regular sampling structure comprised of $1/4^2 = 6.25\%$ of the pixels in the high-resolution image, corresponding to a constant importance map with 6.25% of the samples when adaptive sampling is used. Figure 8 shows that the 4x-upsampling network (red) performs equally good as the adaptive pipeling using a constant importance map (orange). However, when the samples are placed adaptively according to the inferred importance map (green), the reconstruction quality is significantly increased at the same number of samples.

### 5.4 Generalizability

The importance and reconstruction networks are trained solely on Ejecta. To test how well the networks generalize, they are applied to a number of data sets that were never seen during training. We use a Richtmyer-Meshkov (RM) simulation at $1024 \times 1024 \times 960$, CT scans of a human skull (Skull) at $256^3$, an aneurism (Aneurism) at $256^3$, a bug (Bug) at $416^2 \times 247$, and a human body (Thorax) at $256^2 \times 942$, as well as a jet stream simulation (Jet) at $256^3$. Quantitative statistics for RM, Skull and novel views of Ejecta are given in Figure 9. Reconstructed images as well as SSIM and LPIPS statistics for all data sets are shown in Figure 10.

The pipeline generalizes well to new data sets and views, and it performs better than the baseline method using gradient magnitude-based importance mapping and pull-push-based inpainting. In particular, the network pipeline produces a tighter spread of the quantitative measures in general, indicating less significant outliers
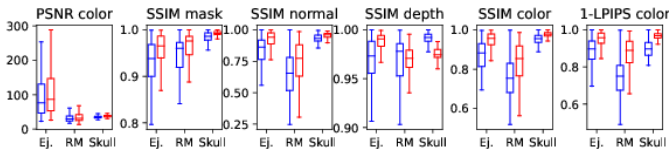
Fig. 9: Quality statistics for novel views of Ejecta and new data sets RM and Skull (see Figure 10). Baseline method (blue) refers to gradient magnitude-based importance mapping and pull-push-based inpainting. Results of the proposed network pipeline are shown in red.

in the reconstructed values. The network shows lower scores only for the depth maps reconstructed from sparse samples of RM and Skull. We attribute this to different zoom levels in the renderings and the training images, yet these inaccuracies do not affect the quality of the reconstructed color images. For reconstruction, we also analyzed the quality of other inpainting algorithms such as PDE-based methods. Notably, these methods are not differentiable and, thus, cannot be used for end-to-end training in combination with the importance network, yet they can be used for sole sparse image reconstruction. A comparison to the pull-push algorithm, however, does not show any perceivable differences. The results further indicate that the network pipeline can reconstruct images at high fidelity from only 5% of the samples that are used to render the data sets at full pixel resolution. In particular sharp edges are well preserved, since the network has learned to increase the sample density along them.

## 6 APPLICATION TO DVR

The proposed network pipeline can be applied to images that are rendered via Direct Volume Rendering (DVR), i.e., volume ray-casting using an emission-absorption model along the rays of sight. In contrast to isosurface ray-casting, not only one single ray-surface intersection point is rendered, but the colors of many sample points along the rays are blended using $\alpha$-compositing to account for volumetric attenuation.

### 6.1 Training and Validation

The importance and reconstruction networks receive RGB$\alpha$ images as input, and the network pipeline outputs the reconstructed high-resolution RGB$\alpha$ images. Interestingly, we observed a noticeable increase in quality when the gradients at the sample points along the view rays are used by the importance and reconstruction network. The normalized gradients in $[-1, 1]^3$ along a single ray are treated as emission and blended according to the volume rendering integral, just as blending the RGB colors. The resulting gradient map is then used as an additional input channel. Since the average gradients indicate, to a certain extent, whether two rays step through vastly different or similar regions, the gradient map serves as an additional coherence indicator. When only a single isosurface is rendered, the resulting values converge against the values in the normal map.

For training and validation, random transfer functions (TFs) are generated and used to render Ejecta, with $L_1$ losses on color and alpha in combination with a LPIPS-based perceptual loss (section C). Since the low-resolution input to the importance network is also generated with a TF, the network can learn to select features specific to that TF, even though this was never seen during training. It is important to note that the reconstruction

quality strongly depends on the use of TFs that include a broad range of different colors in the training step. For instance, if the training data only contains desaturated colors, strongly saturated colors during testing cannot be reconstructed.

### 6.2 DVR Results

For novel views of Ejecta and the data sets introduced in sub-section 5.4, Figure 11 shows a qualitative analysis of the results of importance sampling and reconstruction using DVR as well as SSIM and LPIPS statistics. None of these data sets was used in the training and validation phases, and the results have been generated using TFs that were never seen during training. The results indicate that the network pipeline generalizes well to new volumes and TFs, yet the reconstruction quality is affected by the occurring color variations. Especially for Thorax and Aneurism, where the TFs introduce rather small scale color variations in some areas, in these areas the network places the samples rather uniformly and, thus, cannot accurately reconstruct the rendered structured. Overall, it can be seen that the reconstruction problem is significantly more challenging when using DVR samples instead of isosurface samples. When rendering isosurfaces, the shading in the interior of the rendered structures is rather smooth, enabling the network to focus on the silhouettes and internal edges. In DVR, on the other hand, the network needs to learn both the shape and the color texture stemming from the application of a TF.

## 7 PERFORMANCE ANALYSIS

Even though performance improvements are not our main objective, it is interesting to see whether network-based adaptive sampling and image reconstruction can be faster than full-resolution GPU ray-casting, due to the reduced number of samples that need to be taken. The following performance tests were carried out on a workstation running Windows 10 with an Intel Xeon E5-1630 @3.70GHz CPU, 32GB RAM, and an NVidia Titan RTX. All timings are averages over 100 frames with random camera positions, with the screen resolution set to $1024^2$. The ray-caster uses a constant step size of 0.25 voxels and tricubic interpolation.

For some of the data sets shown in Figure 10 and Figure 11, Table 1 lists the times that are required by the pipeline stages and full-resolution volume ray-casting. Only for the larger data sets and DVR can the network pipeline achieve a slightly better performance than the ray-caster. Especially the reconstruction network consumes a significant portion of the overall time, sometimes even more than it requires to render at full resolution. This is because the reconstruction network requires a large amount of data access and arithmetic operations on the GPU, independent of the volume resolution.

On the one hand, the performance of the reconstruction network scales linearly with the number of pixels, and hence quadratically with the screen resolution. Volume rendering, on the other hand, scales quadratically with the screen resolution but also linear in the volume resolution. The sampling stage, even though it also scales in the volume resolution, performs a significantly smaller number of sampling operations than the full-resolution ray-caster. Thus, its overall contribution is negligible, so that performance benefits can be expected with increasing image and volume size. This is demonstrated in Table 2, where versions of RM and Ejecta at $2048^3$ are rendered at different resolution levels and large images sizes. Note that in these experiments an Nvidia Titan RTX graphics card with 24GB of memory was used to keep all data in memory.

