

# 3D Appearance Super-Resolution with Deep Learning

Yawei Li<sup>1</sup>, Vagia Tsiminaki<sup>2</sup>, Radu Timofte<sup>1</sup>, Marc Pollefeys<sup>2,3</sup>, and Luc van Gool<sup>1</sup>

<sup>1</sup>Computer Vision Lab, ETH Zurich, Switzerland

{yawei.li, radu.timofte, vangool}@vision.ee.ethz.ch

<sup>2</sup>Computer Vision and Geometry Group, ETH Zurich, Switzerland, <sup>3</sup>Microsoft, USA

{vagia.tsiminaki, marc.pollefeys}@inf.ethz.ch

## Abstract

We tackle the problem of retrieving high-resolution (HR) texture maps of objects that are captured from multiple view points. In the multi-view case, model-based super-resolution (SR) methods have been recently proved to recover high quality texture maps. On the other hand, the advent of deep learning-based methods has already a significant impact on the problem of video and image SR. Yet, a deep learning-based approach to super-resolve the appearance of 3D objects is still missing. The main limitation of exploiting the power of deep learning techniques in the multi-view case is the lack of data. We introduce a 3D appearance SR (3DASR) dataset based on the existing ETH3D [42], SyB3R [31], Middlebury, and our Collection of 3D scenes from TUM [21], Fountain [51] and Relief [53]. We provide the high- and low-resolution texture maps, the 3D geometric model, images and projection matrices. We exploit the power of 2D learning-based SR methods and design networks suitable for the 3D multi-view case. We incorporate the geometric information by introducing normal maps and further improve the learning process. Experimental results demonstrate that our proposed networks successfully incorporate the 3D geometric information and super-resolve the texture maps.

## 1. Introduction

Retrieving efficiently the appearance information of objects through multi-camera observations is of a great importance for the final goal of creating realistic 3D content. To increase the realism of the reconstructed 3D object a detailed appearance needs to be added on top of geometry. This high quality 3D content is used in applications such as movie production, video games and digital culture heritage preservation. Yet, even with highly accurate 3D geometric reconstruction, simply re-projecting the images onto the geometry does not guarantee detailed appearance coverage.

To regain details from the low-resolution (LR) images,

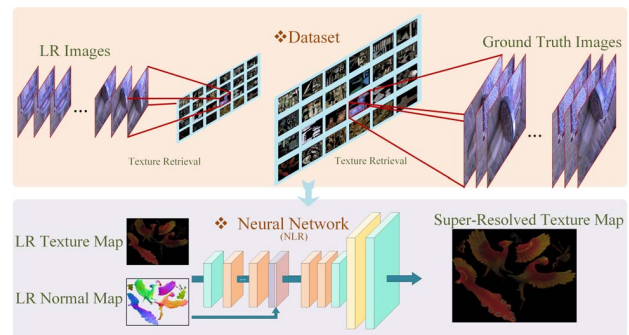


Figure 1: We introduce the 3DASR, a 3D appearance SR dataset and a deep learning-based approach to super-resolve the appearance of 3D objects.

model-based super-resolution (SR) techniques have been introduced in the multi-view case [22, 21, 45]. These methods introduce a single coherent texture space to define a common texture map and they model the captured image as a downgraded version of this high-resolution (HR) texture map. Through image formation model they exploit the visual redundancy of the overlapping views [22, 21] and of video frames [45]. Although these model-based SR techniques recover successfully high quality texture maps, they are computationally demanding.

On the other hand, 2D example-based SR methods have been shown to outperform the model-based methods. The basic assumption of example-based SR is the recurrence of similar patches in different parts of an image or in different images [18]. In particular, recent deep learning-based techniques have been proposed to learn the mapping between the LR and HR images. Different networks are trained on large image datasets that contain pairs of HR and LR images. Super-resolving LR images is then realized with a feed forward step. Yet, a deep learning-based approach to super-resolve the appearance of 3D objects is still missing.

In this paper, our goal is to introduce deep learning techniques into the problem of appearance SR in the multi-view case. To exploit the capacity of 2D deep learning techniques, we first provide a 3D appearance dataset. Similar to

the model-based SR methods, we introduce a common texture space and define a single coherent texture map. This texture map is first mapped onto the geometry. Then the textured surface is projected into the image space. We express the concatenation of these two mappings through the image formation model (Fig. 2). Through this image generation process and using captured images of multiple scaling factors we can then recover the corresponding texture maps. We provide a dataset that contains ground truth HR texture maps together with LR texture maps of down-scaling factor  $\times 2$ ,  $\times 3$ , and  $\times 4$ . The dataset covers both synthetic scenes SyB3R [31] and real scenes ETH3D [42], MiddleBury, and our Collection of 3D scenes from TUM [21], *Fountain* [51] and *Relief* [53]. We then leverage the capacity of 2D learning-based methods [36] and design two architectures suitable for the 3D multi-view case. Similar to [27] we introduce normal maps to capture the local structure of the 3D model and incorporate the 3D geometric information into the 2D SR network. To our knowledge, our work is the first that introduces deep learning approaches for the appearance SR in the multi-view case. Using our provided dataset, we evaluate different texture map SR methods including interpolation-based, model-based, and learning-based. In summary, the contributions of our paper are:

1. a 3D texture dataset that contains pairs of HR and LR textures of 3D objects. With this dataset we facilitate the integration of deep learning techniques into the problem of appearance SR in the multi-view case and we open up a promising novel research direction. We refer to the dataset as 3DASR.
2. the first appearance SR framework that elegantly combines the power of 2D deep learning-based techniques with the 3D geometric information in the multi-view setting.

The rest of the paper is organized as follows. Sec. 2 introduces related works of this paper. Sec. 3 describes how the texture maps are retrieved. Sec. 4 explains the generation process of the dataset. Sec. 5 explores the introduction of normal information into neural networks to super-resolve LR texture maps. Sec. 6 shows the evaluation results of different methods. Sec. 7 concludes the paper.

