# Adversarial Texture Optimization from RGB-D Scans

Jingwei Huang[1,3], Justus Thies[2], Angela Dai[2], Abhijit Kundu[3], Chiyu "Max" Jiang[3,4], Leonidas Guibas[1], Matthias Nießner[2], and Thomas Funkhouser[3]

[1]Stanford University    [2]Technical University of Munich    [3]Google Research    [4]UC Berkeley

**Input Image**    **Geometry**    **Our Reconstructed Textured Model**    **Zhou and Koltun**    **Ours**

Figure 1. Our goal is to reconstruct high-quality textures from an RGB-D scan. Unlike traditional methods which optimize for a parametric color map to reduce misalignment error (Zhou and Koltun [36]), we learn a misalignment-tolerant discriminator, producing sharper textures.

## Abstract

*Realistic color texture generation is an important step in RGB-D surface reconstruction, but remains challenging in practice due to inaccuracies in reconstructed geometry, misaligned camera poses, and view-dependent imaging artifacts. In this work, we present a novel approach for color texture generation using a conditional adversarial loss obtained from weakly-supervised views. Specifically, we propose an approach to produce photorealistic textures for approximate surfaces, even from misaligned images, by learning an objective function that is robust to these errors. The key idea of our approach is to learn a patch-based conditional discriminator which guides the texture optimization to be tolerant to misalignments. Our discriminator takes a synthesized view and a real image, and evaluates whether the synthesized one is realistic, under a broadened definition of realism. We train the discriminator by providing as 'real' examples pairs of input views and their misaligned versions – so that the learned adversarial loss will tolerate errors from the scans. Experiments on synthetic and real data under quantitative or qualitative evaluation demonstrate the advantage of our approach in comparison to state of the art (see Figure 1, right). Our code is publicly available[1] with video demonstration[2].*

---

[1] https://github.com/hjwdzh/AdversarialTexture
[2] https://youtu.be/52xlRn0ESek

## 1. Introduction

The wide availability of consumer range cameras has spurred extensive research in geometric reconstruction of real-world objects and scenes, with state-of-the-art 3D reconstruction approaches now providing robust camera tracking and 3D surface reconstruction [22, 19, 33, 9]. However, producing photorealistic models of real-world environments requires not only geometric reconstruction but also high-quality color texturing. Unfortunately, due to noisy input data, poorly estimated surface geometry, misaligned camera poses, unmodeled optical distortions, and view-dependent lighting effects, aggregating multiple real-world images into high-quality, realistic surface textures is still a challenging problem. In order to overcome these problems, various approaches have been developed to optimize color textures using models to adjust camera poses [36, 15], distort images [4, 15, 36], and balance colors [15, 36]. However, these prior approaches are not expressive enough and/or their optimization algorithms are not robust enough to handle the complex distortions and misalignments commonly found in scans with commodity cameras – and therefore they fail to produce high-quality results for typical scans, as shown in the results from Zhou and Koltun [36] in Figure 1.

To address these issues, we propose a flexible texture optimization framework based on a learned metric that is robust to common scanning errors (right side of Figure 1).
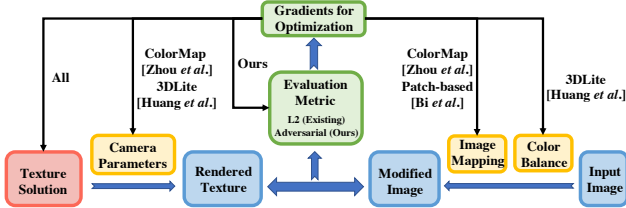
1

Figure 2. All methods target at optimizing a texture solution. Existing methods optimize the texture jointly with camera parameters [36, 15], image mapping [36, 4] or color balance [15]. Instead, we jointly solve texture with an adversarial evaluation metric to tolerate the errors.

The key idea behind our approach is to account for misalignments in a *learned objective function* of the texture optimization. Rather that using a traditional object function, like $L1$ or $L2$, we learn a new objective function (adversarial loss) that is robust to the types of misalignment present in the input data. This novel approach eliminates the need for hand-crafted parametric models for fixing the camera parameters [36, 15], image mapping [4, 36], or color balance [15] (bottom row of Figure 2) and replaces them all with a learned evaluation metric (green box in Figure 2). As such, it adapts to the input data.