TABLE 1: Timings (in milliseconds) of network-based volume rendering (averaged over 100 different views at $1024^2$ target resolution) for data sets shown in Figure 10 and Figure 11. Timings are for rendering the low-resolution input image ($128^2$) and the sparse set of samples (5% and 10% of the target resolution for isosurface rendering and DVR respectively), generating the importance map and sampling pattern, reconstructing the image, and GPU ray-casting at the target resolution.

| | Test case | Rendering | Importance | Reconstruction | GT |
|---|---|---|---|---|---|
| ISO | RM $1024^3$ | 34.3 | 5.8 | 92.0 | 89.4 |
| | Ejecta $512^3$ | 24.3 | 7.4 | 92.1 | 105.8 |
| | Skull $256^3$ | 6.1 | 5.9 | 93.7 | 27.3 |
| DVR | RM $1024^3$ | 51.7 | 5.8 | 91.8 | 158.3 |
| | Ejecta $512^3$ | 46.5 | 5.8 | 92.2 | 224.8 |
| | Thorax $512^3$ | 15.9 | 5.9 | 92.9 | 63.7 |

TABLE 2: Performance scaling w.r.t. image and volume size. Each entry shows the total time of the network pipeline (low-resolution rendering, importance network, sparse sampling, reconstruction network) and the time required by the volume ray-caster at full resolution. All timings are in milliseconds. The cells are colored with a diverging color map, encoding the performance differences from red (superior performance of ray-casting) to blue (superior performance of the network pipeline).

| | | | Screen resolution | | | |
|---|---|---|---|---|---|---|
| | | | 256 | 512 | 1024 | 2048 |
| ISO | Ejecta | 256 | 24/11 | 45/20 | 115/75 | 398/294 |
| | | 512 | 31/16 | 55/31 | 124/106 | 405/400 |
| | | 1024 | 45/47 | 73/104 | 141/221 | 445/655 |
| | | 2048 | 85/278 | 132/806 | 238/1358 | 724/2324 |
| | RM | 256 | 24/4 | 41/11 | 111/42 | 393/173 |
| | | 512 | 31/6 | 51/15 | 116/58 | 397/236 |
| | | 1024 | 45/15 | 71/31 | 132/89 | 411/366 |
| | | 2048 | 87/122 | 124/249 | 211/453 | 646/802 |
| DVR | Ejecta | 256 | 46/17 | 64/46 | 131/162 | 413/574 |
| | | 512 | 59/36 | 75/73 | 144/225 | 447/785 |
| | | 1024 | 86/151 | 103/262 | 180/433 | 594/1274 |
| | | 2048 | 160/575 | 186/1692 | 333/2771 | 1319/4292 |
| | RM | 256 | 32/8 | 49/21 | 117/71 | 398/280 |
| | | 512 | 41/11 | 58/27 | 125/96 | 404/363 |
| | | 1024 | 60/29 | 81/58 | 149/158 | 450/549 |
| | | 2048 | 115/203 | 143/410 | 248/656 | 830/1304 |

It can be seen that for large image sizes—where the GPU is fully utilized by the network—and volume sizes larger than $1024^3$, the network pipeline outperforms the GPU ray-caster. Even though a ray-caster using advanced acceleration schemes can achieve improved performance, we are confident that in these scenarios faster deep-learning hardware and performance-optimized network architectures will let the performance differences grow due to better scalability of the network pipeline.

## 8 CONCLUSION AND FUTURE WORK

In this paper we have introduced and analyzed a network pipeline that learns adaptive screen space sampling and reconstruction for 3D visualization, with the focus on volume rendering applications. For the first time, to our best knowledge, a fully differentiable adaptive sampling pipeline comprised of an importance network, a sampling stage, and a reconstruction network is proposed. Our experiments have shown, that the pipeline learns to determine the locations that are important for an accurate image reconstruction, and achieves high reconstruction quality for a sparse set of samples.

We are particular intrigued about the quality of the results compared to sampling methods that consider explicitly certain feature descriptors. Even without such supervision, the network pipeline can improve on the reconstruction quality, using solely image-based quality losses. We believe that especially for data visualization there is value in the observation that artificial neural networks can learn the relevance of structures for generating visual representations. For sole rendering tasks, on the other hand, superior performance compared to classical volume ray-casting can only be achieved for large image and volume sizes.

The application to DVR opens the interesting question whether the proposed network pipeline can be used beyond adaptive sampling in screen space, and learn where to sample in object space so that the relevant information is conveyed visually. Conceptually this requires end-to-end learning of a mapping from a low-resolution object space representation to a high-resolution visual representation. The ultimate goal is to let the network learn to convert a low-resolution input volume to a compact yet feature-preserving latent-space representation from which a highly accurate view can be inferred.

In particular, we envision a neural volume rendering pipeline, where during training a neural scene representation is build and trained end-to-end with a renderer that learns sampling and color mapping simultaneously. In the future, we will analyse whether a network can learn a suitable color mapping for a given volumetric field. We also see challenging research problems in the area of transfer learning, to infer the most important samples for training, and to generate synthetic volumetric fields to enable training in domains where training data is rare.

## REFERENCES

[1] T. Bashford-Rogers, K. Debattista, and A. Chalmers. Importance driven environment map sampling. *IEEE transactions on visualization and computer graphics*, 20, 11 2013. doi: 10.1109/TVCG.2013.258

[2] S. Belyaev, P. Smirnov, V. Shubnikov, and N. Smirnova. Adaptive algorithm for accelerating direct isosurface rendering on gpu. *Journal of Electronic Science and Technology*, 16:222–231, 01 2018. doi: 10.11989/JEST.1674-862X.71013102

[3] M. Berger, J. Li, and J. A. Levine. A generative model for volume rendering. *IEEE Transactions on Visualization and Computer Graphics*, 25(4):1636–1650, April 2019. doi: 10.1109/TVCG.2018.2816059

[4] M. Bertalmio, A. L. Bertozzi, and G. Sapiro. Navier-stokes, fluid dynamics, and image and video inpainting. In *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. CVPR 2001*, vol. 1, pp. I–I, Dec 2001. doi: 10.1109/CVPR.2001.990497

[5] M. R. Bolin and G. W. Meyer. A perceptually based adaptive sampling algorithm. In *Proceedings of the 25th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH 98, p. 299309. Association for Computing Machinery, New York, NY, USA, 1998. doi: 10.1145/280814.280924

[6] L. Campagnolo, W. Celes, and L. Figueiredo. Accurate volume rendering based on adaptive numerical integration. pp. 17–24, 08 2015. doi: 10.1109/SIBGRAPI.2015.27

[7] M. Chu, Y. Xie, L. Leal-Taixé, and N. Thuerey. Temporally coherent gans for video super-resolution (tecogan). *arXiv:1811.09393*, 2018.