## 2. Related Works

### 2.1. 2D image super-resolution

2D image SR has been extensively studied and it can be classified into three categories, i.e. interpolation-based, model-based, and example-based [40, 17, 48, 18]. Although a comprehensive review of these methods is beyond the scope of this paper, we present the underlying concepts of each of them. Interpolation-based methods [2, 32] increase

the resolution by computing pixel values using the neighbouring information. But leveraging only the local information within the image cannot guarantee the recovery of high-frequency details. Model-based approaches describe the LR image as downgraded version of the HR image and express analytically the forward degradation system. Solving for the inverse problem prior knowledge over the unknown HR image such as smoothness and non-local similarity [8, 34] is imposed. Treating the problem as a stochastic process, maximum likelihood [17] or maximum a posterior [19] approach is followed. Although these methods successfully recover high-frequency details, they require elegant optimization techniques. Most of the times they correspond to iterative approaches that are computationally heavy and time-consuming. Learning-based methods shift this computational burden to the learning phase and using the trained network they super-resolve the image through a feed forward step. Due to the availability of large datasets, carefully designed network architectures can learn the mapping from LR to HR image and achieve state-of-the-art performance [14, 44, 28, 36, 33, 50]. Our work, introduces deep learning-based approach in the multi-view case to retrieve the fine texture of 3D objects.

### 2.2. Texture retrieval

Adding a high quality texture layer onto the 3D geometry plays an essential role in the final realism. This is a challenging step since in the multi-view case there are additional sources of variation that we need to account for, namely occlusions, calibration and reconstruction inaccuracies. Several methods have been proposed in the literature [23] to efficiently exploit all the available color information and to address the aforementioned challenges.

**Single view selection.** To cope with different geometric inaccuracies, several methods use only one view to assign texture to each face. Lempitsky and Ivanon [29] compensate for seams between the boundaries of each face by solving a discrete labeling problem. Gal *et al.* [20] incorporate in their optimization the effect of foreshortening, image resolution, and blur by modifying the weighting function. Waechter *et al.* [46] add an additional smoothness term to penalize inconsistencies between adjacent faces. By choosing a single view, these methods disregard the multiple color information that exists in the multi-view setting.

**Multi-view selection.** To leverage the multiple color information over views, several methods blend the images for each face. Debevec *et al.* [12] reproject and blend view contributions according to visibility and viewpoint-to-surface angle. To capture view dependent shading effects Buehler *et al.* [9] model and approximate the plenoptic function for the scene object. Some hybrid approaches [3, 10] select a single view per face and blend in frequency space views close to texture patch borders. To correct geometric inaccu-

racies, in [52] camera poses are jointly optimized with the photometric consistency. Following the success of patch-based synthesis methods, Bi *et al.* propose a single view-independent texture mapping method that account for geometric misalignment [7]. Generally these methods do not exploit efficiently viewpoint visual redundancy.

**Multi-view texture SR methods.** To retrieve fine appearance details, a handful of texture SR methods have leveraged the SR principle in the multi-view case and compute texture maps with a resolution higher than the input images [25, 39]. Goldlücke *et al.* introduce an image formation model to super-resolve texture maps [22] and to refine the geometry and camera calibration [21]. Tsiminaki *et al.* [45] further improve SR texture quality by exploiting additional temporal redundancy and by uniformly correcting calibration and geometry errors with optical flow. These methods are however computationally expensive.

We alleviate the limitations of these model-based SR by introducing the deep learning-based approaches that have been proven to outperform in the 2D case.

### 2.3. Super-resolution benchmark

In order to be able to use deep learning-based techniques for super-resolving the texture of 3D objects, datasets need to be available. For 2D image SR there are several benchmarking datasets Set5 [6], Set14 [49], Urban100 [24], BSD100 [38] and works [47, 26]. ImageNet [13] has been also used as training dataset in several example based approaches [14, 15]. More recently, DIV2K dataset was introduced to provide higher quality images [1].

Such data are however not available in the multi-view case. We propose in this work a methodology to compute textures of several resolution and we provide a 3D texture dataset, 3DASR, that contains pairs of HR and LR textures of 3D objects.

## 3. Texture Retrieval

### 3.1. Image formation model

The image formation model simulates the generation of the image from the unknown texture map. In Fig. 2, we can distinguish two steps *i.e.*, texture mapping and projection to image space.

**Texture mapping** The texture mapping function  $\mu$  assigns each entity of the texture map (texel) to a 3D point of the geometry. In order to be able to define the texture map and the mapping, we first need to parameterize the geometry in a common space. We assume that the 3D model  $M$  is a known triangulated mesh and thus we can define any UV parameterization. In [4] advanced algorithms that result in space-optimized texture maps are discussed. In this work we use a fixed UV parameterization, described in

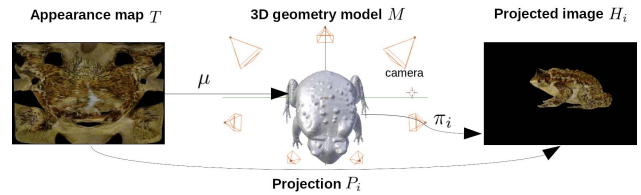


Figure 2: Image formation model.

Subsec. 4.1. Through this mapping function  $\mu$ , a texel  $x$  is mapped to a point  $\mu(x)$  of the 3D mesh model  $M$ .

**Projection to image space** We assume that we know the camera poses and the intrinsic camera parameters. The textured 3D object is then projected into the image space given the known projection matrices. Let  $\pi_i$  be the camera projection matrix at the view point  $i$  and  $H_i$  the corresponding image of resolution  $h \times w$ . The geometric point  $\mu(x)$  is projected to the pixel location  $(\pi_i \circ \mu)(x)$  in the image plane. Let  $T^{h \cdot w}$  and  $H_i^{h \cdot w}$  be the vectorized version of the texture map and the projected image. The image is then expressed as a linear combination of the texture map  $H_i^{h \cdot w} = P_i T^{h \cdot w}$  where  $P$  is a matrix of dimension  $h \cdot w \times h \cdot w$ . To estimate this projection operator several issues need to be addressed. First, two geometric points of the surface might be projected into the same location due the convexity of the geometry and then only the visible color value needs to be selected. Second, this projection step can lead to non-integer locations [35]. Third, the distribution of the projected points in the image space is non-uniform, which means that the points may be sparse for some areas. To combine the contributions of all the projected points falling into the neighborhood of a pixel  $q$  we introduce the Gaussian function as the weighting function. This function takes the location proximity into account, encouraging pixels near the center of  $q$  while penalizing those far away from  $q$ . By combining the contributions of all projected points falling into the neighborhood of a pixel with this Gaussian function we solve for the sparse areas in the image space that can originate due to high curvature regions of the surface.

