

# Perceptual Deep Depth Super-Resolution

Oleg Voynov<sup>1</sup>, Alexey Artemov<sup>1</sup>, Vage Egiазarian<sup>1</sup>, Alexander Notchenko<sup>1</sup>,  
Gleb Bobrovskikh<sup>1,2</sup>, Denis Zorin<sup>3,1</sup>, Evgeny Burnaev<sup>1</sup>

<sup>1</sup>Skolkovo Institute of Science and Technology, <sup>2</sup>Higher School of Economics,

<sup>3</sup>New York University

{oleg.voinov, a.artemov, vage.egiazarian, alexandr.notchenko}@skoltech.ru,

bobrovskikh@gmail.com, dzorin@cs.nyu.edu, e.burnaev@skoltech.ru

[adase.group/3ddl/projects/perceptual-depth-sr](https://adase.group/3ddl/projects/perceptual-depth-sr)

## Abstract

RGBD images, combining high-resolution color and lower-resolution depth from various types of depth sensors, are increasingly common. One can significantly improve the resolution of depth maps by taking advantage of color information; deep learning methods make combining color and depth information particularly easy.

However, fusing these two sources of data may lead to a variety of artifacts. If depth maps are used to reconstruct 3D shapes, e.g., for virtual reality applications, the visual quality of upsampled images is particularly important.

The main idea of our approach is to measure the quality of depth map upsampling using renderings of resulting 3D surfaces. We demonstrate that a simple visual appearance-based loss, when used with either a trained CNN or simply a deep prior, yields significantly improved 3D shapes, as measured by a number of existing perceptual metrics. We compare this approach with a number of existing optimization and learning-based techniques.

## 1. Introduction

RGBD images are increasingly common as sensor technology becomes more widely available and affordable. They can be used for reconstruction of the 3D shapes of objects and their surface appearance. The better the quality of the depth component, the more reliable the reconstruction.

Unfortunately, for most methods of depth acquisition the resolution and quality of the depth component is insufficient for accurate surface reconstruction. As the resolution of the RGB component is usually several times higher and there is a high correlation between structural features of the color image and the depth map (e.g., object edges) it is natural to use the color image for depth map super-resolution, i.e. up-sampling of the depth map. Convolutional neural networks

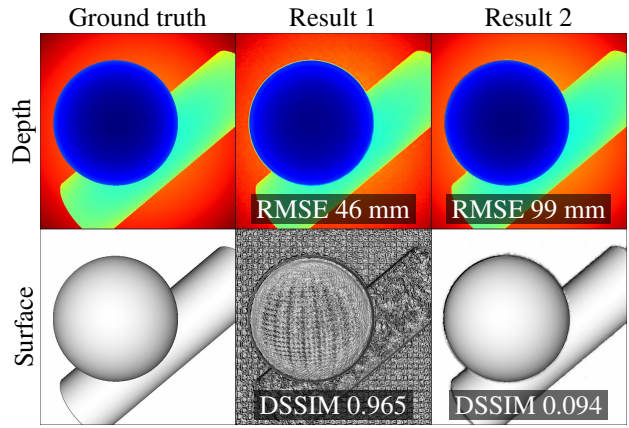


Figure 1: Visually inferior super-resolution result in the middle gets higher score according to direct depth deviation but lower score according to perceptual deviation of the rendered image of the 3D surface. While the surfaces differ significantly, the corresponding depth maps do not capture this difference and look almost identical.

are a natural fit for this problem as they can easily fuse heterogeneous information.

A critical aspect of any upsampling method is the measure of quality it optimizes (i.e., the loss function), whether the technique is data-driven or not. In this paper we focus on applications that require reconstruction of 3D geometry visible to the user, like acquisition of realistic 3D scenes for virtual or augmented reality and computer graphics. In these applications the *visual* appearance of the resulting 3D shape, i.e., how the surface looks when observed under various lighting conditions, is of particular importance.

Most existing research on depth super-resolution is dominated by simple measures based on pointwise deviation of depth values. However, direct pointwise difference of the depth maps do not capture the visual difference between

the corresponding 3D shapes: for example, low-amplitude high-frequency variations of depth may correspond to significant difference in appearance, while conversely, relatively large smooth changes in depth may be perceptually less relevant, as illustrated in Figure 1.

Hence, we propose to compare the rendered images of the surface instead of the depth values directly. In this paper we explore depth map super-resolution using a simple loss function based on visual differences. Our loss function can be computed efficiently and is shown to be highly correlated with more elaborate perceptual metrics. We demonstrate that this simple idea used with two deep learning-based RGBD super-resolution algorithms results in a dramatic improvement of visual quality according to perceptual metrics and an informal perceptual study. We compare our results with six state-of-the-art methods of depth super-resolution that are based on distinct principles and use several types of loss functions.

In summary, our contributions are as follows: (1) we demonstrate that a simple and efficient visual difference-based metric for depth map comparison can be, on the one hand, easily combined with neural network-based whole-image upsampling techniques, and, on the other hand, is correlated with established proxies for human perception, validated with respect to experimental measurements; (2) we demonstrate with extensive comparisons that with the use of this metric two methods of depth map super-resolution, one based on a trainable CNN and the other based on the deep prior, yield high-quality results as measured by multiple perceptual metrics. To the best of our knowledge, our paper is the first to systematically study the performance of visual difference-based depth super-resolution across a variety of datasets, methods, and quality measures, including a basic human evaluation.

Throughout the paper we use the term *depth map* to refer to the depth component of an RGBD image, and the term *normal map* to refer to the map of the same resolution with the 3D surface normal direction computed from the depth map at each pixel. Finally, the *rendering of a depth map* refers to the grayscale image obtained by constructing a 3D triangulation of the height field represented by the depth map, via computing the normal map from this triangulation, and rendering it using fixed material properties and a choice of lighting. This is distinct from a commonly used depth map visualization with grayscale values obtained from the depth values by simple scaling. We describe this in more detail in Section 3.

## 2. Related work

### 2.1. Image quality measures

Quality measures play two important roles in image super-resolution: on the one hand, they are used to formu-

late an optimization functional or a loss function, on the other hand, they are used to evaluate the quality of the results. Ideally, the same function should serve both purposes, however, in some instances it may be optimal to choose different functions for evaluation and optimization. While in the former case the top priority is to capture the needs of the application, in the latter case the efficiency of evaluation and differentiability are significant considerations.

In most works on depth map reconstruction and upsampling a limited number of simple metrics are used, both for optimization and final evaluation. Typically these are scaled  $L_2$  or  $L_1$  norms of depth deviations (see *e.g.* [9]).

Another set of measures introduced in [19, 20] and primarily used for evaluation, not optimization or learning, consists of heuristic measures of various aspects of the depth map geometry: foreground flattening/thinning, fuzziness, bumpiness, etc. Most of them require a very specific segmentation of the image for detection of flat areas and depth discontinuities.

Visual similarity measures, well-established in the area of photo-processing, aim to be consistent with human judgment, in the sense of similarity ordering (which of the two images is more similar to the ground truth?). The examples include (1) the metrics based on simple vision models of *structural similarity* SSIM [52], FSIM [57], MSSIM [53], (2) based on a sophisticated model of low-level visual processing [35], or (3) on convolutional neural networks (see [58] for a detailed overview). The latter use a simple distance measure on deep features learned for an image understanding task, *e.g.*  $L_2$  distance on the features learned for image classification, and have been demonstrated to outperform statistical measures such as SSIM.

### 2.2. Depth super-resolution

Depth super-resolution is closely related to a number of depth processing tasks, such as denoising, enhancement, inpainting, and densification (*e.g.*, [5, 6, 8, 21, 33, 34, 45, 46, 55]). We directly focus on the problem of super-resolution, or more specifically, estimation of high-resolution depth map from a single low-resolution depth map and a high-resolution RGB image.

**Convolutional neural networks** have achieved most impressive performance among learning-based methods in high-level computer vision tasks and recently have been applied to depth super-resolution [22, 30, 39, 43]. One approach [22] is to resolve ambiguity in the depth map upsampling by explicitly adding high-frequency features from high-resolution RGB data. Another, hybrid approach [39, 43] is to add a subsequent optimization stage to a CNN to produce sharper results. Different approaches to CNN-based photo-guided depth super-resolution include linear filtering with CNN-derived kernels [26], deep fusion of time-of-flight depth and stereo images [1], and generative

adversarial networks [62].

These techniques use either  $L_2$  or  $L_1$  norm of the depth differences as the basis of their loss functions, often combined with regularizers of different types. The recent approach of [62] is the closest to ours: it uses the difference of gradients as one of the loss terms to capture some of the visual information. For evaluation, these works report root mean square error (RMSE), mean absolute error (MAE), peak signal-to-noise ratio (PSNR), all applied directly to depth maps, and, rarely [4, 43, 44, 62], perceptual SSIM *also applied directly to depth maps*. In contrast, we propose to measure the perceptual quality of depth map *renderings*.

**Dictionary learning** has also been investigated for depth super-resolution [11, 13, 29], however, compared to CNNs, it is typically restricted to smaller dimensions and as a result to structurally simpler depth maps.

**Variational approach** aims to combine RGB and depth information explicitly by carefully designing an optimization functional, without relying on learning. Most relevant examples employ shape-from-shading problem statement for single-image [14] or multiple-image [38] depth super-resolution. These works include visual difference-related terms in the optimized functional and report normal deviation, capturing visual similarity. While showing impressive results in many cases, they typically require prior segmentation of foreground objects and depend heavily on the quality of such segmentation.

Another strategy to tackle ambiguities in super-resolution is to design sophisticated regularizers to balance the data-fidelity terms against a structural image prior [15, 24, 56]. In contrast to this approach, which requires custom hand-crafted regularized objectives and optimization procedures, we focus on the standard training strategy (*i.e.*, gradient-based optimization of a CNN) while using a loss function that captures visual similarity.

Yet another approach is to choose a carefully-designed model such as [63] featuring a sophisticated metric defined in a space of minimum spanning trees and including an explicit edge inconsistency model. In contrast to ours, such model requires manual tuning of multiple hyperparameters.

### 2.3. Perceptual photo super-resolution

Perceptual metrics have been considered more broadly in the context of photo processing. While convolutional neural networks for photo super-resolution trained with simple mean square or mean absolute color deviation keep demonstrating impressive results [16, 18, 59, 60], it has been widely recognized that pixelwise difference of color image data is not well correlated with perceptual image difference. For this reason, relying on a pixelwise color error may lead to suboptimal performance.

One solution is to instead use the loss function represented by the deviation of the features from a neural net-



Figure 2: Depth map renderings generated with four light directions that we use for metric calculation.

work trained for an image understanding task [25]. This idea can be further combined with an adversarial training procedure to push the super-resolution result to the natural image manifold [28]. Another extension to this idea is to train the neural network to generate images with natural distribution of statistical features [12, 36, 50, 51]. To balance between the perceptual quality and pixelwise color deviation, generative adversarial networks can be used [7, 31, 49].

Another solution is to learn a quality measure from perceptual scores, collected from a human subject study, and use this quality measure as the loss function. Such quality measure may capture similarity of two images [58] or an absolute naturalness of the image [32].

### 3. Metrics

In this section, we discuss visually-based metrics and how they can be used to evaluate the quality of depth map super-resolution and as loss functions. The general principle we follow is to apply comparison metrics to *renderings* of the depth maps to obtain a measure of their difference instead of considering depth maps directly. The difficulty with this approach is that there are infinitely many possible renderings depending on lighting conditions, material properties and camera position. However, we demonstrate that even a very simple rendering procedure already yields substantially improved results. We label visually-based metrics with subscript “v” and the metrics that compare the depth values directly with subscript “d”.

**From depth map to visual representation.** To approximate the appearance of a 3D scene depicted with a certain depth map we use a simple rendering procedure. We illuminate the corresponding 3D surface with monochromatic directional light source and observe it with the same camera that the scene was originally acquired with. We use the diffuse reflection model and do not take visibility into account. For this model, the intensity of a pixel  $(i, j)$  of the rendering  $I$  is proportional to cosine of the angle between the normal at the point of the surface corresponding to the pixel  $\mathbf{n}_{ij}$  and direction to the light source  $\mathbf{e}$ :  $I_{ij} = \mathbf{e} \cdot \mathbf{n}_{ij}$ . We calculate the normals from the depth maps using first-order finite-differences. Any number of vectors  $\mathbf{e}$  can be used to generate a collection of renderings representing the depth map, however, any rendering can be obtained as a



linear combination of three basis ones corresponding to independent light directions. Renderings for different light directions are presented in Figure 2.

**Perceptual metrics.** We briefly describe two representative metrics: a statistics-based DSSIM, and a neural network-based LPIPS. Either of these can be applied to three basis renderings (or a larger sample of renderings) and reduced to obtain the final value. While, in principle, they can also be used as loss functions, the choice of a loss function needs to take stability and efficiency into account, so we opt for a more conservative choice described below.

*Structural similarity index measure (SSIM)* [52] takes into account the changes in the local structure of an image, captured by statistical quantities computed on a small window around each pixel. For each pair of pixels of the compared images  $I_k$ ,  $k = 1, 2$  the luminance term  $\ell$ , the contrast term  $c$  and the structural term  $s$ , each normalized, are computed using the means  $\mu_k$ , standard deviations  $\sigma_k$  and cross-covariance  $\sigma_{12}$  of the pixels in the corresponding local windows. The value of SSIM is then computed as pixelwise mean product of these terms

$$\ell = \frac{2\mu_1\mu_2}{\mu_1^2 + \mu_2^2}, c = \frac{2\sigma_1\sigma_2}{\sigma_1^2 + \sigma_2^2}, s = \frac{\sigma_{12}}{\sigma_1\sigma_2}, \quad (1)$$

$$\text{SSIM}_v(I_1, I_2) = \frac{1}{N} \sum_{ij} \ell_{ij} \cdot c_{ij} \cdot s_{ij},$$

where  $N$  is the number of pixels. Dissimilarity measure can be computed as  $\text{DSSIM}_v(I_1, I_2) = 1 - \text{SSIM}_v(I_1, I_2)$ .

*Neural net-based metrics* rely on the idea of measuring the distance between features extracted from a neural network. Specifically, feature maps  $\mathbf{x}_{k\ell}$ ,  $\ell = 1 \dots L$  with spatial dimensions  $H_\ell \times W_\ell$  are extracted from  $L$  layers of the network for each of the compared images. In the simplest case, the metric value is then computed as pixelwise mean square difference of the feature maps, summed over the layers

$$\text{NN}_v(I_1, I_2) = \sum_{\ell} \frac{1}{H_\ell W_\ell} \sum_{ij} \|\mathbf{x}_{1\ell,ij} - \mathbf{x}_{2\ell,ij}\|_2^2. \quad (2)$$

*Learned perceptual image patch similarity (LPIPS)* [58] adds a learned channel-wise weighting to the above formula and uses 5 layers from Alexnet [27] or VGG [41] or the first layer from SqueezeNet [23] as the CNN of choice.

**Our visual difference-based metric.** While the metrics described above are good proxies for human evaluation of difference between depth map renderings, they are lacking as loss functions due to their complex landscapes. Optimization with DSSIM as the loss function may produce the results actually inferior with respect to *DSSIM itself* compared to a simpler loss function we define below, as illustrated in Figure 3. LPIPS has a complex energy profile typical for neural networks, and having a neural network as the loss function for another may behave unpredictably [61].

The simplest metric capturing the difference between all

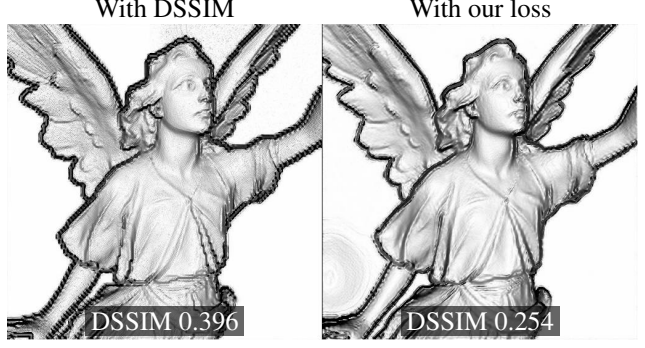


Figure 3: Optimization with DSSIM as the loss function may produce the results inferior with respect to *DSSIM itself* compared to our simpler loss function.

possible renderings of the depth maps  $d_k$  can be computed as the average root mean square deviation of three basis renderings  $\mathbf{e}_m \cdot \mathbf{n}_k$  in an orthogonal basis  $\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3$

$$\text{RMSE}_v(d_1, d_2) = \sqrt{\text{MSE}_v(d_1, d_2)},$$

$$\text{MSE}_v(d_1, d_2) = \frac{1}{3N} \sum_{ij,m} \|\mathbf{e}_m \cdot \mathbf{n}_{1,ij} - \mathbf{e}_m \cdot \mathbf{n}_{2,ij}\|_2^2, \quad (3)$$

similarly to RMS difference of the normal maps.

We found that this simple metric for depth map comparison is efficient and stable as the loss function and at the same time, as we demonstrate in Section 5, it is well correlated with DSSIM and LPIPS, *i.e.*, situations when the value of one metric is high and the value of another is low are unlikely. Our experiments confirm that optimization of this metric also improves both perceptual metrics.

## 4. Methods

We selected eight representative state-of-the-art depth processing methods based on different principles: (1) a purely variational method [14], (2) a bilateral filtering method that uses a high-resolution edge map [54], (3) a dictionary learning method [13], (4) a hybrid CNN-variational method [39], (5) a pure CNN [22], (6) a zero-shot CNN [47], (7) a densification [34] and (8) an enhancement [55] CNNs. Our goals were (a) to modify the methods for using with the visual difference-based loss function, and (b) to compare the results of the modified methods with alternatives of different types. In our experiments the last two methods did not perform well compared to others, so we did not consider them further. We found that two neural network-based methods (5) and (6), that we refer to as MSG and DIP, can be easily modified for using with a visual difference-based loss function, as we explain now.