[8] A. Criminisi, P. Perez, and K. Toyama. Region filling and object removal by exemplar-based image inpainting. *IEEE*
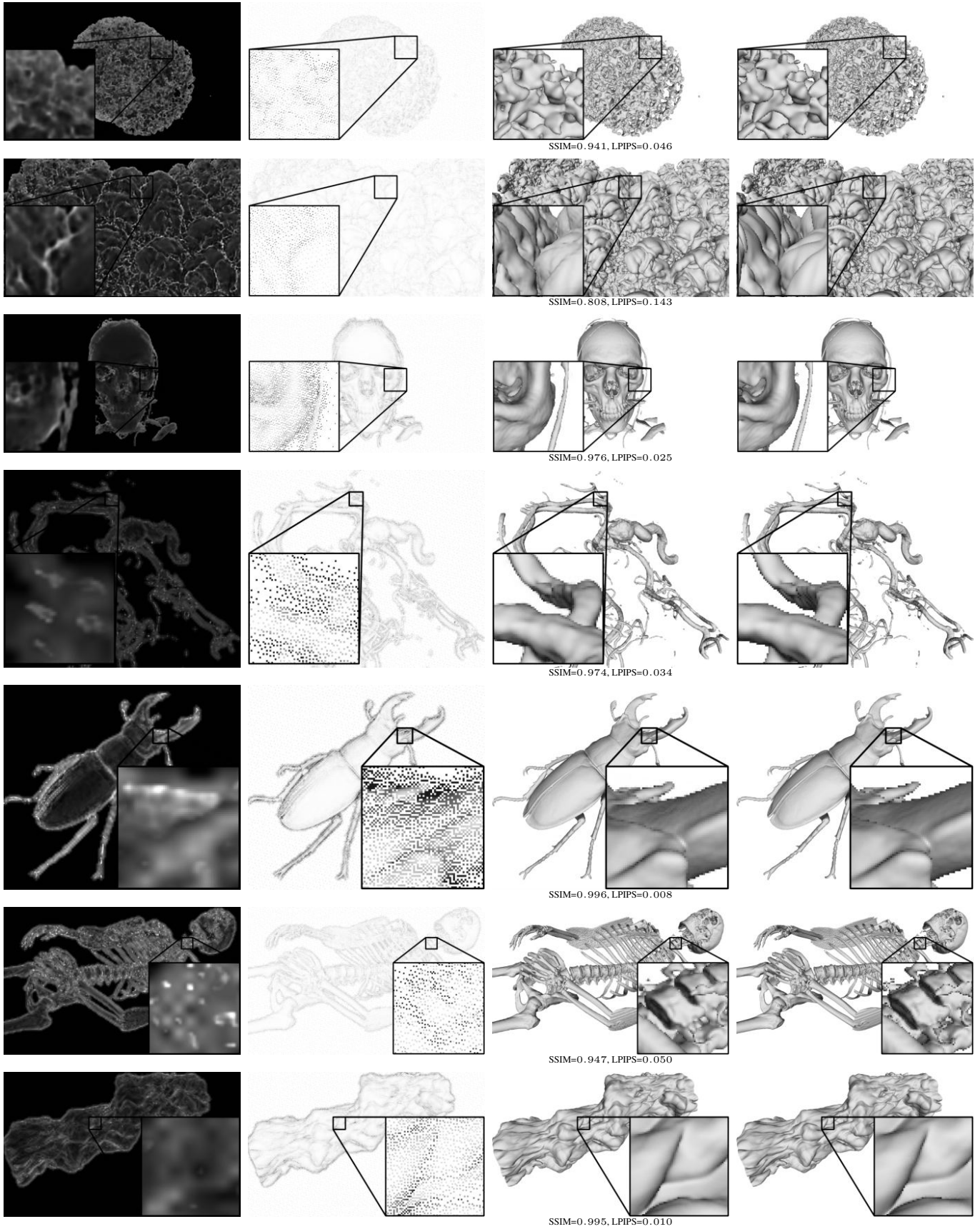
Fig. 10: Visual comparison of adaptive sampling of isosurfaces. From left to right: The importance map and, the sparse set of samples, the inpainted samples, the network output, the ground truth (normals and colors using reconstructed normals). From top to bottom: Novel views of Ejecta, RM, Skull, Aneurism, Bug, Human, Jet. Networks were trained only on Ejecta. Number of samples is $\mu = 5\%$ of the pixels in the output image.
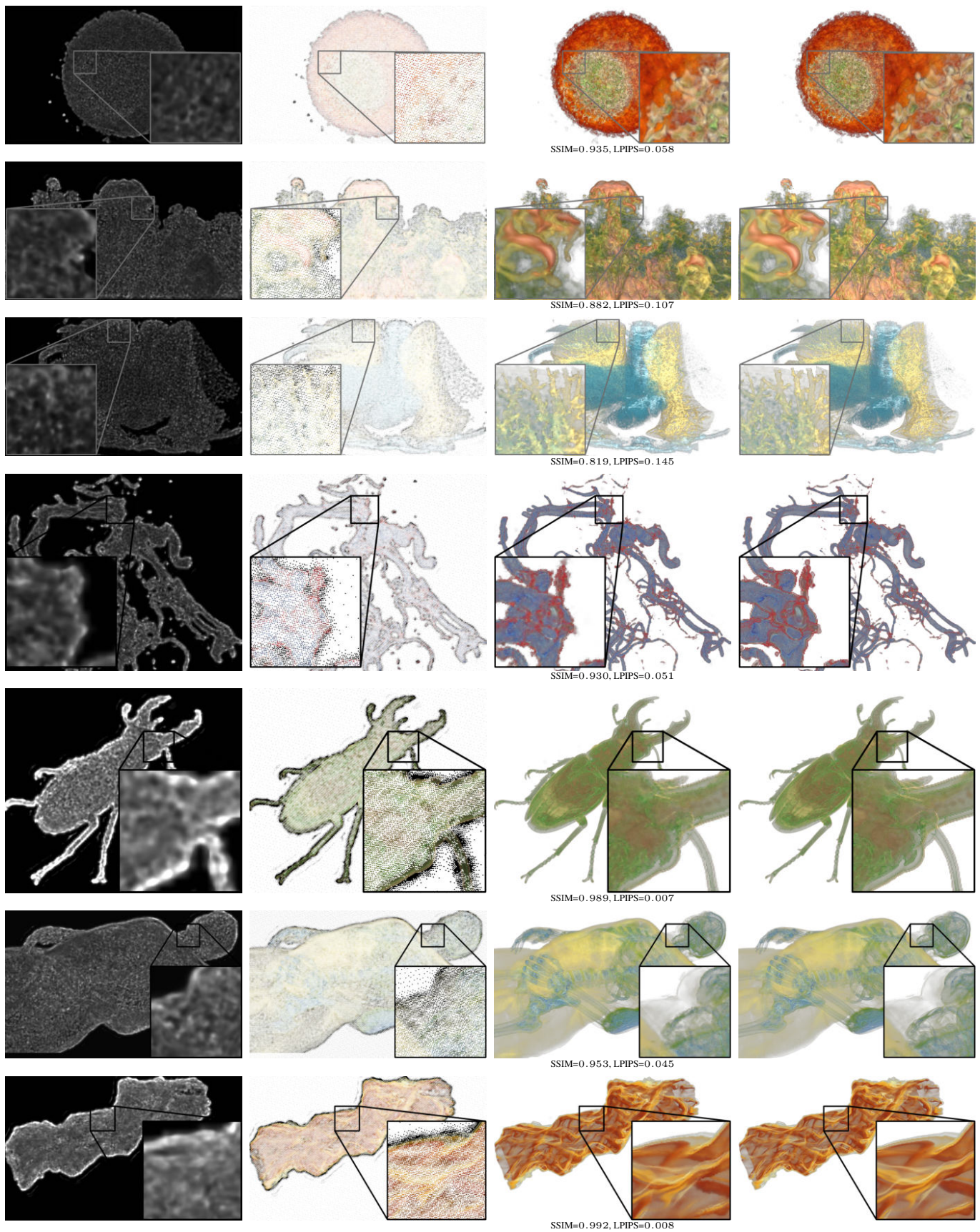
Fig. 11: Visual comparison of adaptive sampling for DVR. Each row shows – from left to right – the importance map, the sparse set of samples, the network output, the ground truth. From top to bottom: Novel views of Ejecta, RM, Thorax, Aneurism, Bug, Human, Jet. Networks were trained only on Ejecta. Number of samples is $\mu = 10\%$ of the pixels in the output image.