### 3.2. Texture retrieval: the inverse process

We retrieve the texture maps by inverting the image formation model. We examine several scaling factors including the ground truth high resolution and down-scaling factor  $\times 2, \times 3, \times 4$ . Given the projection matrices with the multi-view images we compute the corresponding texture maps.

## 4. The Dataset: 3DASR

The 3DASR dataset we provide is based on four existing subsets; one synthetic subset SyB3R [31] and three real subsets EHT3D [42], MiddleBury [43], and Collection of *Bird*, *Beethoven* and *Bunny* from the multi-view dataset of TUM [21], *Fountain* [51] and *Relief* [53]. We follow a generic

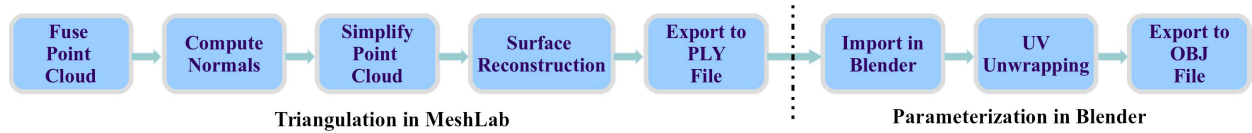


Figure 3: Conversion from point clouds to mesh models and unified parameterization.

pipeline to preprocess all subsets. We compute the triangulated 3D mesh with texture coordinates and vertex normals. We use the images provided by the original dataset as the HR images and we downscale them using scale factors  $\times 2$ ,  $\times 3$ ,  $\times 4$  to compute the corresponding LR images. The projection matrices for the corresponding LR images are derived by RQ matrix decomposition of the original projection matrix and then scaling down the intrinsic parameters.

#### 4.1. The real subsets

ETH3D, Collection, and MiddleBury correspond to real scenes. Regarding ETH3D, we use the training set of the HR multi-view subset that contains 13 scenes. Every scene is provided with multi-view images captured by DSLR cameras, the camera intrinsic and extrinsic parameters, and the ground truth point clouds captured by laser scanners. Collection is a collection of 6 3D scenes. We use the *TempleRing* and *DinoRing* of MiddleBury.

**Mesh: triangulation and UV mapping.** We first compute the triangulated mesh and then unwrap it to define the texture map. Through the UV unwrapping we assign to each vertex a UV coordinate.

For MiddleBury, we use the Multi-View Stereo (MVS) pipeline [41] to reconstruct the meshes. For *Bird*, *Beethoven* and *Bunny* we use the same meshes as in the paper [45] and for *Fountain Relief* the meshes are refined in the work of Maier *et al.* [37]. To ensure low appearance distortion, we use conformal parameterization similar to [22, 45]. We compute a conformal atlas by selecting the algorithm of LSCM [30] that is implemented in Blender.

For ETH3D subset, the provided 3D model is just a point cloud. Therefore, both of the processing steps are needed. Fig. 3 shows the workflow. Note that triangulation is implemented in MeshLab while parameterization is done in Blender. First of all, for most of the scenes, there are multiple point clouds and each of them captures the scene geometry from different viewpoints. Thus, these point clouds are fused to create a fully-fledged scene geometry followed by the computation of normals. The merged point cloud contains tens of millions of points which may become a computation bottleneck for the post-processing. Thus, the point cloud is simplified using Poisson disk sampling [11] which reduces the number of points while maintains the geometric details of the scene. Then the mesh is reconstructed using ball pivoting algorithm [5]. The reconstruction result is exported to a PLY file which is imported into Blender.

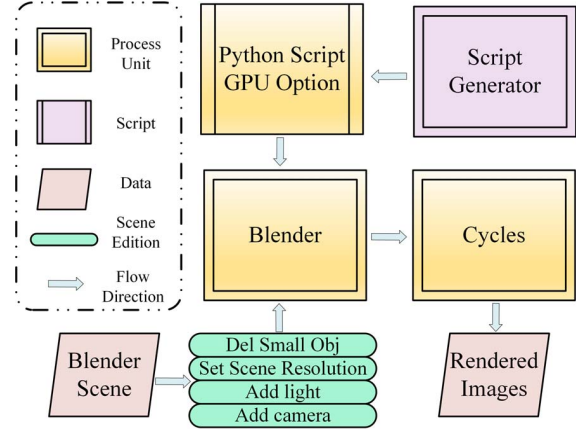


Figure 4: SyB3R image rendering pipeline.

Blender’s UV unwrapping procedure is used for UV parameterization. At last, the triangulated mesh with UV texture coordinates is exported to an OBJ file.

**Images and projection matrices.** We consider the provided images by the original dataset as the HR images and we derive the LR images by down-sampling the HR. The intrinsic and extrinsic parameters are given for ETH3D and MiddleBury. Thus computing the projection matrices is straightforward. For the Collection subset, we use RQ decomposition to compute intrinsic and extrinsic parameters. For all of the three subsets, the projection matrices corresponding to the LR images are derived by down-scaling the intrinsic parameters with  $\times 2$ ,  $\times 3$ , and  $\times 4$  scaling factors.

#### 4.2. The synthetic subset: SyB3R

SyB3R is a synthetic dataset containing four scenes. Each scene contains an accurate geometry mesh model with optimal UV parameterization. The image rendering pipeline is shown in Fig. 4. To speed up the rendering, we add GPU option to the Python script. We edit the synthetic scene by keeping the major object, setting image resolutions, adding lights and cameras. The generated script and altered scene are passed to Blender and Cycles, resulting in the rendered images. The original mesh model of SyB3R contains several separated objects whose texture maps may overlap with each other in the texture space. To address this problem, we only keep the major part of the scene, *i.e.*, the body of *Toad*, the skull of *Skull*, and the single rock of *Geological Sample*. We do not use *Lego Bulldozer* because it consists of many small pieces without meaningful texture.





Figure 5: Rendered images of SyB3R.

**Camera and light.** To capture every surface of the object, 14 cameras are uniformly aligned on the sphere surrounding the object. The focal length of the cameras is 25 mm. The size of the sensor is  $32 \times 18$  mm. To ensure uniform background across the rendered images, 6 lights are added in the scene lighting from the 6 directions of the object.