**MSG** [22] is a deep learning method that uses different strategies to upsample different spectral components of low-resolution depth map. In the modified version of this method, that we denote by **MSG-V**, we replaced the



original loss function with a combination of our visual difference-based metric and mean absolute deviation of Laplacian pyramid  $\text{Lap}_1$  [2] as a regularizer

$$\mathcal{L}(d_1, d_2) = \text{Lap}_1(d_1, d_2) + w \cdot \text{MSE}_v(d_1, d_2). \quad (4)$$

**DIP** [47] is a zero-shot deep learning approach, based on a remarkable observation that, even without any specialized training, the structure of CNN itself may be leveraged for solving inverse problems on images. We note that this approach naturally allows simultaneous super-resolution and inpainting. In this approach, the depth super-resolution problem would be formulated as

$$d_{\theta^*}^{\text{SR}} = \text{CNN}_{\theta^*}, \quad \theta^* = \arg \min_{\theta} \text{MSE}_d(\mathbf{D}d_{\theta}^{\text{SR}}, d^{\text{LR}}), \quad (5)$$

where  $d^{\text{LR}}$  and  $d_{\theta^*}^{\text{SR}}$  are the low-resolution and super-resolved depth maps,  $\text{CNN}_{\theta}$  is the output of the deep neural network parametrised by  $\theta$ ,  $\mathbf{D}$  is the downsampling operator, and  $\text{MSE}_d$  is direct mean square difference of the depth maps. To perform photo-guided super-resolution, we added a second output channel for intensity to the network

$$d_{\theta^*}^{\text{SR}} = \text{CNN}_{\theta^*}^{(1)}, \quad I_{\theta} = \text{CNN}_{\theta}^{(2)}, \quad (6)$$

$$\theta^* = \arg \min_{\theta} \text{MSE}_d(\mathbf{D}d_{\theta}^{\text{SR}}, d^{\text{LR}}) + w_I \cdot \text{Lap}_1(I_{\theta}, I^{\text{HR}}),$$

where  $I^{\text{HR}}$  is the high-resolution photo guidance, and for visually-based version **DIP-V** we further replaced the direct depth deviation  $\text{MSE}_d$  with the function from Equation 4.

We used the remaining four methods (1)-(4) for comparison as-is, as modifying them for a different loss function would require substantial changes to the algorithms.

**SRfS** [14] is a variational method relying on complementarity of super-resolution and shape-from-shading problems. It already includes a visual-difference based term (the remaining methods use depth difference metrics). **EG** [54] approaches the problem via prediction of smooth high-resolution depth edges with Markov random field optimization. It does not use a loss directly, therefore cannot be easily adapted. **DG** [13] is a depth map enhancement method based on dictionary learning that uses depth difference-based fidelity term. It makes a number of modeling choices which may not be suitable for a different loss function, and typically does not perform as well as neural network-based methods. **PDN** [39] is a hybrid method featuring two stages: the first is composed of fully-convolutional layers and predicts a rough super-resolved depth map, and the second performs an unrolled variational optimization, aiming to produce a sharp and noise-free result.

## 5. Experiments

### 5.1. Data

For evaluation we selected a representative and diverse set of 34 RGBD images featuring synthetic, high-quality real and low-quality real data with different levels of geo-

metric and textural complexity. We employed four datasets, most common in literature on depth super-resolution. *ICL-NUIM* [17] includes photo-realistic RGB images along with synthetic depth, free from any acquisition noise. *Middlebury 2014* [40], captured with a structured light system, provides high-quality ground truth for complex real-world scenes. *SUN RGBD* [42] contains images captured with four different consumer-level RGBD cameras: Intel RealSense, Asus Xtion, Microsoft Kinect v1 and v2. *ToF-Mark* [10] provides challenging real-world time-of-flight and intensity camera acquisitions together with an accurate ground truth from a structured light sensor.

In addition, we constructed a synthetic *SimGeo* dataset, that consists of 6 geometrically simple scenes with low- and high-frequency texture, and without any, using Blender. The purpose of *SimGeo* were to reveal artifacts that are not related to the noise or high-frequency geometry in the input data, like false geometric detail caused by color variation on a smooth surface.

We resized and cropped each RGBD image to the resolution of  $512 \times 512$  and generated low-resolution input depth maps with the scaling factors of 4 and 8, that are most common among the works on depth super-resolution. We focused on two downsampling models: Box, *i.e.*, each low-resolution pixel contains the mean value over the “box” neighbouring high-resolution pixels, and Nearest neighbour, *i.e.*, each low-resolution pixel contains the value of the nearest high-resolution pixel. For additional details on our evaluation data and the results for different downsampling models please refer to supplementary material.

### 5.2. Evaluation details

To quantify the performance of the methods, we measured direct RMS deviation of the depth maps (denoted by  $\text{RMSE}_d$ ) and deviation of their renderings with the metrics described in Section 3. For visually-based metrics we calculated their values for three orthogonal light directions, corresponding to the three left-most images in Figure 2, and the value for an additional light direction, corresponding to the right-most image. We then took the worst of the four values. With similar outcomes, we also explored different reducing strategies and a set of different metrics: BadPix and Bumpiness, applied directly to depth values, and BadPix and RMSE applied to separate depth map renderings.

Additionally, we conducted an informal perceptual study using the results on *SimGeo*, *ICL-NUIM* and *Middlebury* datasets, in which subjects were asked to choose the renderings of the upsampled depth maps that look most similar to the ground truth.

### 5.3. Implementation details

We evaluated publicly available trained models for EG, DG, and MSG and trained PDN using publicly available

code; we used the implementation of SRfS provided by the authors; we adapted publicly available implementation of DIP for depth maps, as described in Section 4; we reimplemented MSG-V in PyTorch [37] and trained it according to the original paper using the patches from Middlebury and MPI Sintel [3]. We selected the value of the weighting parameter  $w$  in Equation (4) so that both terms of the loss contribute equally with respect to their magnitudes (see supplementary material for more details).

#### 5.4. Comparison of quality measures

To quantify how well different metrics represent the visual quality of a super-resolved depth map, we compared pairwise correlations of these metrics and calculated the corresponding values of Pearson correlation coefficient. Since LPIPS as a neural network-based perceptual metric has been experimentally shown to represent human perception well, we used its value as the reference. We found that the metrics based on direct depth deviation demonstrate weak correlation with perceptual metrics, as illustrated in Figure 4 for  $RMSE_d$ , and hence are not suitable for measuring the depth map quality when the visual appearance plays an important role. On the other hand, we found that our  $RMSE_v$  correlates well with perceptual metrics, to the same extent they correlate with each other (see Figure 4).

#### 5.5. Comparison of super-resolution methods

In Table 1 and Figure 5 we present the super-resolution results on our SimGeo dataset with the scaling factor of 4; in Table 2 and Figure 6 we present the results on ICL-NUIM and Middlebury datasets with the scaling factors of 4 and 8. We use Box downsampling model in both cases. Please find the additional results in supplementary material or online<sup>1</sup>.

In general, we found that the methods EG, PDN and DG do not recover fine details of the surface, typically over-smoothing the result in comparison to, *e.g.*, Bicubic up-sampling, the methods SRfS and original DIP suffer from false geometry artifacts in case of a smooth textured surface, and original MSG introduces severe noise around the depth edges. As illustrated in Figure 6 and Table 2, all the methods from prior works perform relatively poorly on the images with regions of missing depth measurements (rendered in black), including the ones that inpaint these regions explicitly (SRfS, DG) or implicitly (DIP). The method EG failed to converge on some images.

In contrast, we observed that integration of our visual difference-based loss into DIP and MSG significantly improved the results of both methods qualitatively and quantitatively. The visual difference-based version DIP-V do not suffer from false geometry artifacts as much as the original version. On the challenging images from Middlebury dataset, where it performed simultaneous super-resolution

and inpainting, DIP-V mostly outperformed other methods as measured by the perceptual metrics and was preferred by more than 80% of subjects in the perceptual study. The visual difference-based version MSG-V produces significantly less noisy results in comparison to the original version, in some cases almost without any noticeable artifacts. On the data without missing measurements, including hole-filled “Vintage” from Middlebury, MSG-V mostly outperformed other methods as measured by the perceptual metrics and was preferred by more than 80% of subjects. On SimGeo, ICL-NUIM and Middlebury combined, one of our modified versions, DIP-V or MSG-V, was preferred over the other methods by more than 85% of subjects.

For reference, in Figure 5 we include pseudo-color visualizations of the depth maps. Notice that while the up-sampled depth maps obtained with different methods are almost indistinguishable in this form of visualization, commonly used in the literature on depth processing for qualitative evaluation, the corresponding renderings and, consequently, the underlying geometry varies dramatically.

## 6. Conclusion

We have explored depth map super-resolution with a simple visual difference-based metric as the loss function. Via comparison of this metric with a variety of perceptual quality measures, we have demonstrated that it can be considered a reasonable proxy for human perception in the problem of depth super-resolution with the focus on visual quality of the 3D surface. Via an extensive evaluation of several depth-processing methods on a range of synthetic and real data, we have demonstrated that using this metric as the loss function yields significantly improved results in comparison to the common direct pixel-wise deviation of depth values. We have combined our metric with relatively simple and non-specific deep learning architectures and expect that this approach will be beneficial for other related problems.

We have focused on the case of single regularly sampled RGBD images, but a lot of geometric data has less regular form. The future work would be to adapt the developed methodology to a more general sampling of the depth values for the cases of multiple RGBD images or point clouds annotated with a collection of RGB images.

## Acknowledgements

The work was supported by The Ministry of Education and Science of Russian Federation, grant No. 14.615.21.0004, grant code: RFMEFI61518X0004.

The authors acknowledge the usage of the Skoltech CDISE HPC cluster Zhores for obtaining the results presented in this paper.

<sup>1</sup>[mega.nz/#F!yvRXBABI!pucRoBvtnthzHI!oqsxEvA!y6JmCajS](https://mega.nz/#F!yvRXBABI!pucRoBvtnthzHI!oqsxEvA!y6JmCajS)

	Sphere and cylinder, x4				Lucy, x4				Cube, x4				SimGeo average, x4			
	RMSE <sub>d</sub>	DSSIM <sub>v</sub>	LPIPS <sub>v</sub>	RMSE <sub>v</sub>	RMSE <sub>d</sub>	DSSIM <sub>v</sub>	LPIPS <sub>v</sub>	RMSE <sub>v</sub>	RMSE <sub>d</sub>	DSSIM <sub>v</sub>	LPIPS <sub>v</sub>	RMSE <sub>v</sub>	RMSE <sub>d</sub>	DSSIM <sub>v</sub>	LPIPS <sub>v</sub>	RMSE <sub>v</sub>
SRfS [14]	70	887	1025	417	82	811	781	367	52	934	1036	361	61	711	869	311
EG [54]	55	<u>143</u>	326	<u>130</u>	69	357	426	<u>220</u>	43	<u>113</u>	<u>214</u>	<u>105</u>	53	<u>168</u>	306	<u>136</u>
PDN [39]	157	198	295	150	173	456	<u>368</u>	251	164	156	250	145	162	224	<u>278</u>	165
DG [13]	56	265	372	166	69	523	558	249	44	218	411	139	54	293	<u>420</u>	171
Bicubic	57	189	313	189	72	<u>355</u>	398	267	44	131	287	160	55	197	320	193
DIP [47]	46	965	1062	548	<u>53</u>	827	615	344	45	963	906	530	52	887	893	395
MSG [22]	<u>41</u>	626	859	229	54	444	480	259	<u>29</u>	445	687	176	<u>39</u>	374	569	194
DIP-V	<b>28</b>	560	766	142	<b>44</b>	421	446	223	<b>26</b>	352	613	146	<b>33</b>	313	524	147
MSG-V	99	<b>94</b>	<b>267</b>	<b>96</b>	74	<b>205</b>	<b>251</b>	<b>156</b>	102	<b>70</b>	<b>179</b>	<b>77</b>	96	<b>95</b>	<b>194</b>	<b>99</b>

Table 1: Quantitative evaluation on SimGeo dataset. RMSE<sub>d</sub> is in millimeters, other metrics are in thousandths. Lower values correspond to better results. The best result is in bold, the second best is underlined.

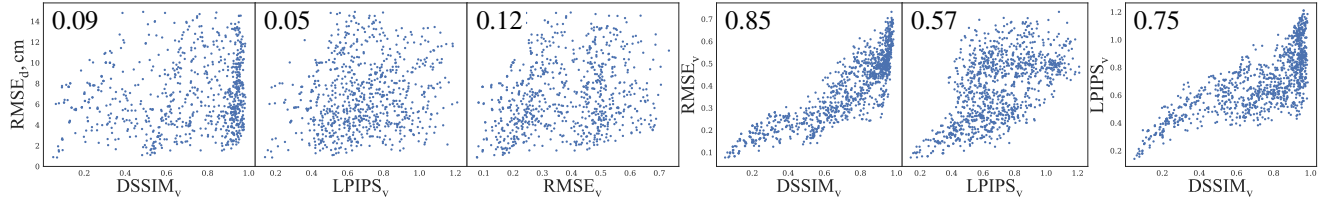


Figure 4: Scatter plots demonstrating correlation of quality measures, and the corresponding values of the Pearson correlation coefficient in the corner. Each point represents one super-resolution result.

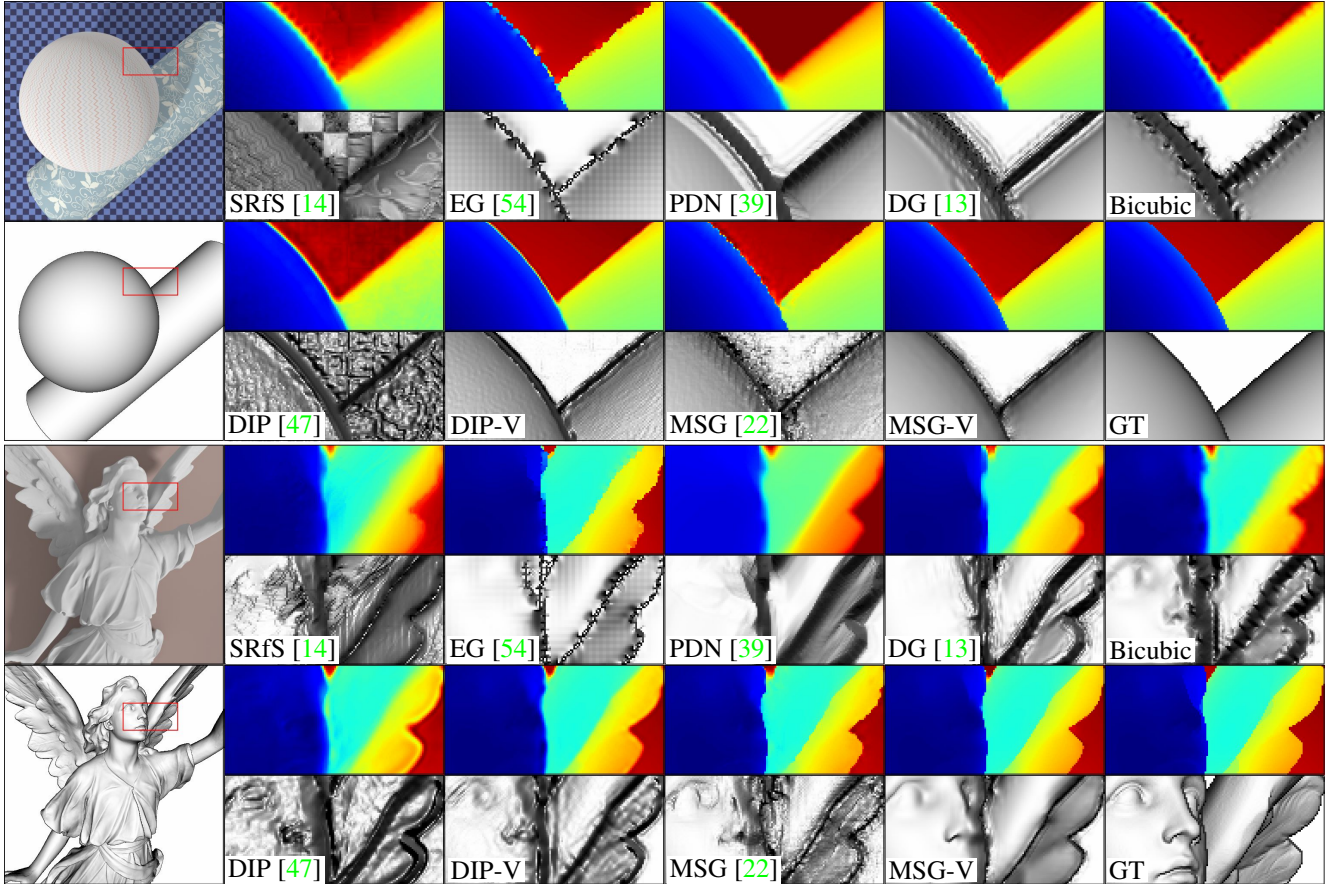


Figure 5: Super-resolution results on “Sphere and cylinder” and “Lucy” from SimGeo with the scaling factor of 4. Depth maps are in pseudo-color and depth map renderings are in grayscale. Best viewed in color.