Inspired by the success of adversarial networks in image synthesis [14], we propose to use a learned conditional discriminator to serve our *objective function*, and jointly optimize the color texture of a reconstructed surface with this discriminator. The condition is a captured image $I_A$ from the source view $V_A$, and the query is either (i) "real:" a second captured image $I_B$ (from an auxiliary view $V_B$) projected onto the surface and then rendered back to $V_A$, or (ii) "fake:" an image of the optimized synthetic texture rendered to view $V_A$. By optimizing the surface texture while jointly training this conditional discriminator, we aim to produce a texture that is indistinguishable from reprojections of captured images from all other views. During the optimization, the discriminator learns invariance to the misalignments and distortions present in the input dataset, while recognizing synthetic artifacts that do not appear in the real images (local blurs and seams). Therefore, the textures optimized to fool the discriminator (ours in Figure 1) appear more realistic than in previous approaches.

Our experiments show that this adversarial optimization framework produces notably improved performance compared to state-of-the-art methods, both quantitatively on synthetic data and qualitatively on real data. Moreover, since it tolerates gross misalignments, we are able to generate realistic textures on CAD models which have been only roughly aligned to 3D scans, in spite of large mismatches in surface geometry. This opens up the potential to produce CAD models with realistic textures for content creation.

## 2. Related Work

RGB-D scanning has a rich history, and many approaches have been proposed for color texture generation.

**View aggregation.** Common texture generation methods [19, 36] average projected input images to generate textures. To reduce blurriness artifacts, some approaches select a single or a few candidate views for each region [11]. Others formulate a multi-label selection energy minimization problem to minimize seam artifacts [21, 27, 30, 32, 15]. For instance, [15] aims at selecting the best view for each region to balance the visual sharpness and color consistency of boundaries between neighboring regions with different views selected, which is modeled as a multi-label graph-cut problem [5]. Our method does not explicitly define the aggregation method, but implicitly aggregates colors from different views based on a learned adversarial metric.

**Parametric color optimization.** Several approaches have been proposed to improve the mapping of input images to textures with parametric models, leveraging both human supervision [12, 23, 24, 34], as well as automatic optimization [3, 25]. Zhou *et al.* [36] propose to optimize a parametric model comprising camera poses and non-rigid grid deformations of input images to minimize an L2 color consistency metric. While these methods are able to fix small misalignments, their deformation models are often not expressive enough to handle many real-world distortions, particularly those due to largely approximate surface geometry. In contrast to a hand-crafted deformation model, we learn a distortion-tolerant adversarial loss.

**Patch-based color optimization.** Patch-based image synthesis strategies have been proposed for color texture optimization [4]. Rather than non-rigid image warping, they re-synthesize the input image with the nearest patch [26] to handle misalignments. However, general misalignment cannot be accurately modeled by translating patches, and the L2 loss is not robust to color, lighting or sharpness differences. Our method optimizes the discriminator to cover all these problems without requiring explicit re-synthesis.

**Neural textures.** Recently, neural rendering approaches have been proposed to synthesize a feature map on a surface that can be interpreted by a deep network to produce novel image views. For instance, [29] stores appearance information as high-dimensional features in a neural texture map associated with the coarse geometry proxy and decodes to color when projected to novel views. [28] stores the appearance information as high-dimensional features in volumes, and [2] uses features stored with points. These methods rely on the representation power of generative networks at rendering times to obtain novel viewpoints, which limits their applicability in standard graphics pipelines.
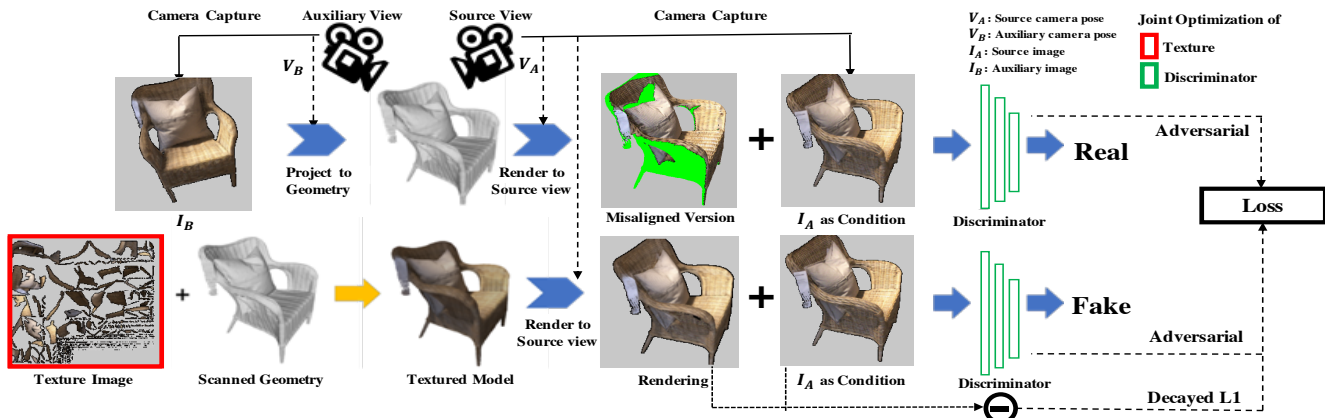
Figure 3. Texture Generation. From an input RGB-D scan, we optimize for both its texture image and a learned texture objective function characterized by a discriminator network. The discriminator operates on reprojections of input color images in order to maintain robustness to various misalignments. We randomly pick a pair of input images, *source* and *auxiliary*, and synthesize the fake and real examples from the source view, conditioned on the re-projected source image. The texture image and discriminator are trained in an alternating process.