**Rendering images** The resolution of HR images is  $3888 \times 2592$  while the resolution of the LR images is calculated by dividing the HR width and height with respective scaling factors. Knowing the focal length, image resolution, principal point, rotation matrix and translation vector, the  $3 \times 4$  camera projection matrix is computed. As stated by the authors [31], the rendering time can be multiple hours per image due to the high computational load of the image synthesis process. Thus, we use GPU to render the images. Examples of rendered images are shown in Fig. 5.

### 4.3. Texture maps

After generating these data, we can now use the texture retrieval algorithm and compute the texture maps of 4 different resolutions. Fig. 6 shows the texture maps of the 24 different scenes.

## 5. Learning-Based Methods

Our 3DASR dataset contains pairs of HR and LR texture maps which resemble two dimensional images. This allows us to make use of state-of-the-art 2D deep learning-based image SR methods. Such an integration is however not without its own source of difficulties. Being in the multi-view setting, the geometric information needs also to be encoded. The texture domain has its own characteristics compare to natural images. It is thus important to adapt the 2D SR deep learning-based method to this new domain. We incorporate the 3D geometric information through the normals and we show how to guide the learning process.

### 5.1. Normal information

Normal coordinates can be normalized and stored as pixel colors in normal maps (Fig. 8) which have the same support as the texture maps. These normal maps capture the local structure of the surface. We thus use them into the network to introduce the 3D geometric information. We store them as PNG images with 4 channels. The first 3 channels store the normalized normal coordinates and the fourth alpha channel is a mask that shows the support of the texture map, namely, where texel information is available.

### 5.2. Network architecture

The next essential step is to incorporate the normal maps and adjust the neural network to the multi-view setting. There are two main approaches. The first is to use them directly as input information by concatenating them with the texture maps. The second approach is to interpret them as high-level features and concatenate them with feature maps computed at specific layers of the network. We follow the second approach due to the following two considerations. First, the normal maps encode 3D geometric information and can indeed be seen as high-level feature maps. Second, in the case where the normal maps were used as input, the whole network should be trained from scratch. Given the small size of our 3DASR dataset this would lead to over-fitting. Thus, by introducing them at higher layers we train only the few last layers of the network, fine-tune the lower ones and avoid this way over-fitting.

In order to examine the importance of the geometric information in the performance of the training, we compute the normals in both spaces of the low and high resolution texture maps. We call them LR and HR normal maps accordingly. We use EDSR [36] as a case study network to show the adaption of the network. We thus provide two difference versions, one where the LR normal maps are added before the upsampling layer and a second where the HR normal maps are added after the upsampling layer.

### 5.3. Implementation details

The architecture of the two adapted networks is shown in Fig. 7a and Fig. 7b, which we name as NLR and NHR, representing the utilization of LR and HR normal maps. In Fig. 7a, LR normal maps are concatenated with the feature maps after the 30th ResBlock. The following two ResBlocks and the upsampling layer learn representation from the combined feature map. In Fig. 7b, upsampling layer is moved before the two fine-tuning ResBlocks and the HR normal maps are added directly after the upsampling layer. Four additional convolutional layers follow the two ResBlocks. The number of feature maps after the concatenation becomes 260 which is the sum of the original 256 channels and the additional 4 channels of the normal map.

We name the layers from the starting convolutional layer to the 30th ResBlock as the body part of the network. The remaining layers are referred to as the tail part. The parameters of the body part are loaded from pretrained EDSR model and fine-tuned to adapt to the texture domain while those of the tail part are randomly initialized and trained from scratch. Thus, a larger learning rate  $10^{-4}$  is used to train the tail parameters while a smaller one  $10^{-5}$  is used to fine-tune the body parameters. We also directly fine-tune the EDSR model without any architecture modification. An in-between learning rate  $2.5 \times 10^{-5}$  is used. To train the CNN, the mask is used to identify the active areas of the

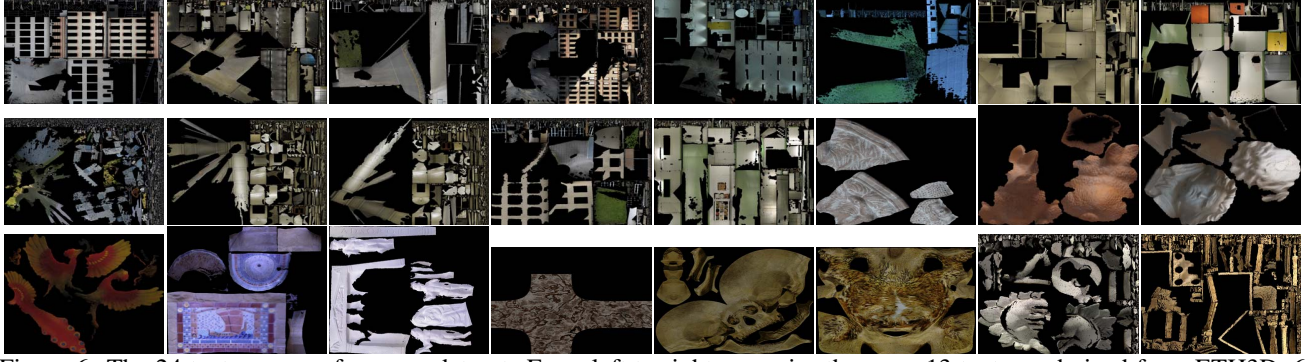


Figure 6: The 24 texture maps from our dataset. From left to right row-wise there are 13 textures derived from ETH3D, 6 from Collection, 3 from SyB3R, and 2 from MiddleBury.