	Plant						Vintage						Recycle						Umbrella					
	DSSIM <sub>v</sub>		LPIPS <sub>v</sub>		RMSE <sub>v</sub>		DSSIM <sub>v</sub>		LPIPS <sub>v</sub>		RMSE <sub>v</sub>		DSSIM <sub>v</sub>		LPIPS <sub>v</sub>		RMSE <sub>v</sub>		DSSIM <sub>v</sub>		LPIPS <sub>v</sub>		RMSE <sub>v</sub>	
	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8
SRFS [14]	658	692	632	649	280	309	721	749	631	<u>634</u>	346	382	715	772	610	623	376	410	843	853	797	831	397	443
EG [54]	568		677		255																			
PDN [39]	574	612	659	699	269	305	663	714	706	700	319	350	635	<b>701</b>	523	589	364	457	799	<b>828</b>	847	882	367	452
DG [13]	611	622	745	785	268	291	666	669	796	840	290	<u>300</u>	696	<u>719</u>	602	617	<u>328</u>	<u>383</u>	846	878	781	856	399	457
Bicubic	<u>562</u>	<u>610</u>	688	763	249	290	<u>558</u>	<u>649</u>	602	729	<u>258</u>	302	<b>575</b>	721	<u>474</u>	576	329	398	<b>749</b>	<u>837</u>	747	886	<u>323</u>	<u>380</u>
DIP [47]	919	880	764	723	490	437	953	965	910	872	656	687	871	923	576	605	434	500	915	<u>953</u>	737	<u>722</u>	467	528
MSG [22]	571	645	<u>582</u>	<b>495</b>	<u>234</u>	285	708	785	<b>510</b>	<b>610</b>	292	364	741	869	624	661	485	550	834	896	<u>678</u>	787	442	496
DIP-V	694	707	<b>463</b>	<u>555</u>	262	<u>276</u>	804	884	<u>579</u>	674	343	435	<b>575</b>	735	<b>388</b>	<b>485</b>	<b>273</b>	<b>332</b>	796	854	<b>604</b>	<b>598</b>	<b>318</b>	<b>352</b>
MSG-V	<b>524</b>	<b>575</b>	639	720	<b>194</b>	<b>236</b>	<b>536</b>	<b>643</b>	670	702	<b>211</b>	<b>268</b>	<u>603</u>	737	520	<u>564</u>	368	473	<u>778</u>	842	800	890	348	427

Table 2: Quantitative evaluation on ICL-NUIM and Middlebury datasets. All metrics are in thousandths. Lower values correspond to better results. The best result is in bold, the second best is underlined.

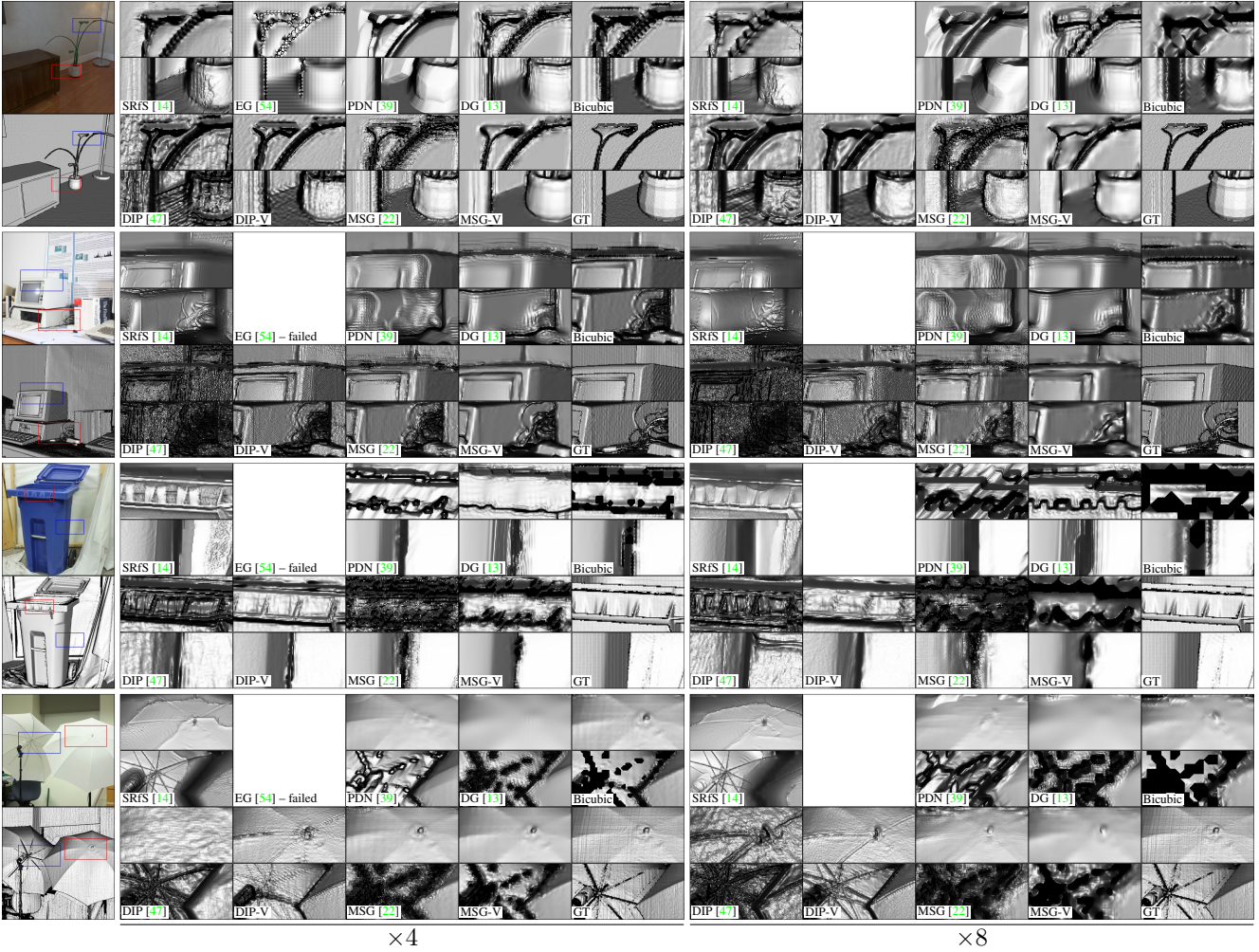


Figure 6: Depth map renderings corresponding to super-resolution results on “Plant” from “ICL-NUIM” and “Vintage”, “Recycle” and “Umbrella” from Middlebury datasets with the scaling factor of 4 on the left and the scaling factor of 8 on the right. Best viewed in large scale.

## Supplementary material

### A. Additional evaluation details

In the literature on range image processing, the term *depth* is used to denote three different types of range data:

- *disparity*, presented in, *e.g.*, Middlebury dataset, *i.e.*, the difference in image location of a feature within two stereo images;
- *orthogonal depth*, presented in, *e.g.*, SUN-RGBD dataset, *i.e.*, the distance from a point in the 3D-space to the image plane;
- *perspective depth*, presented in, *e.g.*, low-resolution scans in ToFMark dataset, *i.e.*, the distance from a point in the 3D-space to the camera.

We use the term *depth map* to denote any data of this kind, however, in our experiments we evaluated each super-resolution method on the range type that it was designed for. For evaluation of the disparity processing methods on the datasets that do not provide disparity maps, we calculated virtual disparity images with the baseline of 20 cm.

Here we describe the quality measures that we considered in addition to the ones discussed in the main text. We recall that we label the metrics that compare the depth values directly with subscript “d”, and the visually-based metrics with subscript “v”.

*BadPix* is the fraction of measurements with absolute deviation larger than a certain threshold  $\tau$

$$\text{BadPix}_d(\tau|d_1, d_2) = \frac{1}{N} |\{ij : |d_{1,ij} - d_{2,ij}| > \tau\}|,$$

or the fraction of measurements with relative deviation larger than a threshold

$$\text{BadPix}_d(\tau\%|d_1, d_2) = \frac{1}{N} |\{ij : \left| \frac{d_{1,ij} - d_{2,ij}}{d_{2,ij}} \right| > \frac{\tau}{100}\}|,$$

where  $d_1$  and  $d_2$  are the compared depth maps,  $ij$  represents individual pixels, and  $N$  is the number of pixels. We considered *BadPix* for depth map comparison with absolute thresholds of 1, 5, and 10 cm and relative thresholds of 1, 5, and 10%. We also considered this metric for comparison of depth map renderings with the absolute thresholds of 1, 5, and 10 each divided by 255 (which correspond to deviations by the respective numbers of shades of gray in 8-bit grayscale images).

*Bumpiness*, introduced in [20] for piece-wise planar regions and generalized in [19] for arbitrary smooth surfaces, is a measure of surface smoothness with respect to a reference. It is calculated as

$$\text{Bumpiness}_d(d_1, d_2) = \frac{1}{N} \sum_{ij} \min(0.05, \|H_{d_1-d_2}(i, j)\|_F) \cdot 100,$$

where  $\|\cdot\|_F$  is Frobenius norm and  $H_f(i, j)$  is the Hessian matrix of the function  $f$ , calculated at point  $(i, j)$ . We used the original implementation of this metric. Since this metric includes some parameter values, presumably, specific for the original evaluation dataset, we converted the depth maps to disparity using the camera intrinsics of this dataset.

We used the implementation of SSIM from *scikit-image* [48] and the original implementation of LPIPS from [58].

In addition to our  $\text{RMSE}_v$  we considered RMS difference of two rendered images without averaging over the basis renderings, *i.e.*, calculated for a single lighting condition. We denote this metric as  $\text{RMSE}_v^1$ : for a light direction  $e$  and a pair of normal maps  $\mathbf{n}_1, \mathbf{n}_2$  it is calculated as

$$\text{RMSE}_v^1(d_1, d_2) = \sqrt{\frac{1}{N} \sum_{i,j} \|e \cdot \mathbf{n}_{1,ij} - e \cdot \mathbf{n}_{2,ij}\|_2^2}.$$

### B. Comparison of quality measures

In Figures 11-16 we compare the relations between different subsets of quality measures. We present pair-wise correlations of the metrics in the form of scatter plots in the lower half of the figure and Pearson and Spearman correlation coefficients in the upper half of the figure. For reference, on the diagonal of the figure we also include kernel density estimates of metric value distributions for each super-resolution method. The distributions for the modified methods DIP-v and MSG-v are represented with the dashed black and solid black curves respectively.

On the depth maps with missing measurements, the methods that do not inpaint the regions with the missing measurements (including MSG-v) sometimes produced severe outliers around these regions. To minimize the influence of such outliers on the results of the metric comparison, we limited the value of  $\text{RMSE}_d$  to a maximum of 0.5 meters. Among the collected super-resolved images, 8% exceeded this threshold.

For each metric, applied to rendered images, we gathered the values of this metric for four different light directions, as described in Section 5.2 of the main text. We then calculated two additional values, the worst and the average of these four. We label the respective versions of the metric with suffixes  $e_1, e_2, e_3, e_4$ , max and avg. For each metric, we found that these six versions are strongly correlated, as illustrated in Figures 11-13, so we further focused on the worst value of each metric.

We also found that different versions of  $\text{RMSE}_v^1$  produce very similar results to our  $\text{RMSE}_v$ , as illustrated in Figure 11. It is consistent with the observation that if  $\text{RMSE}_v$  is bounded by a constant  $C$ , then for *any* choice of the light direction  $e$ ,  $\text{RMSE}_v^1$  is bounded by  $C$ , which can be easily seen from the fact that  $\text{RMSE}_v$  does not depend on the choice of the basis, so we can choose one of the basis light



directions to be equal to  $e$ .

In Figure 14 we compare the metrics of different types: pixel-wise  $\text{RMSE}_d$ ,  $\text{BadPix}_d(5\text{cm})$  and  $\text{BadPix}_d(5\%)$  applied to depth directly; “worst” versions of pixel-wise  $\text{BadPix}_v(5)$  and perceptual  $\text{DSSIM}_v$  and  $\text{LPIPS}_v$ , applied to rendered images; geometrical  $\text{Bumpiness}_d$  and our  $\text{RMSE}_v$ . We found that all three pixel-wise metrics applied to depth directly demonstrate weak correlation with visual and geometrical metrics. Pixel-wise  $\text{BadPix}_v(5)$  applied to rendered images, although strongly correlated with perceptual metrics, is inappropriate for gradient-based optimization. Additional comparison of pixel-wise  $\text{BadPix}_d$  and  $\text{BadPix}_v$  with different thresholds to perceptual  $\text{DSSIM}_v$  and  $\text{LPIPS}_v$  (Figures 15 and 16) leads to the same conclusions.  $\text{Bumpiness}_d$  is also strongly correlated with perceptual metrics but only measures local curvature deviation, while the visual appearance of 3D surface is determined by its local orientation.

### C. Comparison of super-resolution methods

In Tables 4-11 we present the results of quantitative evaluation of super-resolution methods on the datasets SimGeo, ICL-NUIM and Middlebury for Box downsampling model and scaling factors of 4 and 8. In Table 3 we present the average values.  $\text{RMSE}_d$  is in millimeters,  $\text{BadPix}$  is in percents,  $\text{DSSIM}_v$ ,  $\text{LPIPS}_v$  and  $\text{RMSE}_v$  are in thousandths. For all visual metrics except  $\text{RMSE}_v$  the presented value is of the “worst” version. For all metrics the lower value corresponds to the better result. The best results are highlighted in bold and the second best results are underlined.

In addition to metric values, the last three columns of the tables contain the results of the informal perceptual study collected over approximately 250 subjects. In this study, for each scene from SimGeo, ICL-NUIM and Middlebury datasets subjects were shown the renderings of super-resolved depth maps, shuffled randomly, and were asked to choose the renderings, the most and second most similar to ground truth. The renderings calculated with the fourth light direction were used. The values in the columns “User, 1st”, “User, 2nd”, and “Top 2” represent the percentages of the subjects who chose the rendering of the super-resolved depth map, produced by the method in the corresponding method, as the most similar, second most similar, or one of the two most similar to the ground truth respectively. We found that our  $\text{RMSE}_v$  is mostly consistent with human judgements.

### D. Training with $\text{MSE}_v$

Since optimization of  $\text{MSE}_v$  alone is an ill-posed problem, we used a regularization term that penalizes absolute depth deviation. We found that among different regularizers, including  $\text{MSE}_d$ ,  $\text{Lap}_1$  produces the best results. In general, we found that optimization leads to the best results

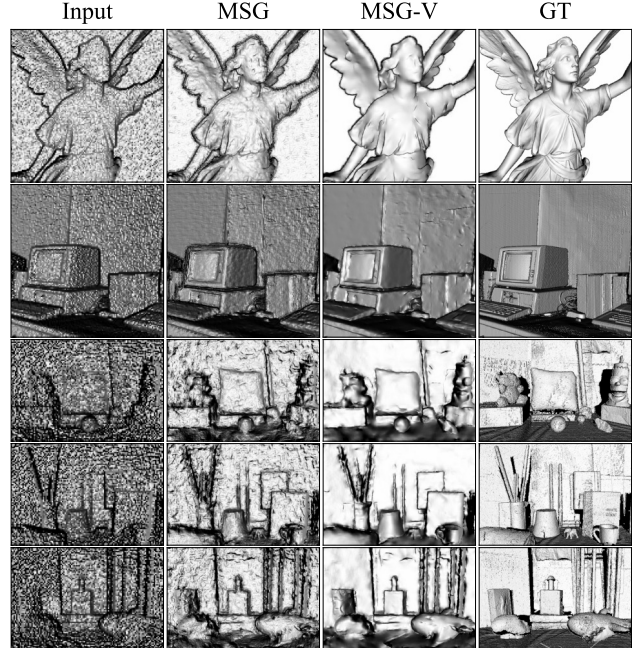


Figure 7:  $\times 4$  super-resolution results produced by the original MSG and MSG-V with our loss, both trained on noisy data. The upper two samples contain synthetic noise, while the lower three from ToFMark dataset represent real noisy ToF measurements. Best viewed in large scale.

if the terms are weighted in such way that geometrically corresponding depth error and angular normal error result in the same magnitudes of terms. The corresponding value of the weighing parameter  $w$  in Equation 4 of the main text is determined by the properties of the training data, such as depth map scaling or field of view of the camera.

### E. Noisy depth measurements in the input

SimGeo, ICL-NUIM and Middlebury datasets were our primary evaluation sets, yielding the most pronounced outcomes, however, these datasets contain only noise-free scenes. As we were interested in evaluation of our approach on a diverse set of RGBD images, we included twelve scenes from SUN RGBD dataset and three scenes from ToFMark dataset that feature real-world noise patterns in our evaluation data. We observed that increased levels of noise are extremely harmful to all non over-smoothing methods, including those modified with our loss, as they fail to produce reasonable super-resolution results, as illustrated in Figures 9-10. To demonstrate that this is not a limitation of our approach, in Figure 7 we present the super-resolution results produced by modified and unmodified versions of MSG, trained on the data with synthetic multiplicative gaussian noise.