## 3. Method

Our goal is to optimize for a color texture that can be used to render a scanned scene using a classical computer graphics pipeline. During the scanning procedure, we obtain color images and their estimated camera poses. These views, along with the reconstructed geometry, are input to our method. To optimize for a texture, we must specify an objective function; in this case, we must account for misalignments of the color images and the reconstructed model. Thus we propose to learn the loss function in conjunction with the texture (see Figure 3). The function is modeled as an adversarial loss using a discriminator network to identify 'real' and 'fake' imagery, and is designed to provide a misalignment-tolerant metric for our texture optimization.

### 3.1. Misalignment-Tolerant Metric

Our key insight is to propose to learn a conditional discriminator as a misalignment-tolerant metric adaptive to the error distribution of the input data. Figure 4(a) shows a 2D example where two observations (b) and (c) are misaligned by 2 units in the horizontal directions, and an L2 loss results in blurry appearance. Ultimately, we aim to synthesize a texture that appears as realistic as either observation. To achieve this, we ask the discriminator to consider both (b) and (c) as real conditioned on either observation. With such a discriminator, the blurred (d) results in a large loss and the texture will instead converge to either (b) or (c).

We extend this intuition to 3D where the geometry is observed from different viewpoints. We then aim to optimize a texture such that local patches of the texture rendered to various views look realistic. Therefore, conditioned on any arbitrary view, we generate real examples by a re-projection from any other view to this view, as shown in Figure 3. Such re-projection can be achieved by projecting the color image

onto the surface and then rendering back to another view. Unlike the simple 2D example, it is highly possible that there is no texture solution so that each local patch perfectly matches the one view from the input images, given camera and geometry error. However, the proposed approach is expected to push those inconsistencies to the smooth textured regions to hide any artifacts that can be easily identified by the discriminator, and thereby producing locally consistent realistic texture solution.

For each optimization iteration, we randomly select two input images, $I_A$ (source image) and $I_B$ (auxiliary image) with corresponding camera poses $V_A$ and $V_B$. The conditioning is $I_A$ from the viewpoint $V_A$, and the 'real' image is $I_B$ projected to the scan geometry and rendered from $V_A$, while the 'fake' image is the synthesized texture rendered from $V_A$. We alternating optimize the texture and discriminator. During texture optimization, we adjust the texture pixel colors to maximize the adversarial loss such that it looks more realistic under the discriminator scoring. During discriminator optimization, we minimize the adversarial loss such that it better classifies real and fake examples. We linearly combine adversarial loss with an L1 loss that decays exponentially as the optimization proceeds, which helps the optimizer find a good initial texture solution.

**Network Architecture** Our framework is adopted from the PatchGAN discriminator architecture proposed by Isola et al. [18]. We choose that framework because it is designed to produce local details that look as realistic as a given set of input images. We use three convolutional layers, resulting in a patch size of $70 \times 70$, which we find suitable for our input images of resolution $640 \times 480$. We apply a PatchGan to evaluate local $70 \times 70$ patches of images rather than the entire image. Patches are selected for discriminator training
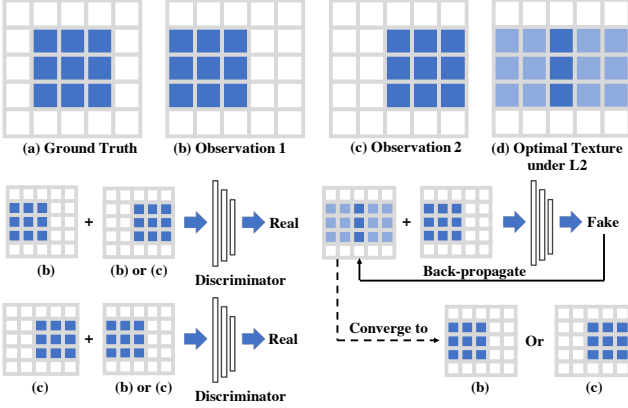
Figure 4. 2D example of a misalignment. (a) shows the ground truth pattern, which is observed with misalignment in (b) and (c); an L2 loss results in blurring (d). We train a discriminator which only accepts (b) and (c) as real examples conditioned on each other, and use it to optimize the texture, which converges to either (b) or (c).

if more than half of the patch is not occluded. Thus, patches used for training have sufficient overlap. Unlike the original, we remove all batch normalization layers and feed a single view example for each optimization iteration, which we empirically found to improve performance. Conditioned on the input view, we ask the discriminator to evaluate the residual of the synthesized example subtracted by the condition input. Finally, since we focus on evaluating foreground regions (pixels corresponding to input geometry), we remove the loss terms for regions where background comprises more than 90% of the receptive field.