*Transactions on Image Processing*, 13(9):1200–1212, Sep. 2004. doi: 10.1109/TIP.2004.833105

[9] J. Danskin and P. Hanrahan. Fast algorithms for volume ray tracing. In *Proceedings of the 1992 Workshop on Volume Visualization*, VVS 92, p. 9198. Association for Computing Machinery, New York, NY, USA, 1992. doi: 10.1145/147130.147155

[10] O. Deussen, M. Spicker, and Q. Zheng. Weighted linde-buzo-gray stippling. *ACM Trans. Graph.*, 36(6), Nov. 2017. doi: 10.1145/3130800.3130819

[11] C. Dong, C. C. Loy, K. He, and X. Tang. Learning a deep convolutional network for image super-resolution. In *European conference on computer vision*, pp. 184–199. Springer, 2014.

[12] A. A. Efros and T. K. Leung. Texture synthesis by non-parametric sampling. In *Proceedings of the Seventh IEEE International Conference on Computer Vision*, vol. 2, pp. 1033–1038 vol.2, Sep. 1999. doi: 10.1109/ICCV.1999.790383

[13] S. Frey, F. Sadlo, K. Ma, and T. Ertl. Interactive progressive visualization with space-time error control. *IEEE Transactions on Visualization and Computer Graphics*, 20(12):2397–2406, 2014.

[14] P. Getreuer. Total variation inpainting using split bregman. *Image Processing On Line*, 2:147–157, 2012. doi: 10.5201/ipol.2012.g-tvi

[15] S. J. Gortler, R. Grzeszczuk, R. Szeliski, and M. F. Cohen. The lumigraph. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pp. 43–54, 1996.

[16] J. Grtler, M. Spicker, C. Schulz, D. Weiskopf, and O. Deussen. Stippling of 2d scalar fields. *IEEE Transactions on Visualization and Computer Graphics*, 25(6):2193–2204, June 2019. doi: 10.1109/TVCG.2019.2903945

[17] B. Guenter, M. Finch, S. Drucker, D. Tan, and J. Snyder. Foveated 3d graphics. *ACM Transactions on Graphics (TOG)*, 31(6):1–10, 2012.

[18] L. Guo, S. Ye, J. Han, H. Zheng, H. Gao, D. Chen, J.-X. Wang, and C. Wang. Spatial super-resolution for vector field data analysis and visualization. In *Proceedings of IEEE Pacific Visualization Symposium*, 2020.

[19] J. H. Halton. Algorithm 247: Radical-inverse quasi-random point sequence. *Commun. ACM*, 7(12):701702, Dec. 1964. doi: 10.1145/355588.365104

[20] J. Han, J. Tao, H. Zheng, H. Guo, D. Z. Chen, and C. Wang. Flow field reduction via reconstructing vector data from 3-d streamlines using deep learning. *IEEE computer graphics and applications*, 39(4):54–67, 2019.

[21] J. Han and C. Wang. Tsr-tvd: Temporal super-resolution for time-varying data analysis and visualization. *IEEE Transactions on Visualization and Computer Graphics*, 26(1):205–215, 2019.

[22] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.

[23] W. He, J. Wang, H. Guo, K. Wang, H. Shen, M. Raj, Y. S. G. Nashed, and T. Peterka. Insitunet: Deep image synthesis for parameter space exploration of ensemble simulations. *IEEE Transactions on Visualization and Computer Graphics*, 26(1):23–33, Jan 2020. doi: 10.1109/TVCG.2019.2934312

[24] H. Igehy and L. Pereira. Image replacement through texture synthesis. In *Proceedings of International Conference on Image Processing*, vol. 3, pp. 186–189 vol.3, Oct 1997. doi: 10.1109/ICIP.1997.632049

[25] S. Iizuka, E. Simo-Serra, and H. Ishikawa. Globally and locally consistent image completion. *ACM Transactions on Graphics (ToG)*, 36(4):1–14, 2017.

[26] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pp. 694–711. Springer, 2016.

[27] J. Kim, J. Kwon Lee, and K. Mu Lee. Accurate image super-resolution using very deep convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1646–1654, 2016.

[28] J. Kim, J. Kwon Lee, and K. Mu Lee. Accurate image super-resolution using very deep convolutional networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.

[29] J. Kim, J. Kwon Lee, and K. Mu Lee. Deeply-recursive convolutional network for image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1637–1645, 2016.

[30] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv:1412.6980*, 2014.

[31] A. Kratz, J. Reininghaus, M. Hadwiger, and I. Hotz. Adaptive screen-space sampling for volume ray-casting. *ZIB-Report*, 2011.

[32] M. Kraus. The pull-push algorithm revisited. *Proceedings GRAPP*, 2:3, 2009.

[33] A. Krizhevsky. One weird trick for parallelizing convolutional neural networks. *arXiv:1404.5997*, 2014.

[34] A. Kuznetsov, N. K. Kalantari, and R. Ramamoorthi. Deep adaptive sampling for low sample count rendering. In *Computer Graphics Forum*, vol. 37, pp. 35–44. Wiley Online Library, 2018.

[35] J. Lawrence, S. Rusinkiewicz, and R. Ramamoorthi. Adaptive numerical cumulative distribution functions for efficient importance sampling. In *Rendering Techniques*, pp. 11–20, 2005.

[36] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4681–4690, 2017.

[37] M. Levoy. Volume rendering by adaptive refinement. *The Visual Computer*, 6(1):2–7, 1990.

[38] S. Lindholm, D. Jönsson, H. Knutsson, and A. Ynnerman. Towards data centric sampling for volume rendering. In *SIGRAD*, 2013.