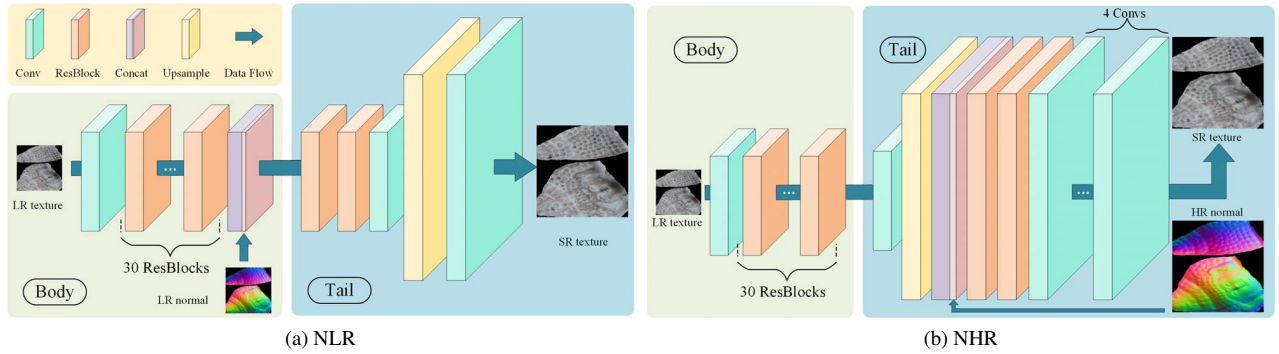


Figure 7: Network structure of (a) NLR and (b) NHR based on the the EDSR [36]. The change in the dimension of the blocks indicates the resolution change of the feature maps. In (a) normal maps are computed in the input low resolution space and are concatenated with the feature map before the upscaling layer. In (b) normal maps are computed in the high resolution space and concatenated after the upscaling layer.

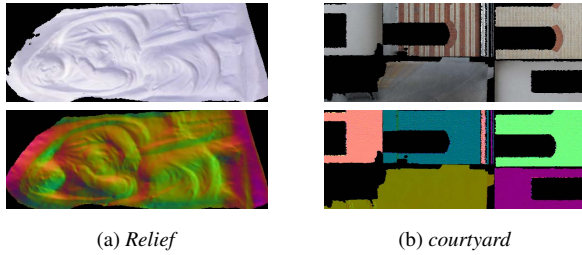


Figure 8: Normal maps capture the local structure of the surface.

texture maps. We crop the texture maps into patches of size  $48 \times 48$  and feed them into the network for training by excluding these ones that have black areas larger than a predefined threshold 50. During inference the CNN is applied on the whole LR texture map.

The provided dataset contains 4 subsets and 24 texture maps in total. Cross-validation is used to get the evaluation result on the whole dataset. That is, we divide the 24 texture map into 2 splits, one for training and one for testing. The texture maps of the 4 subsets are equally distributed to the two splits, thus each with 12 texture maps. In addition,

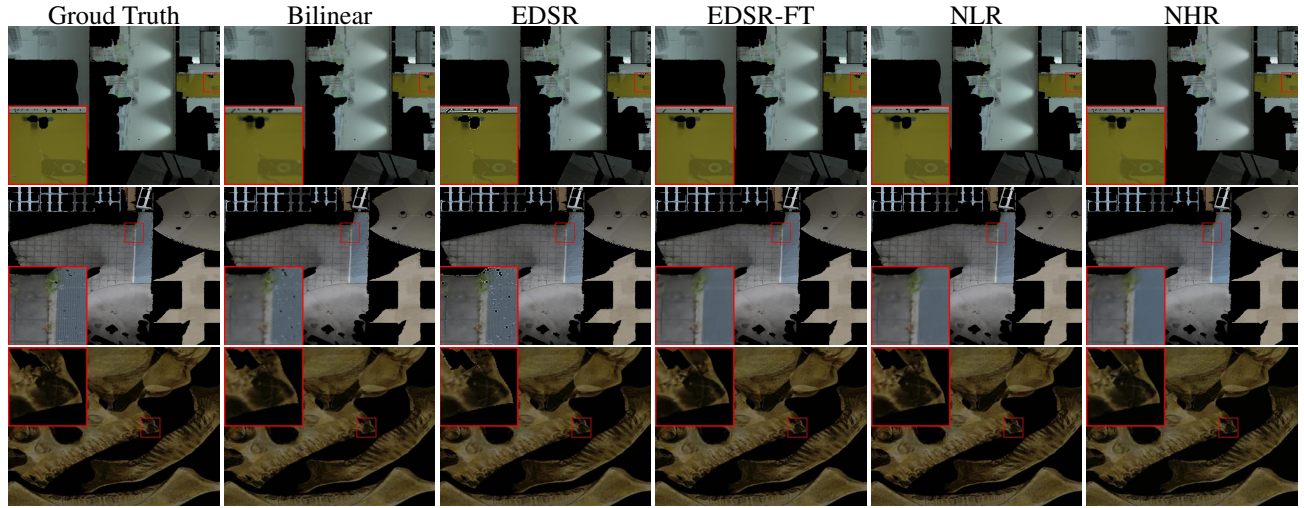
we also try cross-validation within the subset. That is, the training and testing texture maps are from the same subset. The 4 subsets are captured under different conditions and they may have different characteristics. In the case of cross validation within the subset, the training and testing data are from the same subset and they have the same characteristics. In the case of cross-validation on the whole dataset, there are more training data but with different characteristics. A comparison of these two cases can indicate whether subset characteristics or large training set is more important in our problem setting. The networks are trained for 50 epochs for subset cross-validation and 100 epochs for all of the other experiments.

## 6. Results

Using our 3DASR dataset, we compare three main categories; interpolation-based, model-based and learning-based methods for super-resolving the appearance of 3D objects. The interpolation-based methods include nearest, bilinear, bicubic, and Lanczos [16] interpolation. We use the method of Tsiminaki *et al.* [45] as a representative of the model-based category, denoted as HRST. Using the EDSR

Table 1: The PSNR results of different methods for scaling factor  $\times 2$ ,  $\times 3$ , and  $\times 4$ .

Method	ETH3D			Collection			MiddleBury			SyB3R			Average		
	$\times 2$	$\times 3$	$\times 4$	$\times 2$	$\times 3$	$\times 4$	$\times 2$	$\times 3$	$\times 4$	$\times 2$	$\times 3$	$\times 4$	$\times 2$	$\times 3$	$\times 4$
Nearest	19.06	16.71	14.68	24.22	19.7	16.92	10.08	7.93	7.08	30.84	27.88	25.82	21.07	18.12	16.0
Bilinear	20.61	18.24	16.32	26.2	21.48	18.84	11.87	8.88	7.77	31.75	28.83	26.9	22.67	19.6	17.56
Bicubic	20.21	17.96	15.88	25.67	21.12	18.29	11.32	8.81	7.73	31.77	28.78	26.73	22.28	19.34	17.16
Lanczos	20.01	17.74	15.69	25.42	20.86	18.07	11.14	8.81	7.81	31.71	28.7	26.63	22.09	19.15	17.0
HRST	16.18	—	16.12	32.29	—	29.63	22.13	—	20.88	27.9	—	26.34	22.17	—	21.17
HRST-CNN	—	—	—	32.24	—	29.9	22.76	—	21.55	—	—	—	—	—	—
EDSR	16.75	14.08	12.03	21.77	17.2	14.24	8.49	7.13	6.61	29.31	26.18	23.81	18.89	15.79	13.61
EDSR-FT	21.13	19.75	18.44	28.25	25.53	24.19	12.73	11.21	9.9	32.78	29.9	28.31	23.66	21.75	20.4
NLR-Sub	21.21	20.11	19.2	28.08	25.0	23.27	14.68	12.37	11.11	32.18	28.84	26.64	23.75	21.78	20.47
NLR	21.31	20.27	19.18	28.38	25.85	24.84	13.67	12.92	12.29	32.57	29.57	27.67	23.85	22.22	21.08
NHR	25.19	23.95	22.7	30.25	28.41	26.27	17.16	17.21	15.63	30.57	27.42	24.39	26.46	24.94	23.22