## 3.2. Texture Optimization

To retrieve a texture, we jointly optimize the texture and the misalignment-tolerant metric. Inspired by the adversarial loss used in Pix2Pix [18], we express our view-conditioned adversarial loss as:

$$
\begin{aligned}
\mathcal{L}_c(T,D) = &\mathbb{E}_{x,y}(\log D(x,y))+ \\
&\mathbb{E}_{x,M_x}(\log(1-D(x,M_x(T)))),
\end{aligned} \tag{1}
$$

where $T$ and $D$ represent the target texture image and the discriminator parameters we are optimizing for. $x$ is the condition, a reprojected color image from the input sequence of captured images. $M_x$ is the fixed texture-to-image mapping given the camera pose associated with $x$. Here, a real example is an image $y$ re-projected to the view of $x$. We optimize $D$ with the objective to correctly identify real examples, misaligned real imagery, and fake examples rendered from the texture as $M_x(T)$. Simultaneously, we optimize the texture $T$ such that it is difficult to be identified as fake when mapped to view of $x$.

Since the adversarial loss alone can be difficult to train, we additionally add an L1 loss to the texture optimization

to provide initial guidance for the optimization:

$$
\mathcal{L}_{L1}(T) = \mathbb{E}_{x,y,M_x}||y-M_x(T)||_1. \tag{2}
$$

Our objective texture solution is:

$$
T^* = \arg\min_T \max_D \mathcal{L}_c(T,D) + \lambda \mathcal{L}_{L1}(T). \tag{3}
$$

During training, we initialize all pixels in texture image to zero and $\lambda = 10$. The high $\lambda$ allows the L1 loss to provide an initial texture, and for every 1000 steps we exponentially decay the lambda by a factor of 0.8. We optimize in alternating fashion for each optimization step, using the Adam optimizer for both the texture and discriminator with learning rates $10^{-3}$ and $10^{-4}$ respectively. For each object or scene, we optimize for 50000 steps to finalize our texture.

## 3.3. Differentiable Rendering and Projection

To enable the optimization of the RGB texture of a 3D model, we leverage a differentiable rendering to generate synthesized 'fake' views. We pre-compute a view-to-texture mapping using pyRender [17], and can then implement the rendering with a differentiable bilinear sampling.

To create the misaligned 'real' images ($I_B$ seen from $V_A$), we compute a reprojection; note that here we do not need to maintain gradient information. For each pixel $\mathbf{P}_A$ in the source image, we need to determine the corresponding pixel $\mathbf{P}_B$ in the auxiliary image, so that a bilinear sampling can be applied to warp image from the $V_B$ to $V_A$. Specifically, for $\mathbf{P}_A$ with depth value $d_A$ from the source depth map, we can determine its 3D location in the source view's space as $\mathbf{p_A} = d_A \mathbf{K}^{-1}\mathbf{P}_A$ where $\mathbf{K}$ is the intrinsic camera matrix. Suppose the transformations from the camera to the world space for the source and the auxiliary views are given as $\mathbf{T}_A$ and $\mathbf{T}_B$, the corresponding 3D and pixel location in the auxiliary view are $\mathbf{p}_B = \mathbf{T}_B^{-1}\mathbf{T}_A\mathbf{p_A}$ and $\mathbf{P}_B = \mathbf{K}\mathbf{p}_B$. The pixel is visible in the auxiliary view if $\mathbf{P}_B$ is in the scope of the image and the difference between z-dimension of $\mathbf{p}_B$ and $d_B$ from the auxiliary depth map is $< \theta_z$. We use $\theta_z = 0.1$ meters for scenes $\theta_z = 0.03$ for object level scanning.