Average performance on SimGeo dataset																	
	RMSE <sub>d</sub>		BadPix <sub>d</sub> (5cm)		BadPix <sub>v</sub> (5)		DSSIM <sub>v</sub>		LPIPS <sub>v</sub>		Bumpiness <sub>d</sub>		RMSE <sub>v</sub>		User, 1st	User, 2nd	Top 2
	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x4	x4
Bicubic	55	79	4.1	7.9	<u>23.2</u>	<u>38.3</u>	197	301	320	427	0.70	0.98	193	234	0.5	7.8	8.3
SRfS [14]	61	88	7.5	14.3	74.2	77.1	711	729	869	865	1.48	1.69	311	328	0.0	0.0	0.0
EG [54]	53		2.2		33.1		<u>168</u>		306		<u>0.54</u>		<u>136</u>		0.2	3.9	4.2
PDN [39]	162	211	99.4	99.1	39.2	45.1	224	<u>264</u>	<u>278</u>	<u>407</u>	0.63	<u>0.79</u>	165	201	1.6	<u>12.3</u>	13.8
DG [13]	54	84	3.0	6.4	35.2	39.1	293	316	420	437	0.69	0.82	171	190	0.2	2.9	3.2
DIP [47]	52	59	8.5	12.5	90.5	92.0	887	880	893	915	2.21	2.77	395	475	0.6	0.9	1.5
MSG [22]	39	<u>39</u>	<u>1.5</u>	3.3	51.9	69.3	374	544	569	713	0.79	0.97	194	242	0.4	3.7	4.0
DIP-v	<b>33</b>	41	1.7	<u>2.3</u>	49.7	67.1	313	491	524	598	0.60	0.88	147	<u>174</u>	<u>8.3</u>	<b>59.4</b>	<u>67.8</u>
MSG-v	96	<b>29</b>	<b>0.7</b>	<b>1.5</b>	<b>14.2</b>	<b>34.6</b>	<b>95</b>	<b>206</b>	<b>194</b>	<b>367</b>	<b>0.34</b>	<b>0.46</b>	<b>99</b>	<b>129</b>	<b>88.1</b>	9.1	<b>97.2</b>
Average performance on ICL-NUIM dataset																	
Bicubic	34	54	2.8	5.5	<u>59.3</u>	<u>64.2</u>	<u>431</u>	490	558	668	1.15	1.32	210	252	5.0	<u>28.3</u>	33.3
SRfS [14]	42	62	5.5	11.0	73.5	76.1	641	664	636	660	1.72	1.83	287	314	0.0	0.0	0.0
PDN [39]	135	165	93.8	82.9	66.2	70.2	480	509	623	650	<u>1.14</u>	<u>1.24</u>	237	264	2.6	10.5	13.1
DG [13]	36	58	4.3	6.4	64.4	65.5	497	505	663	689	1.28	1.32	234	259	0.6	5.5	6.1
DIP [47]	43	56	10.6	14.2	83.6	83.4	812	806	690	690	2.73	2.58	394	389	1.1	0.9	2.0
MSG [22]	<u>25</u>	<b>36</b>	<u>1.6</u>	<u>3.5</u>	64.1	69.0	489	557	<u>510</u>	<u>534</u>	1.27	1.46	210	255	1.1	7.2	8.3
DIP-v	28	<u>40</u>	2.6	3.9	67.8	69.6	516	548	<b>407</b>	<b>503</b>	1.45	1.56	<u>209</u>	<u>236</u>	<u>9.6</u>	<b>31.9</b>	<u>41.4</u>
MSG-v	<b>24</b>	41	<b>1.3</b>	<b>3.1</b>	<b>56.3</b>	<b>61.1</b>	<b>387</b>	<b>437</b>	527	602	<b>0.94</b>	<b>1.06</b>	<b>157</b>	<b>192</b>	<b>79.9</b>	11.8	<b>91.7</b>
Average performance on Middlebury dataset																	
Bicubic	843	1139	10.8	13.9	<b>71.5</b>	<b>76.7</b>	<b>648</b>	<b>748</b>	<u>575</u>	720	<b>0.87</b>	<b>0.76</b>	<b>344</b>	<b>386</b>	4.1	<u>25.3</u>	29.4
SRfS [14]	100	145	21.4	33.6	86.4	89.5	780	810	669	704	1.32	<u>1.28</u>	428	461	0.0	0.0	0.0
PDN [39]	173	225	85.3	76.5	83.4	86.4	744	790	653	711	1.38	1.67	405	467	9.9	<b>28.1</b>	<u>37.9</u>
DG [13]	266	330	15.0	24.5	81.8	84.1	765	784	728	740	1.54	1.73	421	442	0.7	10.6	11.3
DIP [47]	<u>72</u>	<u>104</u>	19.6	24.4	92.4	93.4	927	947	737	717	2.82	2.90	565	592	1.2	5.6	6.8
MSG [22]	228	426	10.8	13.1	81.8	87.2	774	858	649	696	1.96	2.19	477	525	0.2	1.6	1.8
DIP-v	<b>56</b>	<b>87</b>	<b>6.4</b>	<u>10.6</u>	83.1	87.4	728	821	<b>506</b>	<b>568</b>	1.34	1.56	<u>353</u>	<u>409</u>	<b>72.3</b>	18.2	<b>90.5</b>
MSG-v	96	133	<u>7.3</u>	<b>9.2</b>	<u>73.3</u>	<u>79.0</u>	<u>667</u>	<u>757</u>	639	<u>690</u>	<u>1.20</u>	1.35	376	431	<u>10.8</u>	9.9	20.7
Average performance on the scenes without missing measurements (SimGeo, ICL-NUIM, Vintage)																	
	RMSE <sub>d</sub>		BadPix <sub>d</sub> (5cm)		BadPix <sub>v</sub> (5)		DSSIM <sub>v</sub>		LPIPS <sub>v</sub>		Bumpiness <sub>d</sub>		RMSE <sub>v</sub>		User, 1st	User, 2nd	Top 2
	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x4	x4
Bicubic	46	69	3.5	6.9	<u>43.7</u>	<u>53.3</u>	<u>333</u>	415	<u>452</u>	561	0.97	1.19	206	248	3.0	<u>18.9</u>	21.9
SRfS [14]	55	81	7.3	14.1	74.6	77.4	680	701	743	753	1.60	1.75	303	326	0.0	0.0	0.0
PDN [39]	148	187	94.4	90.1	55.0	59.8	378	<u>412</u>	482	<u>542</u>	<u>0.93</u>	<u>1.06</u>	210	241	1.9	10.5	12.4
DG [13]	47	73	3.9	6.7	52.1	54.4	416	430	561	584	1.02	1.10	209	230	0.4	4.0	4.4
DIP [47]	50	62	10.7	15.9	87.6	88.2	857	853	801	808	2.59	2.79	414	452	0.8	0.8	1.7
MSG [22]	<u>33</u>	<u>39</u>	<u>1.7</u>	3.6	59.8	70.3	454	569	547	622	1.08	1.26	210	258	0.7	5.8	6.4
DIP-v	<b>31</b>	43	2.2	<u>3.3</u>	60.8	69.9	444	548	474	560	1.09	1.32	<u>191</u>	<u>223</u>	<u>10.2</u>	<b>45.5</b>	<u>55.8</u>
MSG-v	38	<b>38</b>	<b>1.1</b>	<b>2.6</b>	<b>38.0</b>	<b>50.1</b>	<b>264</b>	<b>346</b>	<b>385</b>	<b>501</b>	<b>0.69</b>	<b>0.81</b>	<b>135</b>	<b>169</b>	<b>82.7</b>	10.9	<b>93.6</b>
Average performance on the scenes with missing measurements (Middlebury excluding Vintage)																	
Bicubic	972	1313	11.8	14.7	<b>71.3</b>	<b>76.6</b>	<b>663</b>	<b>765</b>	<u>570</u>	718	<b>0.77</b>	<b>0.61</b>	<u>358</u>	<b>400</b>	3.8	<u>24.7</u>	28.5
SRfS [14]	100	145	22.2	33.8	86.9	89.8	790	820	676	716	1.26	<u>1.21</u>	441	474	0.0	0.0	0.0
PDN [39]	178	234	88.3	76.1	83.5	86.6	757	803	644	713	1.36	1.69	419	487	<u>11.5</u>	<b>32.8</b>	<u>44.3</u>
DG [13]	298	367	16.3	26.9	82.2	84.7	781	803	716	724	1.55	1.76	442	465	0.9	12.2	13.1
DIP [47]	<u>72</u>	<u>102</u>	18.8	20.7	92.2	93.3	923	943	708	691	2.62	2.68	549	577	1.2	6.4	7.7
MSG [22]	259	488	12.1	14.1	82.0	87.7	785	870	673	711	2.02	2.24	507	552	0.2	0.2	0.5
DIP-v	<b>58</b>	<b>91</b>	<b>7.1</b>	<u>11.4</u>	82.7	87.2	716	811	<b>494</b>	<b>550</b>	1.23	1.41	<b>354</b>	<u>405</u>	<b>80.1</b>	13.8	<b>93.9</b>
MSG-v	107	145	<u>8.1</u>	<b>9.8</b>	<u>73.6</u>	<u>79.2</u>	<u>688</u>	<u>776</u>	634	<u>688</u>	<u>1.19</u>	1.34	404	458	1.3	8.8	10.1
Average performance on SimGeo, ICL-NUIM, Middlebury																	
Bicubic	339	462	6.1	9.3	<u>52.4</u>	<u>60.6</u>	<u>437</u>	<u>525</u>	489	611	<u>0.91</u>	<u>1.00</u>	254	296	3.3	<u>20.7</u>	24.0
SRfS [14]	69	101	12.0	20.3	78.5	81.3	715	738	722	741	1.50	1.58	347	372	0.0	0.0	0.0
PDN [39]	157	202	92.4	85.7	64.0	68.2	498	536	533	596	1.07	1.26	276	319	5.0	17.5	22.5
DG [13]	126	166	7.8	13.1	61.6	64.0	531	548	610	628	1.19	1.31	283	305	0.5	6.6	7.1
DIP [47]	<u>57</u>	75	13.3	17.4	89.1	89.8	878	881	771	771	2.60	2.76	457	491	1.0	2.6	3.6
MSG [22]	104	181	5.0	6.9	66.8	75.8	559	664	587	650	1.37	1.57	304	350	0.5	4.0	4.6
DIP-v	<b>40</b>	<b>58</b>	3.8	<u>5.9</u>	67.7	75.4	530	631	481	<b>557</b>	1.14	1.34	<u>242</u>	<u>280</u>	32.3	<b>35.5</b>	<b>67.8</b>
MSG-v	60	<u>72</u>	<b>3.3</b>	<b>4.8</b>	<b>49.2</b>	<b>59.3</b>	<b>398</b>	<b>482</b>	<b>464</b>	<u>560</u>	<b>0.85</b>	<b>0.98</b>	<b>220</b>	<b>260</b>	<b>57.0</b>	10.2	<u>67.3</u>

Table 3: Quantitative evaluation summary. The best result is in bold, the second best in underlined.

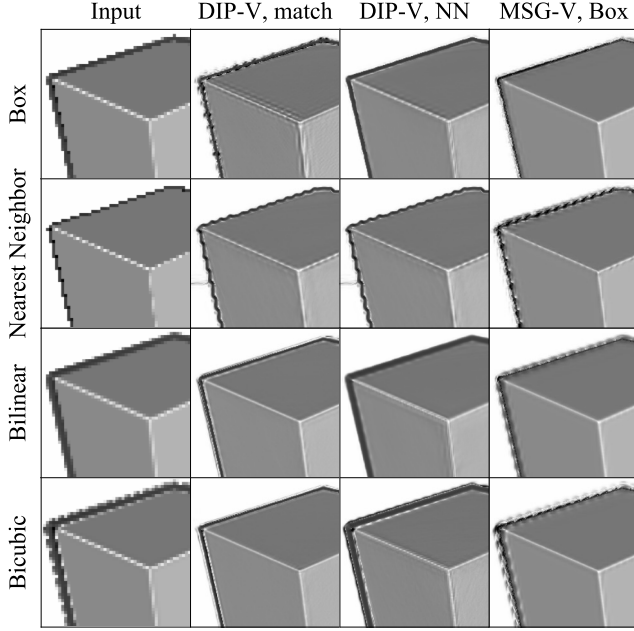


Figure 8:  $\times 4$  super-resolution results for different input downsampling models produced by DIP-V with a matching downsampling model, DIP-V with Nearest Neighbor downsampling model and MSG-V with Box downsampling model. Best viewed in large scale.

## F. Different downsampling models

In Figure 8 we present the results for different downsampling models, used for calculation of low-resolution input. We found that the visual quality remains high when the downsampling model used during training and that of the input match; if this is not the case, the quality deteriorates, as expected.

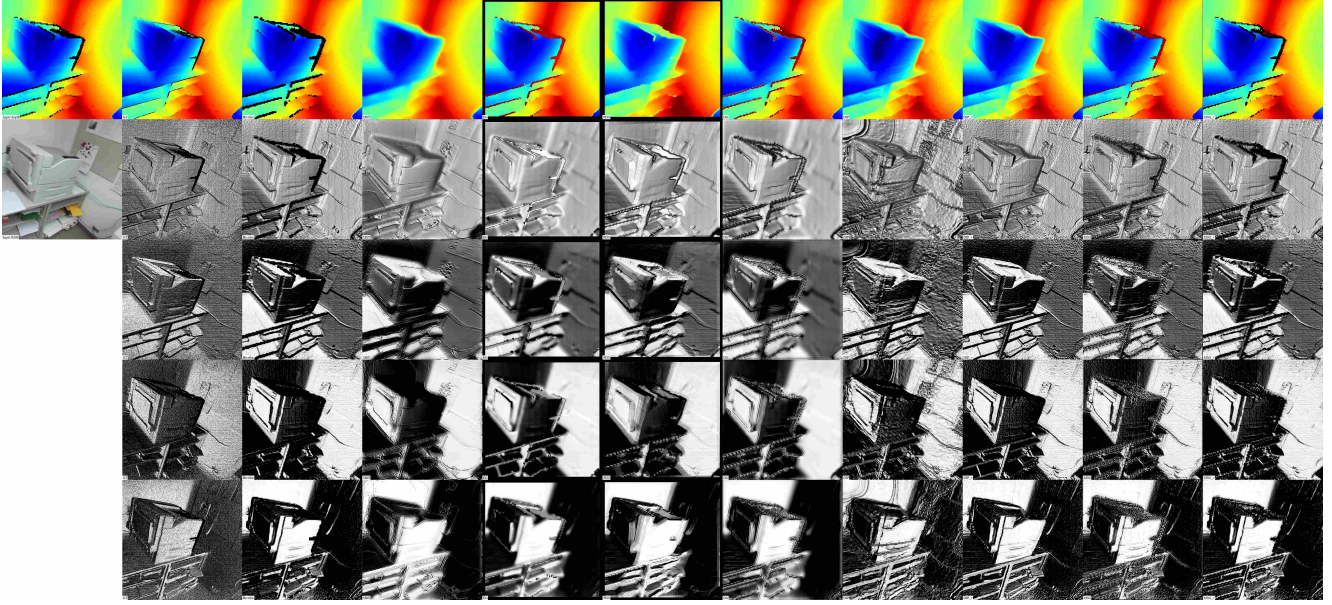


Figure 9: x4 super-resolution results on a Kinect v2 RGBD scan from SUN RGBD dataset. Each visualization is labeled in the bottom left corner. Ground truth is in the 2nd column, DIP-v is in the third from the right, MSG-v in the last one.

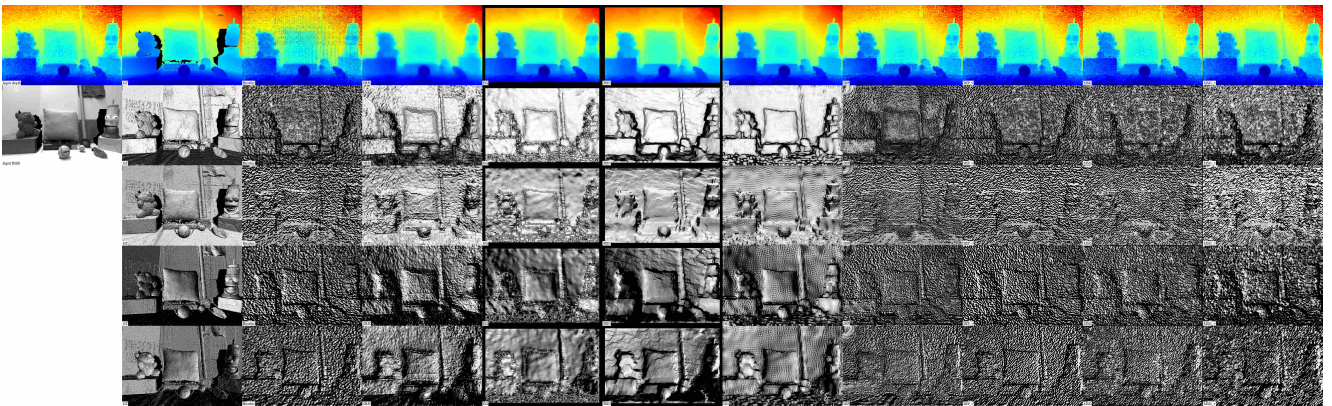


Figure 10: x4 super-resolution results on "Devil" from ToFMark dataset. Each visualization is labeled in the bottom left corner. Ground truth is in the 2nd column, DIP-v is in the third from the right, MSG-v in the last one.



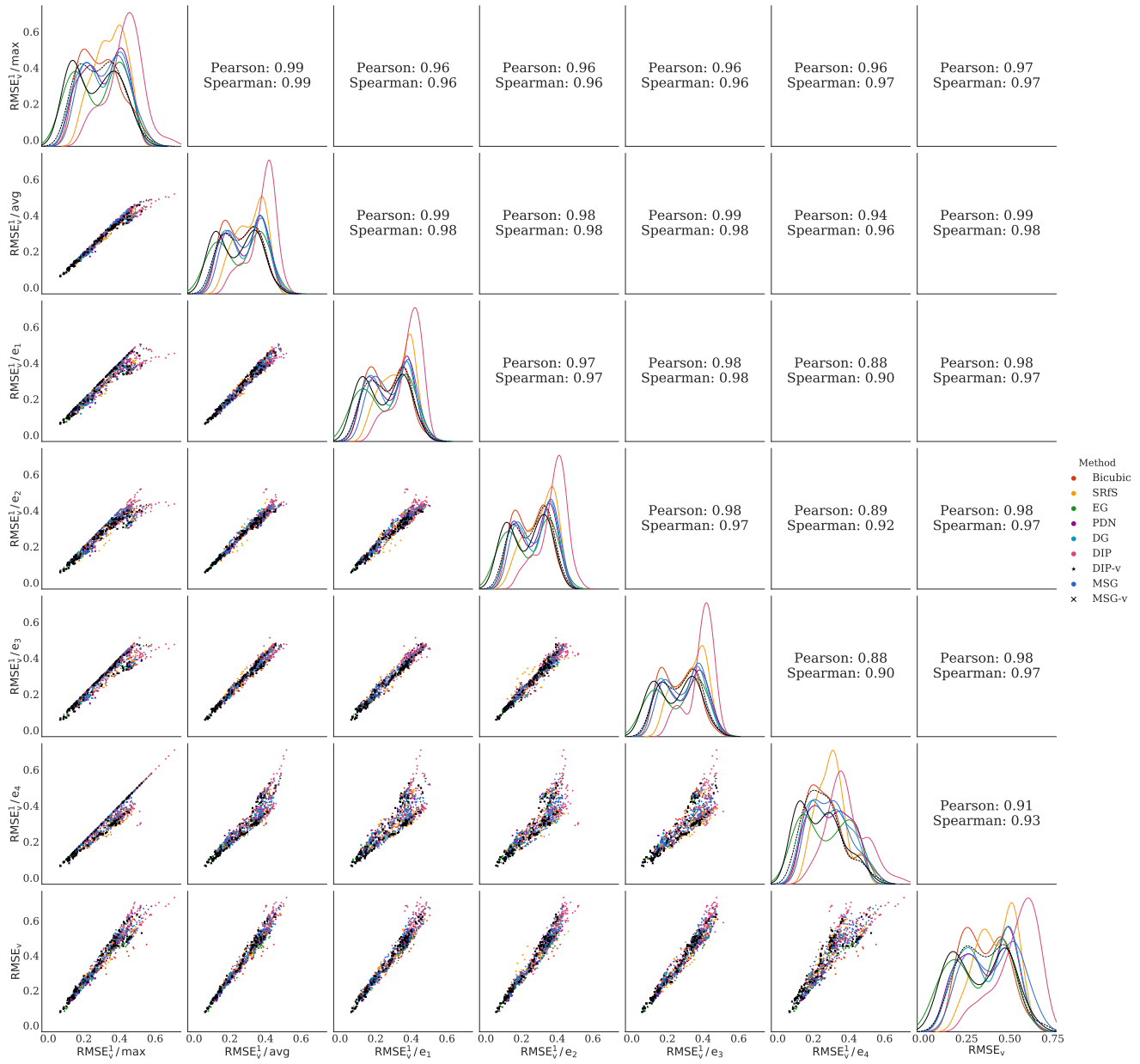


Figure 11: Comparison of different versions of  $RMSE_V^1$  metric and  $RMSE_V$  metric. Best viewed in large scale and in color.