[39] G. Liu, F. A. Reda, K. J. Shih, T.-C. Wang, A. Tao, and B. Catanzaro. Image inpainting for irregular holes using partial convolutions. In *The European Conference on Computer Vision (ECCV)*, September 2018.

[40] P. Longhurst, K. Debattista, and A. Chalmers. A gpu based saliency map for high-fidelity selective rendering. In *Proceedings of the 4th International Conference on Computer Graphics, Virtual Reality, Visualisation and Interaction in Africa*, AFRIGRAPH 06, p. 2129. Association for Computing Machinery, New York, NY, USA, 2006. doi: 10.1145/1108590.1108595

[41] M. Mara, M. McGuire, B. Bitterli, and W. Jarosz. An efficient denoising algorithm for global illumination. In *Proceedings of*

*High Performance Graphics*. ACM, New York, NY, USA, jul 2017. doi: 10.1145/3105762.3105774

[42] K. Myszkowski. The visible differences predictor: applications to global illumination problems. In G. Drettakis and N. Max, eds., *Rendering Techniques '98*, pp. 223–236. Springer Vienna, Vienna, 1998.

[43] K. Novins and J. Arvo. Controlled precision volume integration. In *Proceedings of the 1992 Workshop on Volume Visualization*, VVS 92, p. 8389. Association for Computing Machinery, New York, NY, USA, 1992. doi: 10.1145/147130.147154

[44] J. Painter and K. Sloan. Antialiased ray tracing by adaptive progressive refinement. In *Proceedings of the 16th annual conference on Computer graphics and interactive techniques*, pp. 281–288, 1989.

[45] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, eds., *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran Associates, Inc., 2019.

[46] M. Prantl, L. Vása, and I. Kolingerová. Fast screen space curvature estimation on gpu. In *VISIGRAPP (1: GRAPP)*, pp. 151–160, 2016.

[47] C. R. Alla Chaitanya, A. S. Kaplanyan, C. Schied, M. Salvi, A. Lefohn, D. Nowrouzezahrai, and T. Aila. Interactive reconstruction of monte carlo image sequences using a recurrent denoising autoencoder. *ACM Transactions on Graphics*, 36:1–12, 07 2017. doi: 10.1145/3072959.3073601

[48] M. Ramasubramanian, S. N. Pattanaik, and D. P. Greenberg. A perceptually based physical error metric for realistic image synthesis. In *Proceedings of the 26th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH 99, p. 7382. ACM Press/Addison-Wesley Publishing Co., USA, 1999. doi: 10.1145/311535.311543

[49] J. Rigau, M. Feixas, and M. Sbert. Refinement criteria based on f-divergences. In *Rendering Techniques*, pp. 260–269, 2003.

[50] M. Roberts. The unreasonable effectiveness of quasirandom sequences. http://extremelearning.com.au/unreasonable-effectiveness-of-quasirandom-sequences/, 2020. Accessed: 2020-02-14.

[51] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In N. Navab, J. Hornegger, W. M. Wells, and A. F. Frangi, eds., *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, pp. 234–241. Springer International Publishing, Cham, 2015.

[52] M. S. Sajjadi, R. Vemulapalli, and M. Brown. Frame-recurrent video super-resolution. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6626–6634, 2018.

[53] M. S. M. Sajjadi, B. Scholkopf, and M. Hirsch. Enhancenet: Single image super-resolution through automated texture synthesis. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.

[54] W. Shi, J. Caballero, F. Huszár, J. Totz, A. P. Aitken, R. Bishop, D. Rueckert, and Z. Wang. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1874–1883, 2016.

[55] Y. Tai, J. Yang, and X. Liu. Image super-resolution via deep recursive residual network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3147–3155, 2017.

[56] X. Tao, H. Gao, R. Liao, J. Wang, and J. Jia. Detail-revealing deep video super-resolution. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 4472–4480, 2017.

[57] G. Tkachev, S. Frey, and T. Ertl. Local prediction models for spatiotemporal volume visualization. *IEEE Transactions on Visualization and Computer Graphics*, pp. 1–1, 2019. doi: 10.1109/TVCG.2019.2961893

[58] O. T. Tursun, E. Arabadzhiyska-Koleva, M. Wernikowski, R. Mantiuk, H.-P. Seidel, K. Myszkowski, and P. Didyk. Luminance-contrast-aware foveated rendering. *ACM Trans. Graph.*, 38(4), July 2019. doi: 10.1145/3306346.3322985

[59] Z. Wang, A. C. Bovik, H. R. Sheikh, E. P. Simoncelli, et al. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.

[60] S. Weiss, M. Chu, N. Thuerey, and R. Westermann. Volumetric isosurface rendering with deep learning-based super-resolution. *IEEE Transactions on Visualization and Computer Graphics*, pp. 1–1, 2019. doi: 10.1109/TVCG.2019.2956697

[61] Q. Xu, S. Bao, R. Zhang, R. Hu, and M. Sbert. Adaptive sampling for monte carlo global illumination using tsallis entropy. In *International Conference on Computational and Information Science*, pp. 989–994. Springer, 2005.

[62] J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, and T. S. Huang. Free-form image inpainting with gated convolution. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 4471–4480, 2019.

[63] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018.

[64] Z. Zhou, Y. Hou, Q. Wang, G. Chen, J. Lu, Y. Tao, and H. Lin. Volume upscaling with convolutional neural networks. In *Proceedings of the Computer Graphics International Conference*, pp. 1–6, 2017.

# APPENDIX A
## COMPARISON WITH THE U-NET FOR RECONSTRUCTION

For reconstruction, we also tested different variants (by varying the number of levels and channels at each level) of the U-Net architecture [7]. As one can see in Figure 13, in our application the EnhanceNet vastly outperforms all considered U-Net variants.
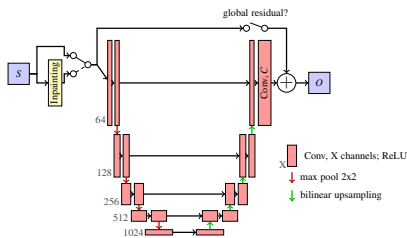


Fig. 12: Reconstruction network based on the U-Net architecture. See 4b for the EnhanceNet architecture we use in our experiments.
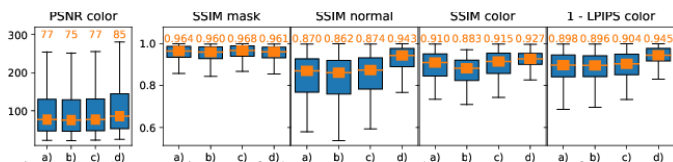


Fig. 13: Comparison of the U-Net and EnhanceNet for sparse image reconstruction: U-Net 4-4 (a), U-Net 5-3 (b), U-Net 5-4 (c), EnhanceNet (d). *a-b* indicates *a* levels and $2^{b+i}$ channels in level *i* (zero-based). The importance network is trained together with the reconstruction network.