## 4. Experiments

**Evaluation Metric**  For evaluation, we adopt several different metrics to measure the quality of the generated texture compared to the ground truth. First, we propose the nearest patch loss as an indicator of how close the patch appearance of the texture is to the ground truth. Specifically, for each pixel $\mathbf{u}$ we extract a $7 \times 7$ patch centered around it in the generated texture and find the L2 distance $d(\mathbf{u})$ between it and the nearest neighbor patch in the ground truth texture. We define the nearest patch loss as the average of all $d(\mathbf{u})$. Second, we adopt the perceptual metric [35] to evaluate perceptual quality. Finally, we propose to measure the difference between generated textures and ground
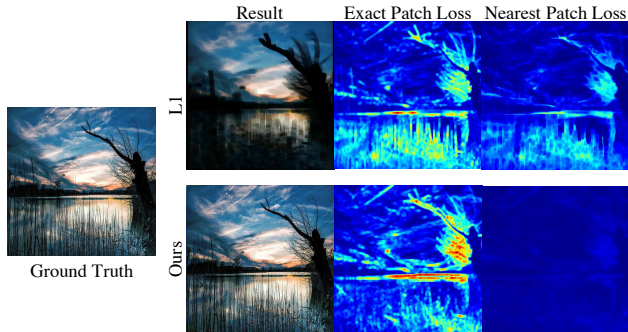
Figure 5. Texture Generation on 2D. The texture provided by our approach is visually closer to the ground truth image while avoiding blurring artifacts such as those introduced by an L1 loss. An exact patch loss favors alignment over perceptual similarity, while the nearest patch loss is a more robust metric.

| Camera / Geometry | Patch | Perceptual | Gradient | Sharp |
|---|---|---|---|---|
| L1 | 6.01 / 6.65 | 0.256 / 0.261 | 0.10 / 0.065 | 0.041 / 0.037 |
| ColorMap | 6.34 / 6.49 | 0.252 / 0.287 | 0.14 / 0.061 | 0.036 / 0.036 |
| Sharpest | 6.06 / 6.06 | 0.385 / 0.407 | 0.06 / 0.036 | 0.087 / 0.069 |
| 3DLite | 5.96 / 5.92 | 0.236 / 0.281 | 0.06 / 0.058 | 0.027 / 0.025 |
| VGG | 6.97 / 7.09 | 0.390 / 0.388 | 0.15 / 0.094 | 0.061 / 0.065 |
| **Ours** | **5.81 / 5.29** | **0.193 / 0.181** | **0.03 / 0.028** | **0.022 / 0.018** |

Table 1. Evaluation of different methods on our 3D synthetic dataset averaged across different levels of camera pose and geometry errors.

truth according to sharpness [31] and the average intensity of image gradients, in order to evaluate how robust the generated textures are to blurring artifacts without introducing noise artifacts. Note that standard image quality metrics such as the mean square error, PSNR [10] or SSIM [6] are ill-suited, as they assume perfect alignment between target and the ground truth [35].

**Synthetic 2D Example** We first verify the effectiveness of our method with a synthesized 2D example. We aim to optimize for a 2D image, given input observations with 2D micro-translation errors. We use an image resolution of $512 \times 512$ and translation error $\in [-16, 16]^2$. During texture optimization, we randomly select one observation as the source and another observation as the auxiliary, and optimize the target image to be more realistic under the current discriminator. Figure 5 shows the resulting image optimized with our approach in comparison to a naive L1 loss.

Visually, our optimized image is sharper and perceptually closer to the ground truth while an L1 loss results in blurry effect from aggregating multiple misaligned observations. In this simple setting, we evaluate the exact patch loss for each pixel quantitatively as the L2 distance of patches centered at this pixel between the generated image and the same one in the ground truth. The exact overall exact patch loss is the L2 norm of exact patch losses for all pixels. We additionally evaluate the nearest patch loss. Optimization with the L1 loss achieves 10.7 exact patch loss while ours is 11.3. However, we achieve 1.53 nearest patch loss, which is smaller than L1 as 7.33. This suggests that our method prefers realistic misalignment to blur. We successfully derive an image where every local patch is nearly identical to a misaligned version of the patch in the ground truth image.

**Synthetic 3D Example** In order to quantitatively evaluate our 3D texture generation, we create a synthetic dataset

of 16 models randomly selected from ShapeNet [7] across different categories. These shapes typically contain sharp edges and self-occlusion boundaries, complexities reflecting those of real-world objects. Since we aim to address arbitrary texturing, we enrich the appearance of these shapes by using 16 random color images from the internet as texture images. To create virtual scans of the objects, we uniformly sample $> 900$ views on a unit hemisphere by subdividing an icosahedron, from which we render the textured geometry as observed color images. To simulate misalignment, we associate each rendered image with a slightly perturbed camera pose, and to simulate geometry errors, we apply random perturbations to the geometric model. We use a set of errors increasing from $n = 1$ to $n = 4.5$, and refer to the supplemental material for additional detail regarding generating camera and geometry perturbations.

In Table 1, we study the effect of varying camera and geometry errors in this synthetic 3D setting. We report evaluation metrics for our approach as well as several state-of-the-art texture optimization methods, including methods based on an L1 loss and texturing using sharpest frame selection [31]. Our approach outperforms all other methods, as it avoids blurring effects often seen with L1 and ColorMap [36], and it avoids seams and over-sharpness introduced by methods relying on sharpness selection (3DLite [15] and sharpest frame selection). VGG [20] aggregates views by blending deep features, which is insufficient for handling misalignment artifacts. Two example scenes with increasing errors in camera and geometry are shown in Figure 6.