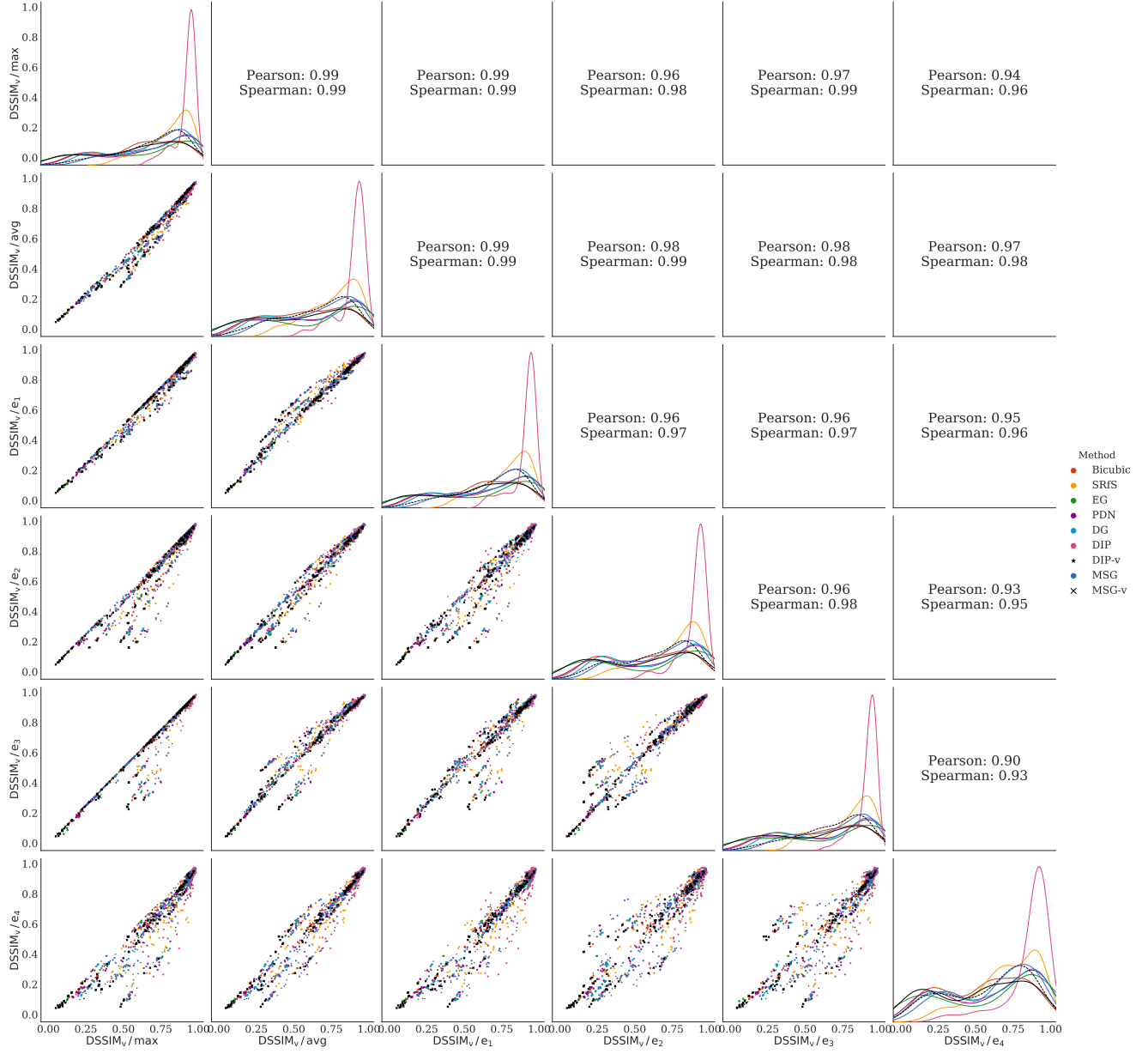


Figure 12: Comparison of different versions of  $DSSIM_v$  metric. Best viewed in large scale and in color.

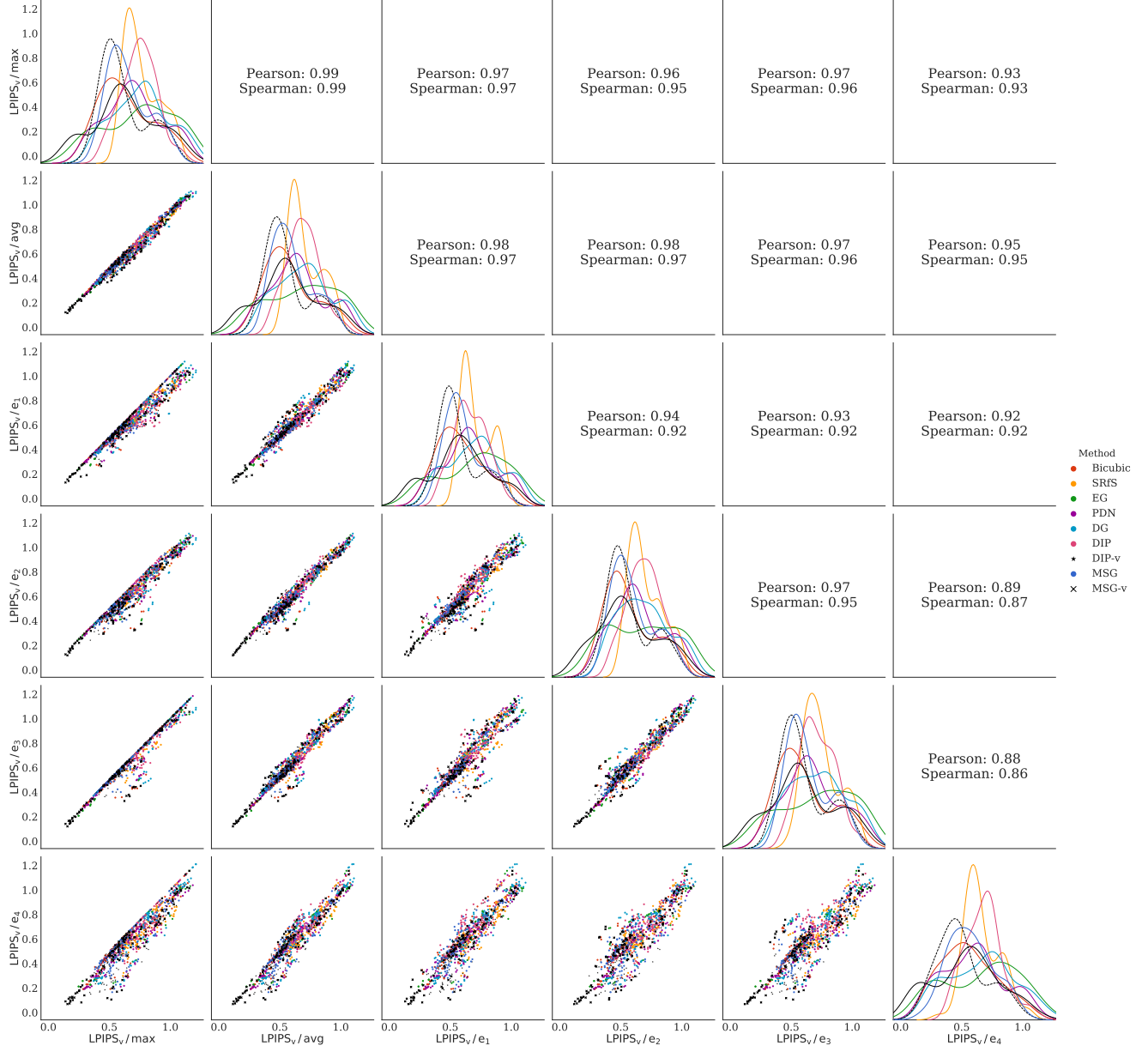


Figure 13: Comparison of different versions of LPIPS<sub>v</sub> metric. Best viewed in large scale and in color.



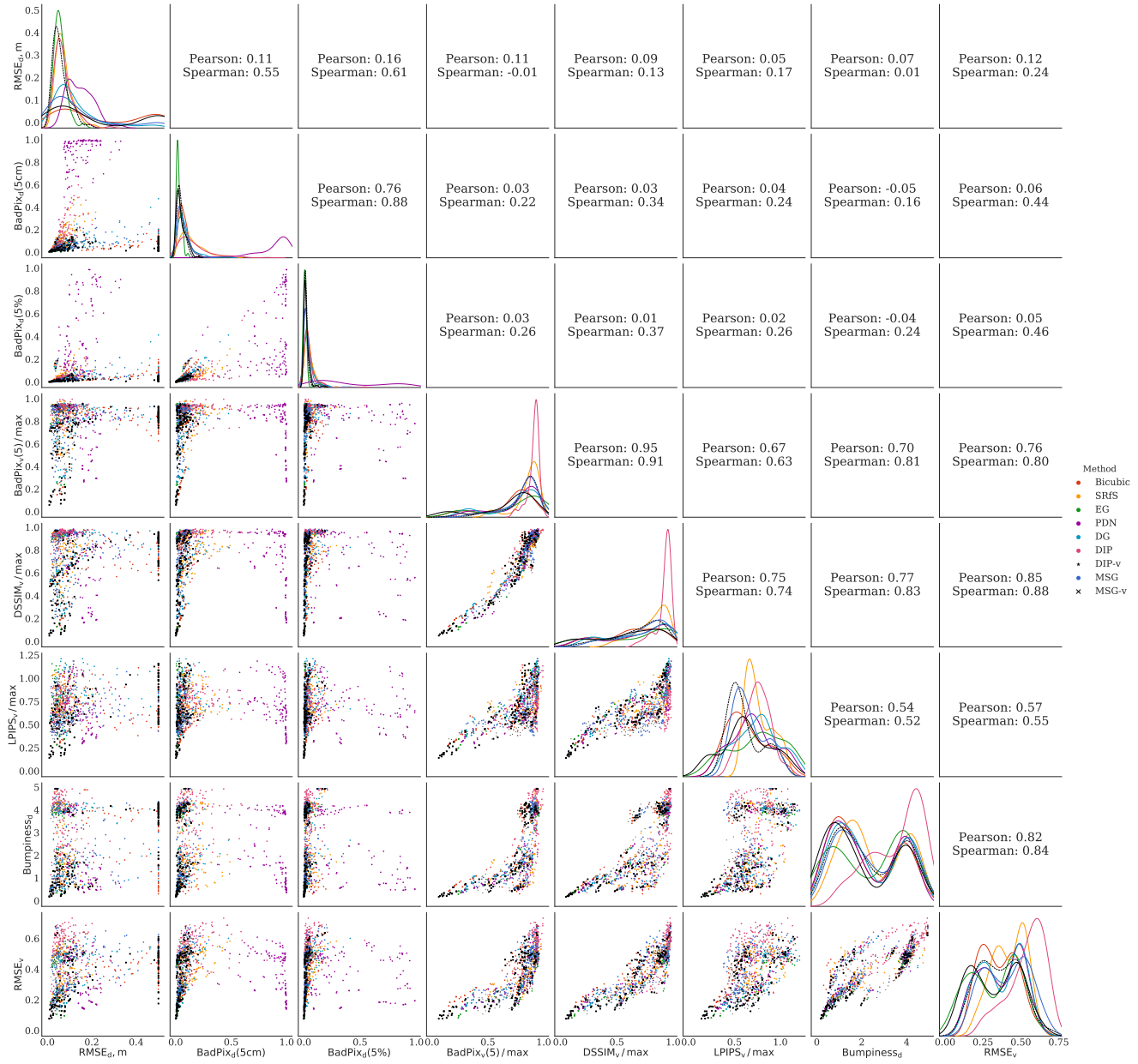
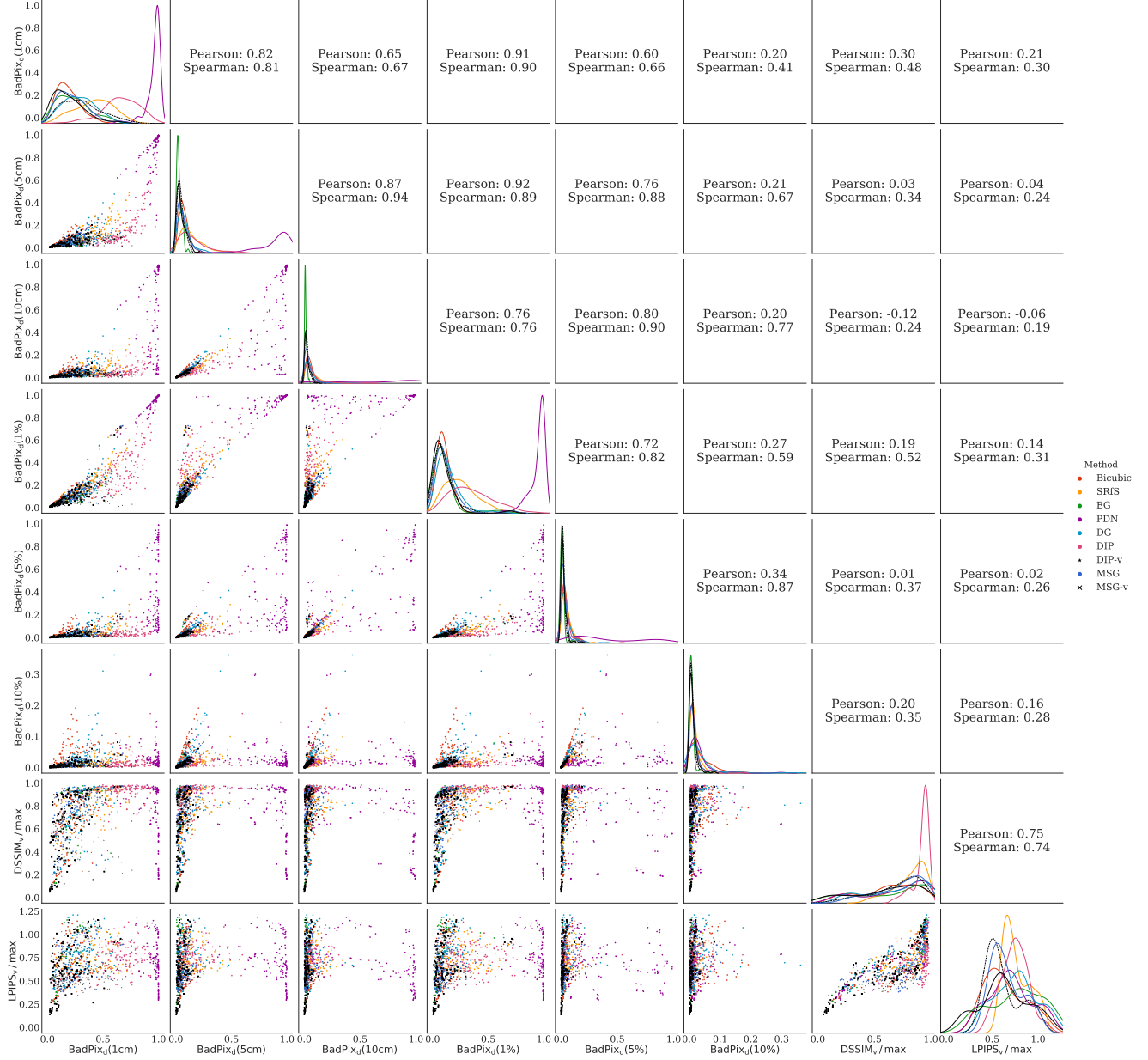


Figure 14: Comparison of metrics of different types. Best viewed in large scale and in color.