Figure 9: The visual results of *pipes*, *terrace*, and *Skull* for scaling factor  $\times 2$ .

network as a base model, we introduce several modifications of it. There are in total 6 different cases. EDSR: We use the pretrained network EDSR and directly test it on our data. EDSR-FT: We fine-tune the pretrained EDSR on our 3DASR dataset without architecture modification and using whole set cross-validation. NLR-Sub: We incorporate LR normal map into EDSR and use subset cross-validation. NLR: We incorporate LR normal map into EDSR as in Fig. 7a and use whole set cross-validation. NHR: We incorporate HR normal map into EDSR as in Fig. 7b and use whole set cross-validation. HRST-CNN: We use EDSR as a post-processing step of the super-resolved texture maps of HRST. In this scenario, the upsampling layer of EDSR is replaced with ordinary convolutional layers.

### 6.1. Objective metrics

We compute PSNR metrics in the active regions of the texture domains, that is, on the set of texels in the texture domain that is actually mapped to the 3D model. For the purpose of benchmarking, these metrics can also be computed in the image domain by reprojecting the texture maps into the image space. According to the PSNR values of Table 1, we can draw the following conclusions.

**Interpolation based methods.** Among the interpolation-based methods, bilinear interpolation achieves better results than bicubic and Lanczos interpolation, which contradicts the 2D image interpolation. This can be probably explained by the fact that the texture and the ordinary image domains have different characteristics. In the 2D image SR, LR image is modeled as bicubic down-sampled version of the HR image, which favors advanced interpolation methods. In the multi-view setting, due to the several sources of variability, the LR and HR texture maps might be not strictly aligned.

**Fine-tuning learning-based methods.** The texture domain knowledge is different than the image domain. The fine-tuning of EDSR-FT incorporates the characteristics of the texture compare to the pretrained EDSR model. Thus, algorithms need to be adapted to the specific domain.

**LR vs. HR normal maps.** We incorporate the 3D geometric information of the multi-view setting through the normal maps and we compare to the simple case of fine-tuned EDSR-FT. According to the PSNR values, the geometric information improves the quality of the reconstructed texture maps. We then validate its importance by comparing the two cases of NLR and NHR. The PSNR values increase even more when we express this geometric information with higher precision. NHR case, where HR



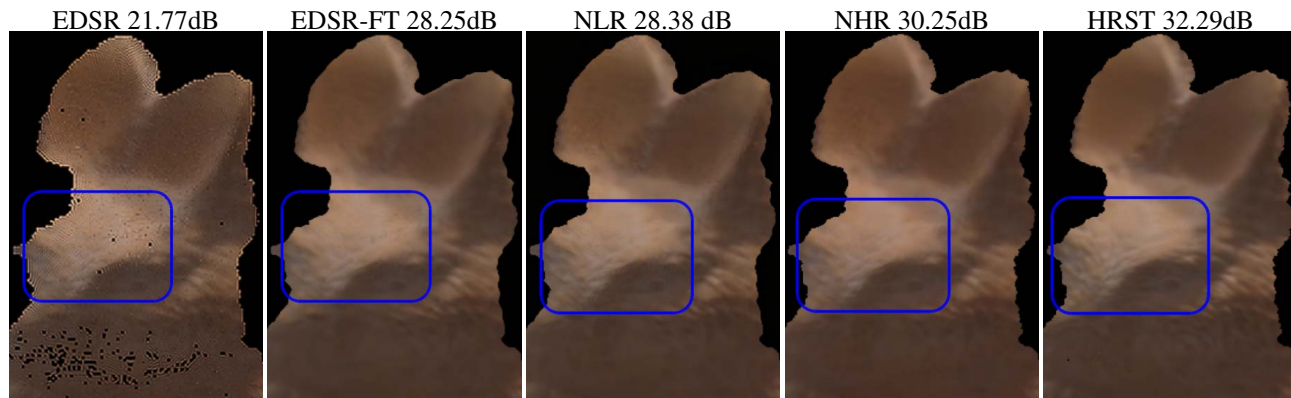


Figure 10: PSNR (dB) and close-ups of the super-resolved texture map of *Bunny* for scaling factor  $\times 2$ . While adding gradually the characteristics of the domain more details are recovered. The NHR achieves the highest PSNR value among the deep learning-based approaches while it stays still below the model-based HRST. Note that with NHR, texture SR is a **feed forward step**, while with HRST is an **iterative approach**.

normal maps are used outperforms NLR. Thus, HR normal maps capture more geometric details and improve the performance.

**Subset characteristics vs. training data size.** NLR-Sub uses cross-validation on the subset while NLR on the whole set. In the case of NLR-Sub, the subset characteristics are respected while in the case of NLR not. The main advantage of NLR is that more data are used for training (12 HR texture maps). The high PSNR values of the NLR compared to NLR-Sub indicates that the training data size is more important than subset characteristics to this task. Furthermore, the PSNR gap between NLR and NLR-Sub on ETH3D is larger than that on MiddleBury and Collection. This is because ETH3D is a relatively larger dataset than MiddleBury and Collection. Thus, even if subset cross-validation is used, NLR-Sub does not diverge a lot from NLR on ETH3D dataset. Therefore, we conclude that although each subset may have its own characteristics, training data size stands out as a major factor.