We additionally study the behavior of all methods in this experiment using the perceptual metric [35] in Figure 8. Although the performance drops for all methods with the increase of camera/geometry errors, our approach maintains the best perceptual quality as the errors increase. Figure 7 shows a qualitative comparison; our approach maintains a sharp result while ColorMap produces increasingly blurriness as the error increases.

**Alternative Discriminators?** We analyze the design choices for our misalignment-tolerant conditional discriminator in Figure 9. Removing the auxiliary view (b) and thus relying only on the source view to provide 'real' examples to the discriminator (similar to pix2pix [18]) renders the
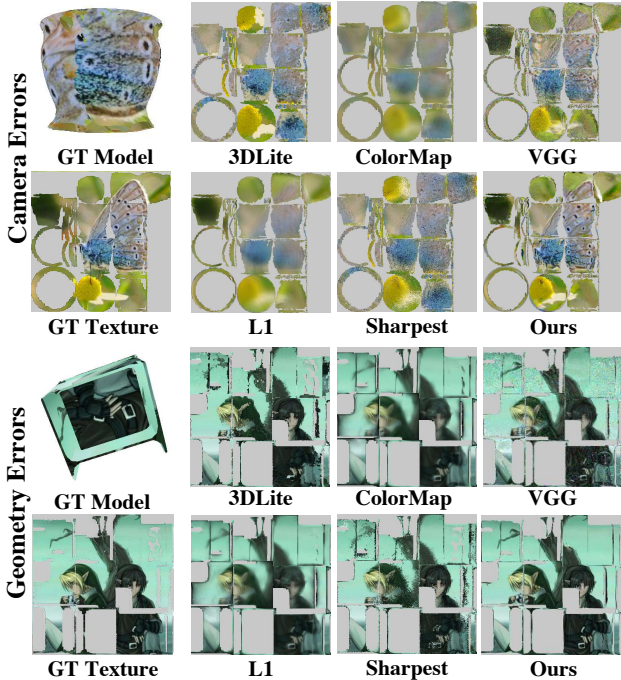
Figure 6. Texture generation in case of high camera or geometry errors. ColorMap [36] suffers from blurring, and Sharpest or 3DLite [15] selection lead to inconsistent boundaries or breaks structures. VGG [20] aggregates views by blending deep features with noises, which is not sufficient for handling misalignment artifacts. Ours is visually closest to the ground truth.
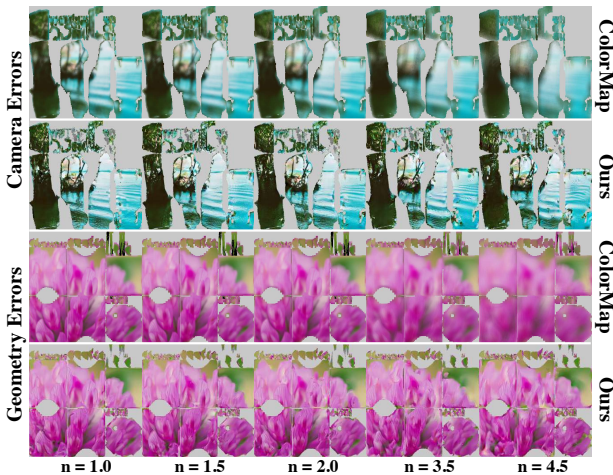


Figure 7. Texture generation under increasing camera or geometry errors. ColorMap [36] produces more blurry results under camera/geometry error while ours maintains sharp textures.

metric unable to handle misalignments. We also evaluate a general discriminator that classifies whether a generated patch is real or fake among entire input view sets without any condition (c), resulting in ambiguity as to where real patches come from. Our conditional discriminator leverag-
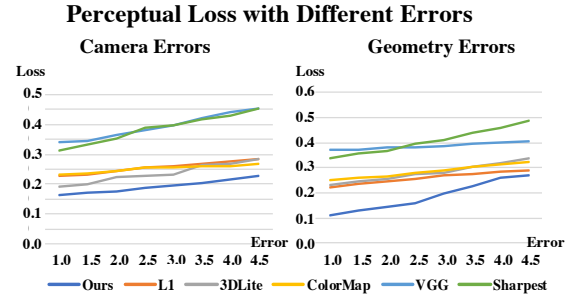