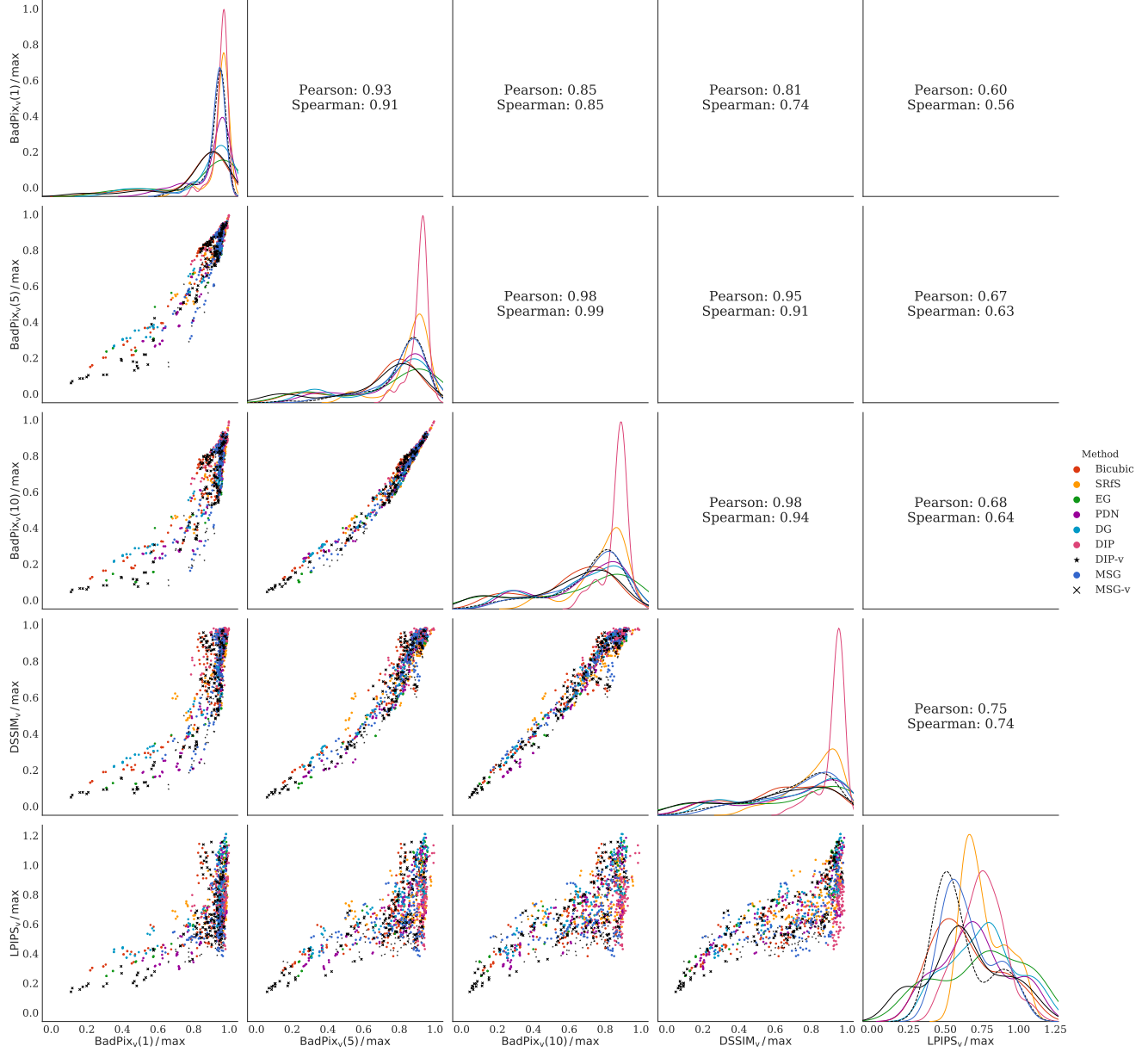


Figure 16: Comparison of different pixel-wise metrics applied to rendered images and perceptual metrics. Best viewed in large scale and in color.



Cube, high-frequency texture																	
	RMSE <sub>d</sub>		BadPix <sub>d</sub> (5cm)		BadPix <sub>v</sub> (5)		DSSIM <sub>v</sub>		LPIPS <sub>v</sub>		Bumpiness <sub>d</sub>		RMSE <sub>v</sub>		User, 1st	User, 2nd	Top 2
	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x4	x4
Bicubic	44	63	2.7	5.2	<u>15.0</u>	<b>27.3</b>	131	<u>204</u>	287	395	0.43	0.61	160	188	0.7	<u>13.2</u>	14.0
SRfS [14]	52	75	6.2	12.1	89.2	80.3	934	818	1036	938	1.73	1.67	361	339	0.0	0.0	0.0
EG [54]	43		1.2		25.4		<u>113</u>		<u>214</u>		<u>0.35</u>		<u>105</u>		0.7	9.6	10.3
PDN [39]	164	219	99.6	99.4	27.5	<u>29.5</u>	156	<b>186</b>	250	<b>368</b>	0.39	<u>0.49</u>	145	171	0.7	4.4	5.1
DG [13]	44	67	1.9	4.2	26.4	30.1	218	240	411	437	0.44	0.55	139	<u>159</u>	0.7	7.4	8.1
DIP [47]	45	48	6.4	8.5	93.5	92.5	963	947	906	918	2.98	2.50	530	494	0.0	0.7	0.7
MSG [22]	<u>29</u>	38	1.0	2.6	60.1	77.9	445	653	687	877	0.79	0.98	176	233	0.0	0.0	0.0
DIP-v	<b>26</b>	<u>36</u>	<u>0.8</u>	<u>1.6</u>	56.2	60.8	352	413	613	653	0.64	0.89	146	162	<u>5.9</u>	<b>58.1</b>	<u>64.0</u>
MSG-v	102	<b>20</b>	<b>0.3</b>	<b>0.7</b>	<b>9.3</b>	51.0	<b>70</b>	316	<b>179</b>	676	<b>0.20</b>	<b>0.39</b>	<b>77</b>	<b>125</b>	<b>91.2</b>	6.6	<b>97.8</b>

Cube, no texture																	
Bicubic	44	63	2.7	5.2	<u>15.0</u>	<u>27.3</u>	131	204	287	395	0.43	0.61	160	188	0.0	5.9	5.9
SRfS [14]	43	63	2.1	4.5	53.4	51.7	516	476	754	728	0.67	0.89	219	228	0.0	0.0	0.0
EG [54]	43		1.2		25.4		128		<u>282</u>		0.35		<u>105</u>		0.0	2.9	2.9
PDN [39]	164	219	99.6	99.4	26.3	29.3	162	<u>185</u>	314	<u>353</u>	0.38	0.49	145	171	0.7	4.4	5.1
DG [13]	44	67	1.9	4.2	26.4	30.1	218	240	411	437	0.44	0.55	139	159	0.0	2.2	2.2
DIP [47]	72	56	23.2	17.3	94.3	99.1	912	980	1026	1133	2.05	4.22	434	683	0.0	0.7	0.7
MSG [22]	29	<u>26</u>	1.0	1.7	30.7	49.8	199	314	509	642	0.42	0.47	157	171	0.0	<u>6.6</u>	6.6
DIP-v	<u>26</u>	35	<u>0.8</u>	<u>1.4</u>	15.1	45.4	<u>95</u>	237	347	478	<u>0.28</u>	<u>0.35</u>	111	<u>107</u>	<u>1.5</u>	<b>76.5</b>	<u>77.9</u>
MSG-v	<b>9</b>	<b>19</b>	<b>0.3</b>	<b>0.4</b>	<b>6.0</b>	<b>13.0</b>	<b>50</b>	<b>73</b>	<b>141</b>	<b>213</b>	<b>0.17</b>	<b>0.21</b>	<b>77</b>	<b>82</b>	<b>97.8</b>	0.7	<b>98.5</b>

Table 4: Quantitative evaluation on “Cube” with different RGBs from SimGeo dataset. The best result is in bold, the second best is underlined.

Sphere and cylinder, high-frequency texture																	
	RMSE <sub>d</sub>		BadPix <sub>d</sub> (5cm)		BadPix <sub>v</sub> (5)		DSSIM <sub>v</sub>		LPIPS <sub>v</sub>		Bumpiness <sub>d</sub>		RMSE <sub>v</sub>		User, 1st	User, 2nd	Top 2
	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x4	x4
Bicubic	57	82	4.1	8.1	<u>20.1</u>	<u>36.7</u>	189	294	313	<u>420</u>	0.67	0.98	189	234	0.0	0.7	0.7
SRfS [14]	70	102	12.1	24.6	91.9	91.8	887	865	1025	1008	2.43	2.41	417	403	0.0	0.0	0.0
EG [54]	55		2.4		30.4		<u>143</u>		326		<u>0.50</u>		<u>130</u>		0.0	1.5	1.5
PDN [39]	157	197	99.3	98.9	40.7	54.1	198	<b>242</b>	<u>295</u>	461	0.60	<u>0.77</u>	150	187	1.5	9.6	11.0
DG [13]	56	87	3.2	6.3	30.9	<b>35.2</b>	265	<u>285</u>	372	<b>386</b>	0.66	<u>0.77</u>	166	<u>180</u>	0.0	1.5	1.5
DIP [47]	46	69	3.9	27.2	97.0	99.2	965	975	1062	1014	4.01	4.80	548	696	1.5	2.9	4.4
MSG [22]	<u>41</u>	<u>41</u>	1.4	3.6	72.6	85.6	626	820	859	960	0.98	1.43	229	314	0.0	0.0	0.0
DIP-v	<b>28</b>	43	<u>1.2</u>	<u>2.4</u>	69.2	86.0	560	850	766	832	0.56	1.45	142	242	<u>32.4</u>	<b>52.9</b>	<u>85.3</u>
MSG-v	99	<b>37</b>	<b>0.6</b>	<b>2.0</b>	<b>14.3</b>	53.0	<b>94</b>	334	<b>267</b>	583	<b>0.29</b>	<b>0.55</b>	<b>96</b>	<b>164</b>	<b>64.7</b>	<u>30.9</u>	<b>95.6</b>
Sphere and cylinder, no texture																	
Bicubic	57	82	4.1	8.1	<u>20.2</u>	36.8	189	294	325	437	0.67	0.98	190	233	0.0	0.7	0.7
SRfS [14]	59	85	4.6	8.6	51.4	70.8	430	619	657	766	0.77	1.25	193	256	0.0	0.0	0.0
EG [54]	56		2.4		30.9		<u>160</u>		383		0.50		<u>128</u>		0.0	1.5	1.5
PDN [39]	157	197	99.3	98.9	38.0	44.1	202	<u>218</u>	<u>294</u>	<u>386</u>	0.58	0.76	150	186	5.9	<u>17.6</u>	23.5
DG [13]	57	87	3.2	6.4	31.0	<u>35.3</u>	265	284	396	409	0.66	0.78	165	180	0.7	2.2	2.9
DIP [47]	49	56	5.0	5.5	85.6	81.6	856	662	927	723	1.01	0.96	244	249	1.5	0.0	1.5
MSG [22]	40	<u>37</u>	<u>1.4</u>	3.1	45.6	64.5	288	444	509	610	0.65	0.76	183	218	0.7	0.7	1.5
DIP-v	<u>35</u>	39	<u>1.4</u>	<u>1.8</u>	41.0	72.6	210	523	517	643	<u>0.47</u>	<u>0.70</u>	130	<u>141</u>	<u>9.6</u>	<b>64.7</b>	<u>74.3</u>
MSG-v	<b>14</b>	<b>27</b>	<b>0.7</b>	<b>1.3</b>	<b>8.5</b>	<b>18.0</b>	<b>77</b>	<b>93</b>	<b>174</b>	<b>200</b>	<b>0.27</b>	<b>0.32</b>	<b>96</b>	<b>110</b>	<b>81.6</b>	12.5	<b>94.1</b>
Sphere and cylinder, low-frequency texture																	
Bicubic	57	82	4.1	8.1	<u>20.1</u>	36.7	189	294	313	420	0.67	0.98	189	234	0.0	2.2	2.2
SRfS [14]	62	91	6.8	14.9	74.9	81.0	691	738	961	956	1.38	1.65	311	335	0.0	0.0	0.0
EG [54]	54		2.4		30.4		<u>160</u>		377		<u>0.50</u>		129		<u>0.7</u>	7.4	8.1
PDN [39]	157	197	99.3	98.9	37.9	44.5	202	<u>219</u>	<u>299</u>	397	0.58	0.76	150	186	<u>0.7</u>	<u>36.0</u>	36.8
DG [13]	56	87	3.2	6.3	30.9	<u>35.2</u>	265	285	372	<u>386</u>	0.66	0.77	166	180	0.0	3.7	3.7
DIP [47]	49	52	8.0	4.9	85.5	84.7	796	812	821	924	1.19	1.18	267	250	0.0	0.0	0.0
MSG [22]	41	<u>41</u>	<u>1.3</u>	3.0	39.6	66.2	264	458	493	612	0.64	0.74	181	213	0.0	1.5	1.5
DIP-v	<u>38</u>	42	1.7	<u>2.2</u>	48.0	60.4	238	351	456	516	<u>0.50</u>	<u>0.61</u>	<u>128</u>	<u>152</u>	<u>0.7</u>	<b>47.8</b>	<u>48.5</u>
MSG-v	<b>16</b>	<b>26</b>	<b>0.7</b>	<b>1.2</b>	<b>8.5</b>	<b>17.5</b>	<b>76</b>	<b>92</b>	<b>156</b>	<b>181</b>	<b>0.27</b>	<b>0.31</b>	<b>97</b>	<b>100</b>	<b>97.8</b>	1.5	<b>99.3</b>

Table 5: Quantitative evaluation on “Sphere and cylinder” with different RGBs from SimGeo dataset. The best result is in bold, the second best is underlined.

Lucy																	
	RMSE <sub>d</sub>		BadPix <sub>d</sub> (5cm)		BadPix <sub>v</sub> (5)		DSSIM <sub>v</sub>		LPIPS <sub>v</sub>		Bumpiness <sub>d</sub>		RMSE <sub>v</sub>		User, 1st	User, 2nd	Top 2
	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x4	x4
Bicubic	72	103	6.8	13.0	48.8	<u>65.0</u>	<u>355</u>	<u>519</u>	398	497	1.37	1.74	267	328	<u>2.2</u>	<u>24.3</u>	<u>26.5</u>
SRfS [14]	82	113	13.2	20.8	84.6	87.1	811	857	781	792	1.90	2.28	367	407	0.0	0.0	0.0
EG [54]	69		3.5		56.2		357		426		<u>1.05</u>		<u>220</u>		0.0	0.7	0.7
PDN [39]	173	234	99.0	98.8	64.9	68.9	456	535	368	480	1.24	1.47	251	303	0.0	1.5	1.5
DG [13]	69	108	4.9	11.0	65.5	68.6	523	562	558	565	1.28	1.50	249	281	0.0	0.7	0.7
DIP [47]	<u>53</u>	75	4.7	11.4	87.4	95.2	827	908	615	778	2.02	2.93	344	478	0.7	0.7	1.5
MSG [22]	54	<u>53</u>	<u>2.7</u>	5.4	62.9	71.7	444	577	480	578	1.30	1.42	259	306	1.5	13.2	14.7
DIP-v	<b>44</b>	55	4.6	<u>4.4</u>	69.0	77.5	421	574	446	<u>468</u>	1.15	<u>1.27</u>	223	<u>239</u>	0.0	<b>56.6</b>	<u>56.6</u>
MSG-v	74	<b>47</b>	<b>1.6</b>	<b>3.7</b>	<b>38.8</b>	<b>55.0</b>	<b>205</b>	<b>325</b>	<b>251</b>	<b>348</b>	<b>0.82</b>	<b>0.96</b>	<b>156</b>	<b>195</b>	<b>95.6</b>	2.2	<b>97.8</b>

Table 6: Quantitative evaluation on “Lucy” from SimGeo dataset. The best result is in bold, the second best is underlined.