# APPENDIX B
## COMPARISON OF DIFFERENT SAMPLING PATTERN

For deterministic and parallelizable sampling on the GPU, we use a pre-computed sampling pattern in combination with rejection sampling (subsection 3.2). The sampling pattern $P \in [0,1]^{H \times W}$ contains permutations of uniformly distributed numbers in $[0,1]$, $\frac{1}{HW}\{0, ..., HW-1\}$. Here we analyze the four different strategies employed for generating the permutations (Figure 14, top): Random sampling, regular sampling, Halton sampling [2], and plastic sampling [6].

Random sampling generates a random permutation of the numbers in *P*. Regular sampling arranges the pixels in a quad-tree and enumerates them using breath-first traversal to generate the sampling pattern. Random and regular sampling introduce, respectively, largely varying sample densities and a strong bias of the sample distribution towards the top of the image. Both Halton and plastic sampling are deterministic and produce quasi-random sequences with a fairly uniform distribution. As revealed by the quantitative analysis in Figure 15, even though all sampling strategies allow reconstructing the final image at high accuracy, slight differences are noticeable. Halton and plastic sampling lead to superior quality, in particular w.r.t. the variance of the quality metrics. Plastic sampling, designed as a low-discrepancy sampling sequence, shows the lowest variance and slightly higher scores than Halton sampling. We therefore use plastic sampling in our implementation.
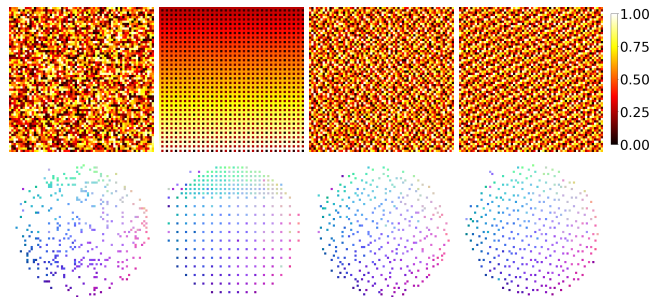


Fig. 14: Comparison of random sampling, regular sampling, Halton sampling and plastic sampling (left to right). Top: The sampling sequences. Bottom: The sequences applied to render a sphere with constant importance of $\mu = 0.1$ (shown are color coded normals at rendered fragments).



Fig. 15: Reconstruction quality using (a) uniform random, (b) regular, (c) Halton (c), and (c) plastic sampling with $\mu = 5\%$ of samples.

# APPENDIX C
## APPLICATION TO DVR

In this section, we provide additional details on how the proposed adaptive sampling pipeline is applied to DVR images, as mentioned in section 6. First we present the changes to the pipeline in terms of input and output channels and the used loss function. Second, we describe how to generate the training data including sampling of transfer functions.

*Input Channels and Loss Function:* First, the input channels to the network pipeline are reinterpreted. For isosurfaces, a mask, normals and depth were passed to the network as input (5 channels per pixel), now color images from the DVR images, together with alpha, depth and normal maps are used as input (8 channels per pixel). The network also only reconstructs color images in RGB$\alpha$ space.

For DVR, depth and normal maps are computed by treating the screen space depth and normal at each sample in object space like a regular color and blended as such with the opacity given by the transfer function (TF). The result is a single depth and normal value per ray which can be interpreted as a weighted average of the depth and normal of all samples along the ray. We found that adding depth and normals as input channels improves the quality of the reconstruction as it provides additional locally consistent information about the curvature of the object.

Second, the loss functions on the individual channels as used for isosurfaces are replaced by losses only on the RGB$\alpha$-color. We apply $L_1$ losses on the color and alpha and an additional LPIPS metric [10] as a perceptual loss, weighted equally:

$$\mathcal{L}_{\text{dvr}} = \mathcal{L}_{1,rgba} + \mathcal{L}_{\text{LPIPS},rgb}. \tag{13}$$

We found that adding a perceptual loss is critical in reconstructing fine details and sharp silhouettes. The networks operates in RGB space, other colorspaces like HSV, XYZ or CIELAB did not improve the result. Furthermore, the training data is augmented by

randomly shuffling the RGB channels. This helps the network to not overfit for a specific color.

Data Set Generation: For training and validation, random transfer functions (TFs) are generated (see below) and Ejecta was used as data set. The test images in the result section use user-generated TFs. Note that since the low-resolution input for the importance network is also generated with a TF, the network can learn to select features specific to that TF, even though it was never seen during training.

To generate meaningful TFs, first a density histogram is computed and then a Gaussian Mixture Model (GMM) is used to cluster densities in an unsupervised manner. GMMs have been previously used to cluster two-dimensional feature points [9], e.g., density and gradient magnitude. Our approach follows the same idea to cluster one-dimensional feature points, i.e., density values. The GMM represents each cluster as a 1D Gaussian function with a certain mean, i.e., the cluster center, and standard deviation, i.e., the cluster spread. To determine the number of components of the GMM, several GMMs with different numbers of components are build and the one with the lowest Bayesian information criterion (BIC) [8] value is selected. BIC penalizes the number of components and prevents overfitting using many components.

After computing the GMM, the number of peaks of the TF is sampled uniformly between 3 and 5. The represented density for each peak is sampled from the computed GMM. Next, a width in density space is sampled uniformly from $[0.005, 0.03]$ and the opacity at that peak from $[0.1, 1.0]$. As colormaps, predefined colormaps from SciVisColor[1] are randomly sampled. The generation process is visualized in Figure 16.
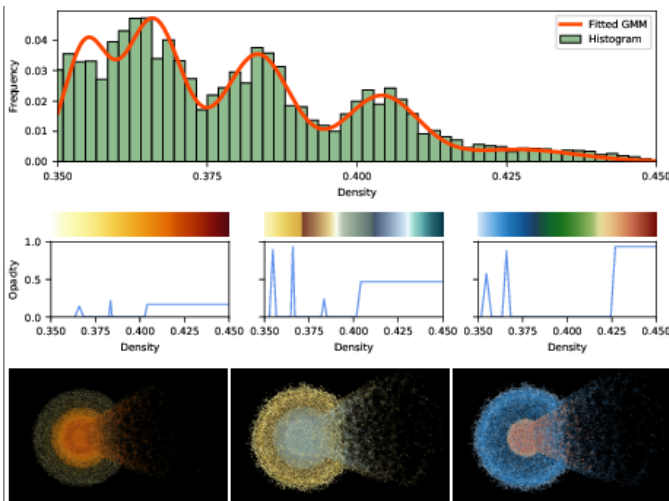


Fig. 16: First row: Histogram of the density values of the Ejecta data set and matched GMM. Second row: three sampled transfer functions with opacity and color. Third row: Renderings from the training data set with those transfer functions.

We note that it is important for the quality of the reconstruction that the color transfer functions in the training data include a broad range of colors. For example, if the training data only contains desaturated colors, strongly saturated colors during testing cannot be reconstructed.