Figure 8. Perceptual loss under increasing camera or geometry errors; we outperform existing methods at various levels of error.
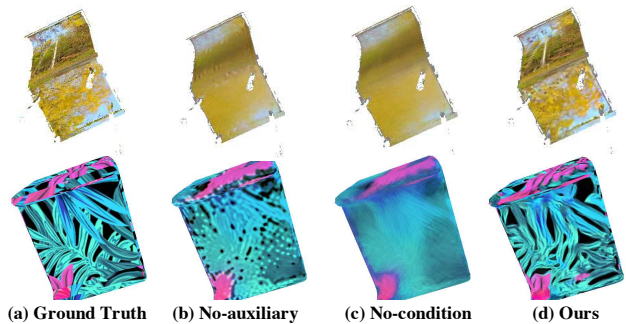


Figure 9. Comparing different discriminator options. (b) removes the auxiliary view from the discriminator, resulting in the lack of robustness to misalignments. (c) removes the condition from the discriminator, resulting in ambiguity in local regions. (d) our conditional discriminator leveraging auxiliary views to provide examples of realistic misalignments enables tolerance to misalignment and generation of textures reflecting input image characteristics.

ing reprojected auxiliary views enables robustness to misalignment, resulting in realistic texturing.

**Real Object Scans** We compare our method to state-of-the-art texturing methods on scanned objects from real environments. We use a structure sensor [1] along with its SLAM system to scan 35 chairs, producing scanned geometry, RGB-D frames and the camera poses ($\approx 500$ frames per scan). The foreground/background for the object in the RGB frames is determined by whether a ray intersects with the reconstructed geometry. Figure 11 (rows 1-4) shows qualitative comparisons. With an L1 loss or ColorMap [36], blur artifacts are induced by misalignment errors. Sharpest selection and 3DLite [15] use sharp region selection, resulting in seams and inconsistent global structures, as shown in the flower, leaf, and chair arms. A VGG loss [20] produces excess noise artifacts. Our approach produces sharp and consistent texturing, including detailed patterns such as the leaves in row 1 and woven structures in rows 2 and 3.

Additionally, we show a quantitative evaluation in Table 2 (first column) by evaluating the perceptual metric [35] for rendered textures against input observed views; our approach achieves the most realistic texturing.

|          | Object | ScanNet | CAD   |
|----------|--------|---------|-------|
| L1       | 0.197  | 0.470   | 0.199 |
| ColorMap | 0.186  | 0.461   | 0.234 |
| Sharpest | 0.222  | 0.510   | 0.260 |
| 3DLite   | 0.185  | 0.445   | 0.238 |
| VGG      | 0.272  | 0.534   | 0.289 |
| Ours     | **0.175** | **0.395** | **0.176** |

Table 2. Mean perceptual loss comparing the input images and rendered textures from different methods. Our method achieves best performance in the real and CAD datasets.

**Real Scene Scans**   To demonstrate the capability of our approach to optimize texture on a larger scale, we run our algorithm on the ScanNet dataset [8], which provides RGB-D sequences and reconstructed geometry of indoor scenes. We evaluate our approach on scenes with ID $\leq 20$ ($\approx 2000 - 3000$ frames per scan) and compare it with the existing state of the arts. Figure 11 (rows 5-9) and Table 2 (middle column) show qualitative and quantitative comparisons. Our method produces texturing most perceptually similar to the observed images; our misalignment-tolerant metrics aids in avoiding blur, increased sharpness, or excess noise produces by other methods due to camera and geometry errors in real-world scans.

**Real to CAD Models**   Since our method can better handle errors from approximate surface geometry, it is possible to consider texturing CAD models using real-world images to attain realistic appearances. While large datasets of 3D CAD models are now available [7], they are often untextured or textured simplistically, resulting in notably different appearance from real-world objects. To test whether our method can be applied in this challenging scenario, we use our collected dataset of real object scans, retrieve similar CAD models from ShapeNet manifold [16], and rigidly align them to the scanned objects. We then replace the scanned geometry with the CAD model and then use the captured color images and estimated poses from the scan to optimize the CAD texture. Qualitative and quantitative evaluation of our approach in comparison to existing state-of-the-art methods are show in Figure 11 (rows 10-13) and Table 2 (right column), respectively. Our approach is able to handle both camera poses errors as well as the synthetic-real geometry differences to produce texturing perceptually very similar to observed imagery, whereas other methods suffer strong blur, noise, and seam artifacts under these errors.