**Model based vs. learning based methods.** The model-based method HRST formulates the texture retrieval problem as an optimization problem. It is a two-stage iterative algorithm and its computational cost increases even more with an increase of geometric complexity. This explains the unstable behaviour of HRST method across the datasets. HRST outperforms NHR on MiddleBury and Collection whereas on ETH3D and SyB3R not. In most of the cases, HRST-CNN enhances the super-resolved texture maps. It is important to note that even in the cases where the model-based method outperforms the deep learning-based approach, the PSNR values are relatively close. More importantly, the deep learning-based approach is a feed-forward step that can be executed in seconds while the model-based is a heavy iterative process.

## 6.2. Visual results

The visual results are shown in Fig. 9 and Fig. 10. Directly upsampling the LR texture maps creates blurring images. EDSR leads to some white texels along the boundaries between the black region and the texture region. While we introduce gradually the characteristics of the domain through the EDSR-FT, NLR, and NHR methods, we successfully recover more visual details.

## 7. Conclusion

We provided 3DASR, a 3D appearance SR dataset<sup>1</sup> that captures both synthetic and real scenes with a large variety of texture characteristics. It is based on four datasets, ETH3D, Collection, MiddleBury, and SyB3R. The dataset contains ground truth HR texture maps and LR texture maps of scaling factors  $\times 2$ ,  $\times 3$ , and  $\times 4$ . The 3D mesh, multi-view images, projection matrices, and normal maps are also provided. We introduced a deep learning-based SR framework in the multi-view setting. We showed that 2D deep learning-based SR techniques can successfully be adapted to the new texture domain by introducing the geometric information via normal maps and achieve relatively similar performance to the model-based methods. This work opens up a novel direction of deep learning-based texture SR methods for the multi-view setting. A necessary next step is to enlarge our dataset either through common augmentation techniques or by following our proposed texture retrieval pipeline to introduce new datasets. The fact that the performance of our deep learning-based SR framework is in some cases (MiddleBury and Collection) below the model-based one indicates that there is still space for more elaborate methods that unify the concepts of model-based SR techniques and the 2D deep learning-based approaches.

<sup>1</sup>The dataset, the evaluation codes, and the baseline models is available at [https://github.com/ofsoundof/3D\\_Appearance\\_SR](https://github.com/ofsoundof/3D_Appearance_SR).



## References

- [1] E. Agustsson and R. Timofte. NTIRE 2017 challenge on single image super-resolution: Dataset and study. In *Proc. CVPRW*, July 2017. 3
- [2] J. Allebach and P. W. Wong. Edge-directed interpolation. In *Proc. ICIP*, volume 3, pages 707–710, 1996. 2
- [3] C. Allène, J.-P. Pons, and R. Keriven. Seamless image-based texture atlases using multi-band blending. In *Proc. ICPR*, pages 1–4, 2008. 2
- [4] L. Balmelli, G. Taubin, and F. Bernardini. Space-optimized texture maps. In *Computer Graphics Forum*, volume 21, pages 411–420, 2002. 3
- [5] F. Bernardini, J. Mittleman, H. Rushmeier, C. Silva, and G. Taubin. The ball-pivoting algorithm for surface reconstruction. *IEEE TVCG*, 5(4):349–359, 1999. 4
- [6] M. Bevilacqua, A. Roumy, C. Guillemot, and M. L. Alberi-Morel. Low-complexity single-image super-resolution based on nonnegative neighbor embedding. In *Proc. BMVC*, 2012. 3
- [7] S. Bi, N. K. Kalantari, and R. Ramamoorthi. Patch-based optimization for image-based texture mapping. *ACM Trans. Graph.*, 36(4), 2017. 3
- [8] A. Buades, B. Coll, and J.-M. Morel. A non-local algorithm for image denoising. In *Proc. CVPR*, volume 2, pages 60–65. IEEE, 2005. 2
- [9] C. Buehler, M. Bosse, L. McMillan, S. Gortler, and M. Cohen. Unstructured lumigraph rendering. In *Proc. SIGGRAPH*, pages 425–432, 2001. 2
- [10] Z. Chen, J. Zhou, Y. Chen, and G. Wang. 3d texture mapping in multi-view reconstruction. In *Proc. ISVC*, pages 359–371, 2012. 2
- [11] M. Corsini, P. Cignoni, and R. Scopigno. Efficient and flexible sampling with blue noise properties of triangular meshes. *IEEE TVCG*, 18(6):914–924, 2012. 4
- [12] P. E. Debevec, C. J. Taylor, and J. Malik. Modeling and rendering architecture from photographs: A hybrid geometry- and image-based approach. In *Proc. SIGGRAPH*, pages 11–20, 1996. 2
- [13] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proc. CVPR*, pages 248–255, 2009. 3
- [14] C. Dong, C. C. Loy, K. He, and X. Tang. Learning a deep convolutional network for image super-resolution. In *Proc. ECCV*, pages 184–199, 2014. 2, 3
- [15] C. Dong, C. C. Loy, K. He, and X. Tang. Image super-resolution using deep convolutional networks. *IEEE PAMI*, 38(2):295–307, 2016. 3
- [16] C. E. Duchon. Lanczos filtering in one and two dimensions. *Journal of Applied Meteorology*, 18(8):1016–1022, 1979. 6
- [17] S. Farsiu, M. D. Robinson, M. Elad, and P. Milanfar. Fast and robust multiframe super resolution. *IEEE TIP*, 13(10):1327–1344, 2004. 2
- [18] W. T. Freeman, T. R. Jones, and E. C. Pasztor. Example-based super-resolution. *IEEE Computer Graphics and Applications*, 22(2):56–65, 2002. 1, 2
- [19] Z. Fu, Y. Li, Y. Li, L. Ding, and K. Long. Frequency domain based super-resolution method for mixed-resolution multi-view images. *JSEE*, 27(6):1303–1314, 2016. 2
- [20] R. Gal, Y. Wexler, E. Ofek, H. Hoppe, and D. Cohen-Or. Seamless montage for texturing models. *Computer Graphics Forum*, 29(2):479–486, 2010. 2
- [21] B. Goldlücke, M. Aubry, K. Kolev, and D. Cremers. A super-resolution framework for high-accuracy multiview reconstruction. *IJCV*, 106(2):172–191, 2014. 1, 2, 3
- [22] B. Goldlücke and D. Cremers. Superresolution texture maps for multiview reconstruction. In *Proc. ICCV*, pages 1677–1684, 2009. 1, 3, 4
- [23] P. S. Heckbert. Survey of texture mapping. *IEEE Computer Graphics and Applications*, 6(11):56–67, 1986. 2
- [24] J.-B. Huang, A. Singh, and N. Ahuja. Single image super-resolution from transformed self-exemplars. In *Proc. CVPR*, pages 5197–5206, 2015. 3
- [25] R. Koch, M. Pollefeys, and L. Van Gool. Multi viewpoint stereo from uncalibrated video sequences. In *Proc. ECCV*, pages 55–71, 1998. 3
- [26] T. Köhler, M. Bätz, F. Naderi, A. Kaup, A. K. Maier, and C. Riess. Benchmarking super-resolution algorithms on real data. *arXiv preprint arXiv:1709.04881*, 2017. 3
- [27] Z. Lahner, D. Cremers, and T. Tung. Deepwrinkles: Accurate and realistic clothing modeling. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 667–684, 2018. 2
- [28] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. P. Aitken, A. Tejani, J. Totz, Z. Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *Proc. CVPR*, volume 2, page 4, 2017. 2
- [29] V. Lempitsky and D. Ivanov. Seamless mosaicing of image-based texture maps. In *Proc. CVPR*, pages 1–6, 2007. 2
- [30] B. Lévy, S. Petitjean, N. Ray, and J. Maillot. Least squares conformal maps for automatic texture atlas generation. *ACM Trans. Graph.*, 21(3):362–371, 2002. 4
- [31] A. Ley, R. Hänsch, and O. Hellwich. Syb3r: A realistic synthetic benchmark for 3d reconstruction from images. In *Proc. ECCV*, pages 236–251, 2016. 1, 2, 3, 5
- [32] X. Li and M. T. Orchard. New edge-directed interpolation. *IEEE TIP*, 10(10):1521–1527, 2001. 2
- [33] Y. Li, E. Eirikur Agustsson, S. Gu, R. Timofte, and L. Van Gool. Carn: Convolutional anchored regression network for fast and accurate single image super-resolution. In *Proc. ECCVW*, volume 4, 2018. 2
- [34] Y. Li, X. Li, and Z. Fu. Modified non-local means for super-resolution of hybrid videos. *CVIU*, 168:64–78, 2018. 2
- [35] Y. Li, X. Li, Z. Fu, and W. Zhong. Multiview video super-resolution via information extraction and merging. In *Proc. ACM MM*, pages 446–450, 2016. 3
- [36] B. Lim, S. Son, H. Kim, S. Nah, and K. M. Lee. Enhanced deep residual networks for single image super-resolution. In *Proc. CVPRW*, volume 1, page 4, 2017. 2, 5, 6
- [37] R. Maier, K. Kim, D. Cremers, J. Kautz, and M. Nießner. Intrinsic3D: High-quality 3D reconstruction by joint appearance and geometry optimization with spatially-varying lighting. In *Proc. ICCV*, volume 4, 2017. 4