1. https://sciviscolor.org/home/colormaps

## APPENDIX D
## PULL-PUSH ALGORITHM

As a baseline method to interpolate the sparse samples, we apply a variation of the push-pull algorithm [1, 3], see Alg. 1 for the pseudo code. The algorithm builds upon the idea of mip-map levels: first, the image is downscaled using bilinear interpolation with weights based on the mask. Then, the image is upscaled again and blended with the values at the finer levels with the mask values at the finer levels. We refer to subsection 3.3 for more details in the context of the adaptive sampling pipeline. The pull-push algorithm can be directly extended to fractional masks as shown in Alg. 1. During the upsampling stage, the mask is not treated binary, i.e. either take the original pixel at the fine level or use the interpolated value from the coarse level, but fractional with a linear interpolation between the original value and the interpolated value. Furthermore, the algorithm consists only of linear pooling and interpolation layers which are easy to differentiate with respect to the input mask. We refer to subsection E.3 for an outline on how to derive the backward pass.

## APPENDIX E
## DIFFERENTIATION OF THE SAMPLING AND RECON-STRUCTION STAGES

The adjoint code for the gradient propagation in the backward pass is automatically generated by PyTorch for the networks, the loss functions and for the sampling function Equation 5. For the pull-push algorithm (Alg. 1), the adjoint code was manually derived and implemented as a custom operation. In this section we provide the fundamentals of the adjoint method to manually derive the adjoint code and show how it can be applied to the sampling function and the pull-push algorithm.

### E.1 Fundamentals of the Adjoint Method

The adjoint method has a long history in Optimal Control Theory, we refer the interested reader to the book by Lions [4] for a complete mathematical introduction. Here, we briefly sketch the fundamentals following the notation by McNamara *et al.* [5].

By ignoring applications to linear systems and differential equations and focussing on chained functions instead, the adjoint method simplifies to an application of the chain rule. Let the algorithm be defined as a concatenation of functions $f_i$ with parameters $w_i$ starting from an input value $x_0$,

$$
\begin{aligned}
x_1 &= f_1(x_0, w_1) \\
x_2 &= f_2(x_1, w_2) \\
&\vdots \\
x_n &= f_n(x_{n-1}, w_n) \\
s &= J(x_n, w_J).
\end{aligned}
\tag{14}
$$

The result $s$ has to be a scalar value, this is crucial for the application of the adjoint method in this simple form. In the context of neural networks, $x_0$ would be the input image, $f_1$ to $f_n$ the network layers with weights $w_i$ and feature vectors $x_i$, $J$ would be the loss function with target image $w_J$ and $s$ the scalar score.

During training, we are interested in the derivatives $\frac{\partial J}{\partial w_i}$ to update the weights or in e.g. $\frac{\partial J}{\partial x_0}$ to update the initial image in a feature-visualization context. First, given the – possibly vector valued – variables $x_i$ and $w_i$, let the *adjoint variables* $\hat{x}_i$ and $\hat{w}_i$

**Algorithm 1** Pseudocode of the pull-push algorithm for power-of-two input images (a version handling non-power-of-two inputs and the adjoint code for computing the derivative with respect to the mask and data are provided in the source code).

---

1: **function** INPAINTING(maskIn : HxW, dataIn : HxWxC)
2:     **if** H $\leq$ 1 or W $\leq$ 1 or all pixels are filled **then**
3:         **return** maskIn, dataIn         ▷ end of recursion
4:     **end if**
    *weighted area downsampling:*
5:     maskLow, dataLow = zeros of shape $\frac{H}{2}$x$\frac{W}{2}$, $\frac{H}{2}$x$\frac{W}{2}$x$C$
6:     **for** $i, j \in \{0, ..., \frac{H}{2} - 1\} \times \{0, ..., \frac{W}{2} - 1\}$ **do**
7:         $N_{\max}, N_{\text{avg}}, d = \mathbf{0}^C$
8:         **for** $a, b \in \{2i, 2i+1\} \times \{2j, 2j+1\}$ **do**   ▷ loop over neighbors in the fine grid
9:             $N_{\max} = \max\{N_{\max}, \text{maskIn}[a, b]\}$
10:           $N_{\text{avg}} \mathrel{+}= \text{maskIn}[a, b]$
11:           $d \mathrel{+}= \text{maskIn}[a, b] \cdot \text{dataIn}[a, b, :]$
12:         **end for**
13:         **if** $N_{\text{avg}} > 0$ **then**
14:             maskLow$[i, j] = N_{\max}$
15:             dataLow$[i, j, :] = d / N_{\text{avg}}$
16:         **end if**
17:     **end for**
    *recursion:*
18:     maskLow, dataLow = INPAINTING(maskLow, dataLow)
    *weighted bilinear upsampling:*
19:     maskOut, dataOut = zeros of shape HxW, HxWxC
20:     **for** $a, b \in \{0, ..., H-1\} \times \{0, ..., W-1\}$ **do**
21:         $N, W = 0, d = \mathbf{0}^C$
22:         $\hat{a} = a \div 2, \hat{b} = b \div 2$     ▷ Integer-division (round down), indices on the coarse grid
23:         $a', b' = -1$ if $a, b$ is even else $+1$
24:         $N = \{(\hat{a}, \hat{b}, \frac{9}{16}), (\hat{a}+a', \hat{b}, \frac{3}{16}), (\hat{a}, \hat{b}+b', \frac{3}{16}), (\hat{a}+a', \hat{b}+b', \frac{1}{16})\}$
25:         **for** $(i, j, w) \in N \cap$ image **do** ▷ loop over neighbors if within bounds
26:             $N \mathrel{+}= w\, \text{maskLow}[i, j]$
27:             $d \mathrel{+}= w\, \text{maskLow}[i, j] \cdot \text{dataLow}[i, j, :]$
28:             $W \mathrel{+}= w$
29:         **end for**
30:         maskOut$[a, b]$ = maskIn$[i, j]$, dataOut$[a, b, :]$ = maskIn$[i, j] \cdot \text{dataIn}[i, j, :]$
31:         **if** $N > 0$ **then** ▷ blend interpolated values with original data
32:             maskOut$[a, b] \mathrel{+}= (1 - \text{maskIn}[i, j])\, N/W$
33:             dataOut$[a, b, :] \mathrel{+}= (1 - \text{maskIn}[i, j])\, d/N$
34:         **end if**
35:     **end for**
36:     **return** maskOut, dataOut
37: **end function**

---

be defined as the gradient of $s$ with respect to $x_i$ and $w_i$, $\hat{x}_i := \nabla_{x_i} s, \hat{w}_i := \nabla_{w_i} s$ as column vectors. Next, we drop the index $i$, as we require it to index the elements in the input and output vectors, and look at a single function $f \in \mathbb{R}^N \times \mathbb{R}^W \to \mathbb{R}^M$ with inputs $x \in \mathbb{R}^N, w \in \mathbb{R}^W$ and output $y \in \mathbb{R}^M$. The adjoint variables are then computed using

$$\hat{x} = J_{f,x}^T(x, w)\hat{y} , \quad \hat{w} = J_{f,w}^T(x, w)\hat{y}. \tag{15}$$

Here, the Jacobian matrix with respect to the different inputs is used, defined as

$$(J_{f,x})_{ij} := \frac{\partial f_i}{\partial x_j} , \quad (J_{f,w})_{ij} := \frac{\partial f_i}{\partial w_j}. \tag{16}$$

As one can see, given the adjoint variable of the output $\hat{y}$, the adjoint method propagates these gradients back through the derivatives of $f$ to the adjoint variables of the inputs $\hat{x}$ and $\hat{w}$. In the context of the chained function Equation 14, this implies that,

starting with gradients on the output $\hat{x}_n$ from the loss function, gradients are first propagated to $\hat{x}_{n-1}, \hat{w}_n$ via $J_{f_i}$, then to $\hat{x}_{n-2}, \hat{w}_{n-1}$, and so on until $\hat{x}_0, \hat{w}_1$ is reached.