Painting																	
	RMSE <sub>d</sub>		BadPix <sub>d</sub> (5cm)		BadPix <sub>v</sub> (5)		DSSIM <sub>v</sub>		LPIPS <sub>v</sub>		Bumpiness <sub>d</sub>		RMSE <sub>v</sub>		User, 1st	User, 2nd	Top 2
	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x4	x4
Bicubic	28	47	2.5	5.6	<u>57.1</u>	64.1	<u>423</u>	514	544	649	0.95	1.15	213	265	4.4	<b>47.8</b>	<u>52.2</u>
SRfS [14]	39	60	6.5	15.9	78.4	81.2	707	722	612	661	1.47	1.55	308	337	0.0	0.0	0.0
EG [54]	36		3.1		61.9		481		720		0.94		231		0.0	3.7	3.7
PDN [39]	151	215	99.3	99.2	65.2	70.2	488	532	669	709	<u>0.89</u>	<u>1.01</u>	237	275	4.4	10.3	14.7
DG [13]	31	49	2.4	5.5	61.9	<u>63.9</u>	503	<u>506</u>	678	700	1.08	1.13	232	272	0.7	3.7	4.4
DIP [47]	30	37	4.0	4.7	80.4	79.5	802	766	630	612	2.18	1.82	362	341	0.0	0.0	0.0
MSG [22]	<u>21</u>	<b>29</b>	<u>1.2</u>	<u>2.2</u>	63.7	67.9	495	570	<u>475</u>	<u>507</u>	0.97	1.12	<u>203</u>	243	2.2	5.1	7.4
DIP-v	22	<u>32</u>	2.3	3.2	70.1	70.3	567	564	<b>386</b>	<b>501</b>	1.07	1.12	210	<u>239</u>	2.9	<u>21.3</u>	24.3
MSG-v	<b>17</b>	34	<b>0.9</b>	<b>1.8</b>	<b>51.4</b>	<b>58.0</b>	<b>354</b>	<b>410</b>	532	607	<b>0.67</b>	<b>0.77</b>	<b>142</b>	<b>170</b>	<b>85.3</b>	8.1	<b>93.4</b>

Sofa																	
	RMSE <sub>d</sub>		BadPix <sub>d</sub> (5cm)		BadPix <sub>v</sub> (5)		DSSIM <sub>v</sub>		LPIPS <sub>v</sub>		Bumpiness <sub>d</sub>		RMSE <sub>v</sub>		User, 1st	User, 2nd	Top 2
	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x4	x4
Bicubic	38	58	1.8	3.6	<u>75.4</u>	<u>77.0</u>	<u>566</u>	<u>616</u>	704	764	<u>2.12</u>	<u>2.33</u>	<u>212</u>	<u>250</u>	3.7	15.4	19.1
SRfS [14]	39	58	2.0	3.5	82.3	88.1	715	832	631	743	2.97	3.45	310	405	0.0	0.0	0.0
EG [54]	42		2.5		79.0		598		767		2.28		213		0.0	8.8	8.8
PDN [39]	86	91	71.0	70.8	83.3	83.0	641	658	784	763	2.40	2.50	260	264	0.7	3.7	4.4
DG [13]	41	63	3.2	4.4	77.7	77.9	624	632	823	855	2.30	<u>2.33</u>	255	263	0.0	5.1	5.1
DIP [47]	45	57	7.1	12.7	93.1	94.0	928	946	758	738	3.91	3.99	518	560	0.0	0.0	0.0
MSG [22]	<b>27</b>	<b>36</b>	1.2	2.3	80.6	85.7	718	791	<u>606</u>	<u>610</u>	2.71	3.22	254	316	0.0	0.7	0.7
DIP-v	<u>27</u>	<u>43</u>	<u>0.9</u>	<u>2.0</u>	79.1	82.5	645	718	<b>414</b>	<b>585</b>	2.67	3.07	215	266	<u>19.1</u>	<b>47.8</b>	<u>66.9</u>
MSG-v	<u>35</u>	44	<b>0.7</b>	<b>1.6</b>	<b>74.0</b>	<b>75.7</b>	<b>537</b>	<b>585</b>	710	759	<b>1.96</b>	<b>2.10</b>	<b>165</b>	<b>196</b>	<b>76.5</b>	<u>18.4</u>	<b>94.9</b>

Plant																	
	RMSE <sub>d</sub>		BadPix <sub>d</sub> (5cm)		BadPix <sub>v</sub> (5)		DSSIM <sub>v</sub>		LPIPS <sub>v</sub>		Bumpiness <sub>d</sub>		RMSE <sub>v</sub>		User, 1st	User, 2nd	Top 2
	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x4	x4
Bicubic	38	58	3.7	6.4	<u>75.9</u>	<u>79.9</u>	<u>562</u>	<u>610</u>	688	763	1.58	1.79	249	290	1.5	<u>22.1</u>	23.5
SRfS [14]	46	65	5.8	9.5	82.9	85.0	658	692	632	649	1.96	2.13	280	309	0.0	0.0	0.0
EG [54]	43		4.5		82.2		568		677		1.64		255		0.0	0.7	0.7
PDN [39]	88	89	94.5	37.8	79.5	82.5	574	612	659	699	<u>1.46</u>	<u>1.60</u>	269	305	4.4	7.4	11.8
DG [13]	40	63	3.9	6.7	79.5	81.1	611	622	745	785	1.67	1.70	268	291	2.2	11.0	13.2
DIP [47]	38	47	6.9	6.1	93.9	92.8	919	880	764	723	4.33	3.95	490	437	0.0	0.7	0.7
MSG [22]	<u>31</u>	<u>44</u>	<u>2.3</u>	<b>3.7</b>	78.0	81.8	571	645	<u>582</u>	<b>495</b>	1.62	1.84	<u>234</u>	285	0.0	11.8	11.8
DIP-v	<u>31</u>	<b>40</b>	4.7	4.8	83.5	84.1	694	707	<b>463</b>	<u>555</u>	2.25	2.21	262	<u>276</u>	<u>11.0</u>	<b>33.1</b>	<u>44.1</u>
MSG-v	<b>27</b>	<u>44</u>	<b>1.8</b>	<u>3.9</u>	<b>74.3</b>	<b>77.8</b>	<b>524</b>	<b>575</b>	639	720	<b>1.31</b>	<b>1.47</b>	<b>194</b>	<b>236</b>	<b>80.9</b>	13.2	<b>94.1</b>

Table 7: Quantitative evaluation on RGBD frames from ICL-NUIM “Living Room” sequence. The best result is in bold, the second best is underlined.

Office																	
	RMSE <sub>d</sub>		BadPix <sub>d</sub> (5cm)		BadPix <sub>v</sub> (5)		DSSIM <sub>v</sub>		LPIPS <sub>v</sub>		Bumpiness <sub>d</sub>		RMSE <sub>v</sub>		User, 1st	User, 2nd	Top 2
	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x4	x4
Bicubic	47	80	4.0	7.8	<u>24.4</u>	<u>34.1</u>	<u>216</u>	<u>285</u>	<u>412</u>	594	0.81	0.95	208	254	<u>19.9</u>	<b>44.1</b>	<u>64.0</u>
SRfS [14]	49	89	5.8	14.4	53.4	54.4	595	593	690	636	1.71	1.66	298	302	0.0	0.0	0.0
PDN [39]	185	185	99.3	90.5	36.5	50.2	250	294	457	518	<u>0.76</u>	<u>0.92</u>	234	272	0.7	3.7	4.4
DG [13]	49	85	9.0	11.9	36.5	37.6	319	330	534	571	1.03	1.05	240	266	0.0	0.7	0.7
DIP [47]	76	109	30.3	48.2	72.1	73.9	726	819	690	797	2.45	2.70	372	408	1.5	1.5	2.9
MSG [22]	<u>35</u>	<b>48</b>	<u>2.4</u>	<u>6.8</u>	35.4	44.5	263	360	415	543	0.83	0.95	<u>199</u>	247	2.2	3.7	5.9
DIP-v	40	<u>65</u>	3.8	7.4	45.4	47.9	311	352	414	<u>504</u>	1.08	1.18	205	<u>235</u>	17.6	<u>25.0</u>	42.6
MSG-v	<b>32</b>	<u>65</u>	<b>1.9</b>	<b>5.3</b>	<b>19.3</b>	<b>29.6</b>	<b>157</b>	<b>224</b>	<b>313</b>	<b>432</b>	<b>0.59</b>	<b>0.72</b>	<b>151</b>	<b>198</b>	<b>58.1</b>	21.3	<b>79.4</b>
Coat rack																	
Bicubic	<u>13</u>	20	1.5	3.0	<u>73.1</u>	<u>75.3</u>	<u>507</u>	539	537	651	0.54	0.60	171	196	0.0	<u>19.1</u>	19.1
SRfS [14]	24	28	3.8	5.4	82.3	80.5	672	556	650	612	0.83	0.57	237	203	0.0	0.0	0.0
EG [54]	<u>13</u>		1.2		77.8		541		550		0.55		186		0.7	7.4	8.1
PDN [39]	140	191	99.6	99.9	77.1	78.0	544	557	621	631	<u>0.48</u>	<u>0.50</u>	178	193	<u>5.1</u>	<b>28.7</b>	<u>33.8</u>
DG [13]	<u>13</u>	20	1.4	3.2	74.4	75.6	530	<u>532</u>	593	621	0.54	0.58	166	201	0.0	9.6	9.6
DIP [47]	15	24	1.9	3.5	85.5	85.4	766	701	625	624	1.15	0.97	256	246	4.4	2.2	6.6
MSG [22]	<b>11</b>	<b>17</b>	<u>0.9</u>	<b>1.6</b>	73.3	76.1	522	546	523	<u>554</u>	0.51	0.55	<u>165</u>	189	2.2	16.2	18.4
DIP-v	<u>13</u>	<b>17</b>	1.8	<u>2.0</u>	75.2	75.5	543	542	<b>422</b>	<b>463</b>	0.62	0.60	171	<u>181</u>	0.7	12.5	13.2
MSG-v	<b>11</b>	<u>18</u>	<b>0.8</b>	2.2	<b>71.4</b>	<b>74.2</b>	<b>482</b>	<b>502</b>	<u>516</u>	563	<b>0.42</b>	<b>0.48</b>	<b>136</b>	<b>161</b>	<b>86.8</b>	4.4	<b>91.2</b>
Displays																	
Bicubic	41	63	3.2	6.4	<u>49.9</u>	<u>54.9</u>	<u>315</u>	<u>374</u>	460	585	0.92	1.08	208	256	0.7	<u>21.3</u>	22.1
SRfS [14]	53	75	9.0	17.3	61.9	67.3	500	591	599	659	1.35	1.60	288	328	0.0	0.0	0.0
EG [54]	46		5.9		66.7		388		587		0.94		216		0.0	2.9	2.9
PDN [39]	159	220	99.2	99.0	55.4	57.2	381	403	547	580	<u>0.85</u>	<u>0.95</u>	242	275	0.0	9.6	9.6
DG [13]	43	66	5.8	6.7	56.5	56.7	395	406	606	601	1.06	1.10	243	265	0.7	2.9	3.7
DIP [47]	52	60	13.4	9.7	76.9	74.6	732	724	672	645	2.36	2.06	365	344	0.7	0.7	1.5
MSG [22]	<u>26</u>	<b>42</b>	<u>1.7</u>	4.4	53.9	58.0	367	430	461	<u>493</u>	0.97	1.08	204	251	0.0	5.9	5.9
DIP-v	32	45	2.4	<u>4.0</u>	53.7	57.6	336	407	<b>344</b>	<b>409</b>	1.00	1.18	<u>191</u>	<u>221</u>	<u>5.9</u>	<b>51.5</b>	<u>57.4</u>
MSG-v	<b>23</b>	<u>43</u>	<b>1.4</b>	<b>3.5</b>	<b>47.2</b>	<b>51.0</b>	<b>271</b>	<b>324</b>	<u>451</u>	531	<b>0.69</b>	<b>0.80</b>	<b>152</b>	<b>190</b>	<b>91.9</b>	5.1	<b>97.1</b>

Table 8: Quantitative evaluation on RGBD frames from ICL-NUIM “Office Room” sequence. The best result is in bold, the second best is underlined.

	Vintage																
	RMSE <sub>d</sub>		BadPix <sub>d</sub> (5cm)		BadPix <sub>v</sub> (5)		DSSIM <sub>v</sub>		LPIPS <sub>v</sub>		Bumpiness <sub>d</sub>		RMSE <sub>v</sub>		User, 1st	User, 2nd	Top 2
	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x4	x4
Bicubic	67	98	4.6	9.0	<u>72.8</u>	<b>77.3</b>	<u>558</u>	<u>649</u>	602	729	1.51	1.64	<u>258</u>	302	5.9	<u>28.7</u>	34.6
SRfS [14]	101	145	16.8	32.3	83.7	87.2	721	749	631	<u>634</u>	1.64	1.68	346	382	0.0	0.0	0.0
PDN [39]	140	174	67.6	79.0	82.3	85.7	663	714	706	700	1.51	1.57	319	350	0.0	0.0	0.0
DG [13]	72	107	7.1	10.4	79.4	80.1	666	669	796	840	<u>1.50</u>	<u>1.52</u>	290	<u>300</u>	0.0	0.7	0.7
DIP [47]	74	117	24.8	46.9	93.6	94.2	953	965	910	872	4.01	4.16	656	687	0.7	0.7	1.5
MSG [22]	<u>41</u>	<b>59</b>	3.2	<u>6.8</u>	80.6	84.6	708	785	<b>510</b>	<b>610</b>	1.62	1.85	292	364	0.0	9.6	9.6
DIP-v	42	67	<u>2.7</u>	<b>5.9</b>	85.2	88.8	804	884	<u>579</u>	674	1.94	2.48	343	435	<u>25.7</u>	<b>44.1</b>	<u>69.9</u>
MSG-v	<b>33</b>	<u>65</u>	<b>2.5</b>	<b>5.9</b>	<b>71.4</b>	<u>77.6</u>	<b>536</b>	<b>643</b>	670	702	<b>1.29</b>	<b>1.43</b>	<b>211</b>	<b>268</b>	<b>67.6</b>	16.2	<b>83.8</b>
Recycle																	
Bicubic	587	880	9.2	16.6	<b>70.6</b>	<b>78.6</b>	<b>575</b>	721	474	576	<b>1.23</b>	<b>1.17</b>	329	398	0.0	<u>11.0</u>	11.0
SRfS [14]	47	72	10.2	22.1	86.1	88.8	715	772	610	623	1.68	<u>1.81</u>	376	410	0.0	0.0	0.0
PDN [39]	95	128	90.5	79.8	84.0	85.7	635	<b>701</b>	523	589	1.66	2.18	364	457	0.0	6.6	6.6
DG [13]	39	82	<u>3.4</u>	11.7	81.6	83.6	696	<u>719</u>	602	617	1.75	1.99	<u>328</u>	<u>383</u>	<u>2.9</u>	<b>65.4</b>	<u>68.4</u>
DIP [47]	<u>29</u>	<u>45</u>	3.9	9.3	91.0	91.8	871	923	576	605	2.95	3.31	434	500	1.5	5.9	7.4
MSG [22]	106	1182	5.8	11.9	82.8	89.6	741	869	624	661	2.60	3.01	485	550	0.7	0.0	0.7
DIP-v	<b>20</b>	<b>34</b>	<b>1.5</b>	<b>4.2</b>	78.9	85.0	<b>575</b>	735	<b>388</b>	<b>485</b>	<u>1.56</u>	1.86	<b>273</b>	<b>332</b>	<b>94.9</b>	3.7	<b>98.5</b>
MSG-v	51	76	3.9	<u>7.9</u>	<u>73.9</u>	<u>82.1</u>	<u>603</u>	737	520	<u>564</u>	1.66	2.02	368	473	0.0	7.4	7.4

Table 9: Quantitative evaluation on samples with small number of missing measurements from Middlebury dataset. The best result is in bold, the second best is underlined.



Umbrella																	
	RMSE <sub>d</sub>		BadPix <sub>d</sub> (5cm)		BadPix <sub>v</sub> (5)		DSSIM <sub>v</sub>		LPIPS <sub>v</sub>		Bumpiness <sub>d</sub>		RMSE <sub>v</sub>		User, 1st	User, 2nd	Top 2
	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x4	x4
Bicubic	1013	1507	6.9	12.1	<b>77.7</b>	<b>80.9</b>	<b>749</b>	<u>837</u>	747	886	<b>0.60</b>	<b>0.60</b>	<u>323</u>	<u>380</u>	<u>5.9</u>	<b>35.3</b>	<u>41.2</u>
SRfS [14]	148	217	19.4	35.5	87.5	90.8	843	853	797	831	0.71	0.78	397	443	0.0	0.0	0.0
PDN [39]	220	287	94.9	89.1	86.6	88.1	799	<b>828</b>	847	882	0.79	1.13	367	452	3.7	<u>22.8</u>	26.5
DG [13]	365	507	9.1	20.3	84.6	87.3	846	878	781	856	0.92	1.36	399	457	0.0	0.7	0.7
DIP [47]	138	<u>145</u>	48.5	21.6	90.5	93.2	915	953	737	<u>722</u>	1.19	1.65	467	528	2.9	16.2	19.1
MSG [22]	292	555	7.4	12.4	84.3	88.1	834	896	<u>678</u>	<u>787</u>	1.27	1.47	442	496	0.0	0.7	0.7
DIP-v	<b>91</b>	<b>129</b>	<b>3.4</b>	<b>5.7</b>	83.4	85.3	796	854	<b>604</b>	<b>598</b>	<u>0.67</u>	0.79	<b>318</b>	<b>352</b>	<b>82.4</b>	8.1	<b>90.4</b>
MSG-v	<u>129</u>	218	<u>5.2</u>	<u>9.7</u>	<u>79.1</u>	<u>82.3</u>	<u>778</u>	842	800	890	0.72	0.89	348	427	5.1	16.2	21.3
Classroom1																	
Bicubic	966	1371	6.7	<u>9.0</u>	<b>75.8</b>	<b>78.3</b>	<b>636</b>	<b>728</b>	<u>581</u>	784	<b>0.41</b>	<b>0.30</b>	<u>268</u>	<b>295</b>	12.5	<b>37.5</b>	<u>50.0</u>
SRfS [14]	135	202	18.5	28.5	82.6	85.7	761	781	718	756	0.62	<u>0.62</u>	332	363	0.0	0.0	0.0
PDN [39]	239	324	96.0	91.0	81.5	82.9	739	759	751	807	0.62	0.76	279	342	<u>16.9</u>	<u>26.5</u>	43.4
DG [13]	307	503	8.8	16.8	82.0	82.7	743	762	766	812	0.74	0.87	313	337	0.0	1.5	1.5
DIP [47]	<u>96</u>	<u>145</u>	17.0	22.4	94.4	94.6	956	952	789	751	1.94	2.12	540	557	0.0	1.5	1.5
MSG [22]	297	408	7.3	10.0	81.2	83.8	723	810	626	<u>604</u>	0.90	1.01	351	391	0.0	0.7	0.7
DIP-v	<b>69</b>	<b>117</b>	<b>4.1</b>	9.3	81.0	86.0	700	789	<b>516</b>	<b>537</b>	0.64	0.86	<b>266</b>	<u>327</u>	<b>64.0</b>	18.4	<b>82.4</b>
MSG-v	127	203	<u>5.4</u>	<b>8.4</b>	<u>76.9</u>	<u>79.4</u>	<u>678</u>	<u>735</u>	739	803	<u>0.60</u>	0.64	283	330	2.2	11.8	14.0

Table 10: Quantitative evaluation on samples with small number of missing measurements from Middlebury dataset. The best result is in bold, the second best is underlined.