- [38] D. Martin, C. Fowlkes, D. Tal, and J. Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *Proc. ICCV*, volume 2, pages 416–423, July 2001. 3
- [39] K. Nakamura, H. Saito, and S. Ozawa. Generation of 3d model with super resolved texture from image sequence. In *Proc. IEEE SMC*, volume 2, pages 1406–1411, 2000. 3
- [40] S. C. Park, M. K. Park, and M. G. Kang. Super-resolution image reconstruction: a technical overview. *IEEE Signal Processing Magazine*, 20(3):21–36, 2003. 2
- [41] J. L. Schönberger, E. Zheng, M. Pollefeys, and J.-M. Frahm. Pixelwise view selection for unstructured multi-view stereo. In *Proc. ECCV*, 2016. 4
- [42] T. Schöps, J. L. Schönberger, S. Galliani, T. Sattler, K. Schindler, M. Pollefeys, and A. Geiger. A multi-view stereo benchmark with high-resolution images and multi-camera videos. In *Proc. CVPR*, volume 3, 2017. 1, 2, 3
- [43] S. M. Seitz, B. Curless, J. Diebel, D. Scharstein, and R. Szeliski. A comparison and evaluation of multi-view stereo reconstruction algorithms. In *Proc. CVPR*, volume 1, pages 519–528, 2006. 3
- [44] R. Timofte, V. De Smet, and L. Van Gool. A+: Adjusted anchored neighborhood regression for fast super-resolution. In *Proc. ACCV*, pages 111–126. Springer, 2014. 2
- [45] V. Tsiminaki, J.-S. Franco, and E. Boyer. High resolution 3D shape texture from multiple videos. In *Proc. CVPR*, pages 1502–1509, 2014. 1, 3, 4, 6
- [46] M. Waechter, N. Moehrle, and M. Goesele. Let there be color! large-scale texturing of 3d reconstructions. In *Proc. ECCV*, pages 836–850, 2014. 2
- [47] C.-Y. Yang, C. Ma, and M.-H. Yang. Single-image super-resolution: A benchmark. In *Proc. ECCV*, pages 372–386, 2014. 3
- [48] J. Yang, J. Wright, T. S. Huang, and Y. Ma. Image super-resolution via sparse representation. *IEEE TIP*, 19(11):2861–2873, 2010. 2
- [49] R. Zeyde, M. Elad, and M. Protter. On single image scale-up using sparse-representations. In *Proc. Curves and Surfaces*, pages 711–730, 2010. 3
- [50] Y. Zhang, K. Li, K. Li, L. Wang, B. Zhong, and Y. Fu. Image super-resolution using very deep residual channel attention networks. In *Proc. ECCV*, pages 286–301, 2018. 2
- [51] Q. Zhou and V. Koltun. Color map optimization for 3d reconstruction with consumer depth cameras. *ACM Trans. Graph.*, 33(4):155:1–155:10, 2014. 1, 2, 3
- [52] Q.-Y. Zhou and V. Koltun. Color map optimization for 3d reconstruction with consumer depth cameras. *ACM Trans. Graph.*, 33(4):155, 2014. 3
- [53] M. Zollhöfer, A. Dai, M. Innmann, C. Wu, M. Stamminger, C. Theobalt, and M. Nießner. Shading-based Refinement on Volumetric Signed Distance Functions. *ACM Trans. Graph.*, 34(4):96, 2015. 1, 2, 3