To provide custom differentiable operations, two functions have to be provided: first, the forward code $y \leftarrow f(x, w)$ with input $x$ and parameter $w$, and second, the backward code to compute $\hat{x}$ and $\hat{w}$ from $\hat{y}$, possibly using $x, w$ from the forward pass again to compute the Jacobian.

## E.2 Backward Pass of the Sampling Function

Using the theory above, we now present the adjoint code for the differentiable sampling from subsection 3.2. This serves to highlight what is differentiated and how the gradients are propagated. Note that these functions are implemented based on PyTorch functions, PyTorch can automatically compute the derivatives.

The differentiable sampling stage takes the importance map $I$ as input and produces the image of sparse samples $S$. In the framework of Equation 14, this can be seen as block of functions that is cut out in the middle. As parameters, the target mean $\mu$ and lower bound $l$, the sampling steepness $\alpha$, the sample pattern $P$ and the target image $T$ are used. Note that no optimization with respect to these parameters is performed, their respective adjoint variables are unused. To recapitulate, the sampling is performed using the following steps:

$$\mu_I = \frac{1}{WH} \sum_{ij} I_{ij} , \quad I^{(1)} = I \tag{17a}$$

$$I'_{ij} = \min\left\{1, l + I_{ij}^{(1)} \frac{\mu - l}{\mu_I + \varepsilon}\right\} \tag{17b}$$

$$S_{ij} = \text{sig}(\alpha(I'_{ij} - P_{ij}))T_{ij} \text{ with } \text{sig}(x) = \frac{1}{1 + e^{-x}}. \tag{17c}$$

Note that the second and third function act on each pixel $ij$ of the images independently. Therefore, we use them as per-element functions to simplify the notation of the derivatives. Using the matrix notation from Equation 15, this would imply a diagonal Jacobian. Furthermore, $T_{ij}$ and $S_{ij}$ return the vector of channels at the specified location. In order to stay within the presented framework of the adjoint method, if variables are used by a function and later again by another function, these variables are passed through as additional outputs ($I^{(1)} = I$).

For the backward pass, we are given the gradients of the output $\hat{S}$ from the backward pass of the reconstruction. This equates to $\hat{x}_n$ in Equation 14. Then the gradients are propagated through the sampling algorithm in reverse order:

$$\begin{aligned}
\hat{T}_{ij} &= \text{sig}(\alpha(I'_{ij} - P_{ij}))\hat{S}_{ij} \\
\hat{P}_{ij} &= -(\alpha T_{ij}\, \text{sig}'\left(\alpha(I'_{ij} - P_{ij})\right))^T \hat{S}_{ij} \\
\hat{I}'_{ij} &= (\alpha T_{ij}\, \text{sig}'\left(\alpha(I'_{ij} - P_{ij})\right))^T \hat{S}_{ij}
\end{aligned} \tag{18a}$$

with $\text{sig}'(x) = \frac{d}{dx}\, \text{sig}(x) = \text{sig}(x)\, \text{sig}(-x)$

$$\hat{I}_{ij}^{(1)} = \begin{cases} \frac{\mu - l}{\mu_I + \varepsilon}\hat{I}'_{ij} , & I'_{ij} < 1 \\ 0 , & I'_{ij} \geq 1 \end{cases} \tag{18b}$$

$$\hat{\mu}_l = \sum_{ij} \left( \begin{cases} \frac{I_{ij}^{(1)}(l - \mu)}{(\mu_I + \varepsilon)^2}\hat{I}'_{ij} , & I'_{ij} < 1 \\ 0 , & I'_{ij} \geq 1 \end{cases} \right)$$

derivatives for $\mu, l, \varepsilon$ are omitted

$$\hat{I}_{ij} = \hat{I}_{ij}^{(1)} + \frac{1}{WH}\hat{\mu}_l \tag{18c}$$

## E.3 Backward Pass of the Pull-Push Algorithm

As one can see in the previous section, deriving the adjoint code is done mechanically by deriving each line of code with respect to the inputs. This, however, produces a vastly longer code, therefore, we only outline the steps to derive the adjoint code of the pull-push algorithm Alg. 1. The full source code is available in the online repository.

The algorithm is a recursive algorithm with three stages: the downsampling to the coarse level, the recursive call, and the up-sampling and interpolation with the fine level. During the backward pass, the order is reversed. First, the adjoint of the upsampling and interpolation at the finest level. Then the adjoint of the recursive call, which itself is the adjoint of upsampling, recursion and downsampling. And lastly the adjoint of the downsampling.

## REFERENCES

[1] S. J. Gortler, R. Grzeszczuk, R. Szeliski, and M. F. Cohen. The lumigraph. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pp. 43–54, 1996.

[2] J. H. Halton. Algorithm 247: Radical-inverse quasi-random point sequence. *Commun. ACM*, 7(12):701702, Dec. 1964. doi: 10.1145/355588.365104

[3] M. Kraus. The pull-push algorithm revisited. *Proceedings GRAPP*, 2:3, 2009.

[4] J. L. Lions. *Optimal control of systems governed by partial differential equations*. Springer, 1971.

[5] A. McNamara, A. Treuille, Z. Popović, and J. Stam. Fluid control using the adjoint method. *ACM Transactions On Graphics (TOG)*, 23(3):449–456, 2004.

[6] M. Roberts. The unreasonable effectiveness of quasirandom sequences. http://extremelearning.com.au/unreasonable-effectiveness-of-quasirandom-sequences/, 2020. Accessed: 2020-02-14.

[7] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In N. Navab, J. Hornegger, W. M. Wells, and A. F. Frangi, eds., *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, pp. 234–241. Springer International Publishing, Cham, 2015.

[8] G. Schwarz et al. Estimating the dimension of a model. *The annals of statistics*, 6(2):461–464, 1978.

[9] Y. Wang, W. Chen, J. Zhang, T. Dong, G.-Y. Shan, and X. Chi. Efficient volume exploration using the gaussian mixture model. *Visualization and Computer Graphics, IEEE Transactions on*, 17:1560 – 1573, 12 2011. doi: 10.1109/TVCG.2011.97

[10] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018.