Playroom																	
	RMSE <sub>d</sub>		BadPix <sub>d</sub> (5cm)		BadPix <sub>v</sub> (5)		DSSIM <sub>v</sub>		LPIPS <sub>v</sub>		Bumpiness <sub>d</sub>		RMSE <sub>v</sub>		User, 1st	User, 2nd	Top 2
	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x8	x4	x4	x4
Bicubic	1263	1744	14.4	20.4	<b>72.0</b>	<b>76.9</b>	<b>684</b>	<b>783</b>	<u>509</u>	675	<b>0.80</b>	<b>0.52</b>	<u>386</u>	<u>441</u>	0.0	2.2	2.2
SRfS [14]	97	151	26.9	42.1	88.1	91.2	802	829	663	715	<u>1.24</u>	<u>1.08</u>	493	540	0.0	0.0	0.0
PDN [39]	181	253	85.0	69.7	86.3	89.2	820	862	583	656	1.54	1.88	472	543	<u>25.7</u>	<b>61.8</b>	<u>87.5</u>
DG [13]	425	133	22.4	25.0	85.4	86.0	845	826	779	691	1.96	1.63	519	469	1.5	2.2	3.7
DIP [47]	<u>58</u>	<u>91</u>	18.4	20.0	93.0	93.2	941	937	647	<u>612</u>	3.09	2.86	602	592	1.5	2.9	4.4
MSG [22]	433	349	16.1	22.3	85.8	89.9	855	911	685	705	2.51	2.74	576	616	0.0	0.0	0.0
DIP-v	<b>49</b>	<b>83</b>	<b>5.4</b>	<b>12.2</b>	83.8	88.5	728	847	<b>459</b>	<b>530</b>	1.29	1.52	<b>357</b>	<b>433</b>	<b>70.6</b>	<u>27.2</u>	<b>97.8</b>
MSG-v	112	166	<u>9.4</u>	<u>15.5</u>	<u>75.2</u>	<u>80.1</u>	<u>721</u>	810	565	615	1.46	1.67	453	510	0.0	3.7	3.7
Backpack																	
Bicubic	985	1078	14.3	11.5	<b>62.7</b>	<b>69.4</b>	<b>639</b>	<b>730</b>	<u>564</u>	692	<b>0.60</b>	<b>0.45</b>	<b>392</b>	<b>424</b>	2.2	<u>34.6</u>	36.8
SRfS [14]	69	83	18.9	25.5	89.9	89.9	831	847	630	651	1.37	1.26	500	505	0.0	0.0	0.0
PDN [39]	173	207	81.6	65.2	80.4	85.4	770	820	609	719	1.59	1.96	519	553	<u>3.7</u>	<b>37.5</b>	<u>41.2</u>
DG [13]	325	465	26.5	39.9	77.0	82.0	765	808	650	696	1.62	2.07	529	545	0.7	2.9	3.7
DIP [47]	<u>41</u>	<u>67</u>	<u>8.2</u>	17.9	93.2	94.5	943	984	766	692	3.36	2.95	639	645	1.5	11.8	13.2
MSG [22]	211	170	15.1	10.4	76.5	86.9	762	856	671	723	2.11	2.31	577	609	0.7	0.0	0.7
DIP-v	<b>38</b>	<b>62</b>	<b>6.2</b>	12.9	82.5	88.7	677	768	<b>457</b>	<b>496</b>	1.31	1.43	<u>409</u>	<u>448</u>	<b>90.4</b>	5.1	<b>95.6</b>
MSG-v	113	89	10.8	<b>5.9</b>	<u>65.3</u>	<u>72.5</u>	<u>663</u>	<u>752</u>	<u>577</u>	<u>635</u>	<u>1.07</u>	<u>1.15</u>	462	480	0.0	5.9	5.9
Jadeplant																	
Bicubic	1017	1297	19.3	18.4	<b>68.8</b>	<b>75.7</b>	<u>695</u>	<u>788</u>	<u>545</u>	696	<b>0.97</b>	<b>0.62</b>	<b>449</b>	<b>464</b>	2.3	<u>27.7</u>	30.0
SRfS [14]	105	143	39.5	48.9	87.2	92.7	787	839	637	719	1.96	1.70	551	583	0.0	0.0	0.0
PDN [39]	161	205	81.8	62.0	82.4	88.0	778	849	551	<u>625</u>	1.95	2.20	512	572	<u>19.1</u>	<b>41.4</b>	<u>60.5</u>
DG [13]	326	512	27.8	47.6	82.6	86.8	791	823	718	670	2.28	2.66	567	601	0.0	0.5	0.5
DIP [47]	<b>70</b>	<u>121</u>	<u>16.7</u>	32.8	91.2	92.3	913	911	735	764	3.19	3.21	615	638	0.0	0.5	0.5
MSG [22]	216	263	21.1	<u>17.8</u>	81.5	87.6	796	880	751	783	2.73	2.90	614	649	0.0	0.0	0.0
DIP-v	84	<u>121</u>	21.8	24.0	86.6	89.5	820	870	<b>542</b>	654	1.92	1.99	<u>503</u>	<u>535</u>	<b>78.2</b>	20.5	<b>98.6</b>
MSG-v	109	<b>117</b>	<b>13.9</b>	<b>11.2</b>	<u>71.0</u>	<u>79.0</u>	<b>688</b>	<b>781</b>	605	<b>622</b>	<u>1.61</u>	<u>1.66</u>	507	<u>529</u>	0.5	8.2	8.6

Table 11: Quantitative evaluation on samples with large number of missing measurements from Middlebury dataset. The best result is in bold, the second best is underlined.

## References

- [1] Gianluca Agresti, Ludovico Minto, Giulio Marin, and Pietro Zanuttigh. Deep learning for confidence information in stereo and tof data fusion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 697–705, 2017. 2
- [2] Piotr Bojanowski, Armand Joulin, David Lopez-Pas, and Arthur Szlam. Optimizing the latent space of generative networks. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 600–609, Stockholmssan, Stockholm Sweden, 10–15 Jul 2018. PMLR. 5
- [3] Daniel J. Butler, Jonas Wulff, Garrett B. Stanley, and Michael J. Black. A naturalistic open source movie for optical flow evaluation. In A. Fitzgibbon et al. (Eds.), editor, *European Conf. on Computer Vision (ECCV)*, Part IV, LNCS 7577, pages 611–625. Springer-Verlag, Oct. 2012. 6
- [4] Baoliang Chen and Cheolkon Jung. Single depth image super-resolution using convolutional neural networks. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1473–1477. IEEE, 2018. 3
- [5] Zhao Chen, Vijay Badrinarayanan, Gilad Drozdov, and Andrew Rabinovich. Estimating depth from rgb and sparse sensing. *CoRR*, abs/1804.02771, 2018. 2
- [6] Xinjing Cheng, Peng Wang, and Ruigang Yang. Depth estimation via affinity learned with convolutional spatial propagation network. In *European Conference on Computer Vision*, pages 108–125. Springer, Cham, 2018. 2
- [7] Manri Cheon, Jun-Hyuk Kim, Jun-Ho Choi, and Jong-Seok Lee. Generative adversarial network-based image super-resolution using perceptual content losses. In Laura Leal-Taixé and Stefan Roth, editors, *ECCV Workshops*, pages 51–62, Cham, 2019. Springer International Publishing. 3
- [8] Nathaniel Chodosh, Chaoyang Wang, and Simon Lucey. Deep convolutional compressed sensing for lidar depth completion. *arXiv preprint arXiv:1803.08949*, 2018. 2
- [9] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. In *Advances in neural information processing systems*, pages 2366–2374, 2014. 2
- [10] David Ferstl, Christian Reinbacher, Rene Ranftl, Matthias Rüther, and Horst Bischof. Image guided depth upsampling using anisotropic total generalized variation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 993–1000, 2013. 5
- [11] David Ferstl, Matthias Ruther, and Horst Bischof. Variational depth superresolution using example-based edge representations. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 513–521, 2015. 3
- [12] Muhammad Waleed Gondal, Bernhard Schölkopf, and Michael Hirsch. The unreasonable effectiveness of texture transfer for single image super-resolution. In *ECCV*, pages 80–97. Springer, 2018. 3
- [13] Shuhang Gu, Wangmeng Zuo, Shi Guo, Yunjin Chen, Chongyu Chen, and Lei Zhang. Learning dynamic guidance for depth image enhancement. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3769–3778, 2017. 3, 4, 5, 7, 8, 11, 20, 21, 22, 23
- [14] Bjoern Haefner, Yvain Quéau, Thomas Möllenhoff, and Daniel Cremers. Fight ill-posedness with ill-posedness: Single-shot variational depth super-resolution from shading. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 164–174, 2018. 3, 4, 5, 7, 8, 11, 20, 21, 22, 23
- [15] Bumsu Ham, Minsu Cho, and Jean Ponce. Robust guided image filtering using nonconvex potentials. *IEEE transactions on pattern analysis and machine intelligence*, 40(1):192–207, 2018. 3
- [16] Wei Han, Shiyu Chang, Ding Liu, Mo Yu, Michael Witbrock, and Thomas S Huang. Image super-resolution via dual-state recurrent networks. In *Proc. CVPR*, 2018. 3
- [17] Ankur Handa, Thomas Whelan, John McDonald, and Andrew J. Davison. A benchmark for RGB-D visual odometry, 3D reconstruction and SLAM. In *IEEE Intl. Conf. on Robotics and Automation, ICRA*, Hong Kong, China, May 2014. 5
- [18] Muhammad Haris, Greg Shakhnarovich, and Norimichi Ukita. Deep back-projection networks for super-resolution. In *Proc. CVPR*, 2018. 3
- [19] Katrin Honauer, Ole Johannsen, Daniel Kondermann, and Bastian Goldluecke. A dataset and evaluation methodology for depth estimation on 4d light fields. In *Asian Conference on Computer Vision*, pages 19–34. Springer, 2016. 2, 9
- [20] Katrin Honauer, Lena Maier-Hein, and Daniel Kondermann. The hci stereo metrics: Geometry-aware performance analysis of stereo algorithms. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2120–2128, 2015. 2, 9
- [21] Jiashen Hua and Xiaojin Gong. A normalized convolutional neural network for guided sparse depth upsampling. In *IJ-CAI*, pages 2283–2290, 2018. 2
- [22] Tak-Wai Hui, Chen Change Loy, and Xiaoou Tang. Depth map super-resolution by deep multi-scale guidance. In *European Conference on Computer Vision*, pages 353–369. Springer, 2016. 2, 4, 7, 8, 11, 20, 21, 22, 23
- [23] Forrest N. Iandola, Song Han, Matthew W. Moskewicz, Khalid Ashraf, William J. Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and 0.5mb model size. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 4
- [24] Zhongyu Jiang, Yonghong Hou, Huanjing Yue, Jingyu Yang, and Chunping Hou. Depth super-resolution from rgb-d pairs with transform and spatial domain regularization. *IEEE Transactions on Image Processing*, 27(5):2587–2602, 2018. 3
- [25] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, 2016. 3
- [26] Beomjun Kim, Jean Ponce, and Bumsu Ham. Deformable Kernel Networks for Joint Image Filtering. working paper or preprint, Oct. 2018. 2
- [27] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural net-

- works. In *Advances in neural information processing systems*, pages 1097–1105, 2012. 4
- [28] Christian Ledig, Lucas Theis, Ferenc Huszar, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, and Wenzhe Shi. Photo-realistic single image super-resolution using a generative adversarial network. In *Proc. CVPR*, pages 105–114, 2017. 3
- [29] Beichen Li, Yuan Zhou, Yeda Zhang, and Aihua Wang. Depth image super-resolution based on joint sparse coding. *Pattern Recognition Letters*, 2018. 3
- [30] Yijun Li, Jia-Bin Huang, Narendra Ahuja, and Ming-Hsuan Yang. Deep joint image filtering. In *European Conference on Computer Vision*, pages 154–169. Springer, 2016. 2
- [31] Xiaotong Luo, Rong Chen, Yuan Xie, Yanyun Qu, and Cuihua Li. Bi-gans-st for perceptual image super-resolution. In Laura Leal-Taixé and Stefan Roth, editors, *ECCV Workshops*, pages 20–34, Cham, 2019. Springer International Publishing. 3
- [32] Chao Ma, Chih-Yuan Yang, Xiaokang Yang, and Ming-Hsuan Yang. Learning a no-reference quality metric for single-image super-resolution. *Computer Vision and Image Understanding*, 158:1–16, 2017. 3
- [33] Fangchang Ma, Guilherme Venturelli Cavalheiro, and Sertac Karaman. Self-supervised sparse-to-dense: Self-supervised depth completion from lidar and monocular camera. *arXiv preprint arXiv:1807.00275*, 2018. 2
- [34] Fangchang Ma and Sertac Karaman. Sparse-to-dense: Depth prediction from sparse depth samples and a single image. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1–8. IEEE, 2018. 2, 4
- [35] Rafat Mantiuk, Kil Joong Kim, Allan G Rempel, and Wolfgang Heidrich. Hdr-vdp-2: a calibrated visual metric for visibility and quality predictions in all luminance conditions. In *ACM Transactions on graphics (TOG)*, volume 30, page 40. ACM, 2011. 2
- [36] Roey Mechrez, Itamar Talmi, Firas Shama, and Lihi Zelnik-Manor. Learning to maintain natural image statistics. *arXiv preprint arXiv:1803.04626*, 2018. 3
- [37] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. In *NIPS-W*, 2017. 6
- [38] Songyou Peng, Bjoern Haefner, Yvain Queau, and Daniel Cremers. Depth super-resolution meets uncalibrated photometric stereo. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2961–2968, 2017. 3
- [39] David Riegler, Gernot aand Ferstl, Matthias Rüther, and Horst Bischof. A deep primal-dual network for guided depth super-resolution. In *British Machine Vision Conference*. The British Machine Vision Association, 2016. 2, 4, 5, 7, 8, 11, 20, 21, 22, 23
- [40] Daniel Scharstein, Heiko Hirschmüller, York Kitajima, Greg Krathwohl, Nera Nešić, Xi Wang, and Porter Westling. High-resolution stereo datasets with subpixel-accurate ground truth. In *German Conference on Pattern Recognition*, pages 31–42. Springer, 2014. 5
- [41] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014. 4
- [42] Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. Sun rgb-d: A rgb-d scene understanding benchmark suite. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 567–576, 2015. 5
- [43] Xibin Song, Yuchao Dai, and Xueying Qin. Deep depth super-resolution: Learning depth super-resolution using deep convolutional neural network. In *Asian Conference on Computer Vision*, pages 360–376. Springer, 2016. 2, 3
- [44] Xibin Song, Yuchao Dai, and Xueying Qin. Deeply supervised depth map super-resolution as novel view synthesis. *IEEE Transactions on Circuits and Systems for Video Technology*, 2018. 3
- [45] Atsuhiko Tsuchiya, Daisuko Sugimura, and Takayuki Hamamoto. Depth upsampling by depth prediction. In *2017 IEEE International Conference on Image Processing (ICIP)*, pages 1662–1666, Sept 2017. 2
- [46] Jonas Uhrig, Nick Schneider, Lukas Schneider, Uwe Franke, Thomas Brox, and Andreas Geiger. Sparsity invariant cnns. In *IEEE International Conference on 3D Vision (3DV)*, 2017. 2
- [47] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Deep image prior. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 4, 5, 7, 8, 11, 20, 21, 22, 23
- [48] Stefan Van der Walt, Johannes L Schönberger, Juan Nunez-Iglesias, François Boulogne, Joshua D Warner, Neil Yager, Emmanuelle Gouillart, and Tony Yu. scikit-image: image processing in python. *PeerJ*, 2:e453, 2014. 9
- [49] Thang Vu, Tung M. Luu, and Chang D. Yoo. Perception-enhanced image super-resolution via relativistic generative adversarial networks. In Laura Leal-Taixé and Stefan Roth, editors, *ECCV 2018 Workshops*, pages 98–113, Cham, 2019. Springer International Publishing. 3
- [50] Xintao Wang, Ke Yu, Chao Dong, and Chen Change Loy. Recovering realistic texture in image super-resolution by deep spatial feature transform. In *Proc. CVPR*, June 2018. 3
- [51] Yifan Wang, Federico Perazzi, Brian McWilliams, Alexander Sorkine-Hornung, Olga Sorkine-Hornung, and Christopher Schroers. A fully progressive approach to single-image super-resolution. In *CVPR Workshops*, June 2018. 3
- [52] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 2, 4
- [53] Zhou Wang, Eero P Simoncelli, and Alan C Bovik. Multiscale structural similarity for image quality assessment. In *The Thirty-Seventh Asilomar Conference on Signals, Systems & Computers*, 2003, volume 2, pages 1398–1402. IEEE, 2003. 2
- [54] Jun Xie, Rogerio Schmidt Feris, and Ming-Ting Sun. Edge-guided single depth image super resolution. *IEEE Transactions on Image Processing*, 25(1):428–438, 2016. 4, 5, 7, 8, 11, 20, 21, 22



- [55] Shi Yan, Chenglei Wu, Lizhen Wang, Feng Xu, Liang An, Kaiwen Guo, and Yebin Liu. Ddnet: Depth map denoising and refinement for consumer depth cameras using cascaded cnns. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 151–167, 2018. 2, 4
- [56] Jingyu Yang, Xinchun Ye, Kun Li, Chunping Hou, and Yao Wang. Color-guided depth recovery from rgb-d data using an adaptive autoregressive model. *IEEE transactions on image processing*, 23(8):3443–3458, 2014. 3
- [57] Lin Zhang, Lei Zhang, Xuanqin Mou, and David Zhang. Fsim: a feature similarity index for image quality assessment. *IEEE transactions on Image Processing*, 20(8):2378–2386, 2011. 2
- [58] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 2, 3, 4, 9
- [59] Yulun Zhang, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong, and Yun Fu. Image super-resolution using very deep residual channel attention networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 286–301, 2018. 3
- [60] Yulun Zhang, Yapeng Tian, Yu Kong, Bineng Zhong, and Yun Fu. Residual dense network for image super-resolution. In *Proc. CVPR*, June 2018. 3
- [61] Hang Zhao, Orazio Gallo, Iuri Frosio, and Jan Kautz. Loss functions for image restoration with neural networks. *IEEE Transactions on Computational Imaging*, 3(1):47–57, 2017. 4
- [62] Lijun Zhao, Huihui Bai, Jie Liang, Bing Zeng, Anhong Wang, and Yao Zhao. Simultaneously color-depth super-resolution with conditional generative adversarial network. *arXiv preprint arXiv:1708.09105*, 2017. 3
- [63] Yifan Zuo, Qiang Wu, Jian Zhang, and Ping An. Minimum spanning forest with embedded edge inconsistency measurement model for guided depth map enhancement. *IEEE Transactions on Image Processing*, 27(8):4145–4159, 2018. 3