**Perceptual Quality**   Although we lack ground truth texturing for objects in real environments, we can compare the perceptual loss [35] of the rendered textured geometry from the corresponding viewpoint. We select 10 views uniformly distributed from the scanning video and render the textured model to compute the mean of the perceptual loss. Table 2 shows the performance of different methods on the
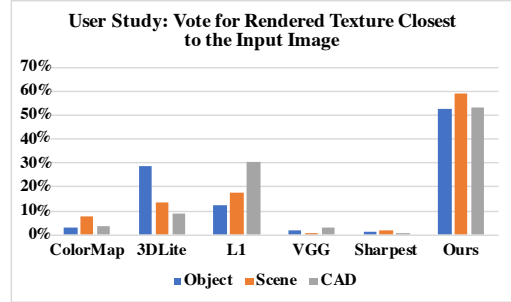


Figure 10. User study. We ask people to vote for the rendered textures from different methods that look closest to the input image.

object scans, scene scans and the CAD models; our method achieves the best performance in these three scenarios.

Additionally, we perform a user study to evaluate the quality of the texture, shown in Figure 10. Our user study comprised 63 participants who were asked to vote for the texture which produced a rendering closest to the input image. For some views, it can sometimes be difficult for users to differentiate between different methods when regions are largely uniform in color. Nevertheless, our method is still notably preferred over other texturing approaches. We provide additional comparisons with [32] and [13] in supplemental C and describe the influence of sparse views for training discriminators in supplemental D.

**Runtime.**   On average, our released implementation takes 7.3 minutes per object and 33.4 minutes per scene on a single TITAN X GPU.

## 5. Conclusion

We have proposed a misalignment-tolerant metric for texture optimization of RGB-D scans, introducing a learned texturing objective function for maintaining robustness to misalignment errors in camera poses and geometry. We represent the learned function as a conditional discriminator trained with an adversarial loss where 'real' examples characterize various misalignment errors seen in the input data. This avoids explicit parametric modeling of scanning errors, and enables our optimization to produce texturing reflective of the realism. Our approach opens up the potential for texturing synthetic CAD models with real-world imagery. It also makes an important step towards creating digital content from real-world scans, towards democratized use, for instance in the context of AR and VR applications.

## Acknowledgements

Figure 11. Visual comparison on object scans, ScanNet [8] scans of scenes, and CAD models aligned with object scans. Due to misalignment errors in both camera pose and geometry, both L1 loss and ColorMap [36] produce blurry artifacts, sharpest selection and 3DLite [15] result in inconsistent regions or breaks in texture structure, and VGG [20] blends learned features resulting in structural artifacts and noise. Our misalignment-tolerant approach produces sharp and consistent textures.

# References

[1] Structure sensor. *https://structure.io*. 6

[2] Kara-Ali Aliev, Dmitry Ulyanov, and Victor Lempitsky. Neural point-based graphics. *arXiv preprint arXiv:1906.08240*, 2019. 2

[3] Fausto Bernardini, Ioana M. Martin, and Holly Rushmeier. High-quality texture reconstruction from multiple scans. *IEEE Transactions on Visualization and Computer Graphics*, 7(4):318–332, 2001. 2

[4] Sai Bi, Nima Khademi Kalantari, and Ravi Ramamoorthi. Patch-based optimization for image-based texture mapping. *ACM Trans. Graph.*, 36(4):106–1, 2017. 1, 2

[5] Yuri Boykov, Olga Veksler, and Ramin Zabih. Fast approximate energy minimization via graph cuts. *IEEE Transactions on pattern analysis and machine intelligence*, 23(11):1222–1239, 2001. 2

[6] Dominique Brunet, Edward R Vrscay, and Zhou Wang. On the mathematical properties of the structural similarity index. *IEEE Transactions on Image Processing*, 21(4):1488–1499, 2011. 5

[7] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015. 5, 7

[8] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5828–5839, 2017. 7, 8

[9] Angela Dai, Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Christian Theobalt. Bundlefusion: Real-time globally consistent 3d reconstruction using on-the-fly surface reintegration. *ACM Transactions on Graphics (TOG)*, 36(3):24, 2017. 1

[10] Johannes F De Boer, Barry Cense, B Hyle Park, Mark C Pierce, Guillermo J Tearney, and Brett E Bouma. Improved signal-to-noise ratio in spectral-domain compared with time-domain optical coherence tomography. *Optics letters*, 28(21):2067–2069, 2003. 5

[11] Arnaud Dessein, William AP Smith, Richard C Wilson, and Edwin R Hancock. Seamless texture stitching on a 3d mesh by poisson blending in patches. In *2014 IEEE International Conference on Image Processing (ICIP)*, pages 2031–2035. IEEE, 2014. 2

[12] Thomas Franken, Matteo Dellepiane, Fabio Ganovelli, Paolo Cignoni, Claudio Montani, and Roberto Scopigno. Minimizing user intervention in registering 2d images to 3d models. *The Visual Computer*, 21(8-10):619–628, 2005. 2

[13] Yanping Fu, Qingan Yan, Long Yang, Jie Liao, and Chunxia Xiao. Texture mapping for 3d reconstruction with rgb-d sensor. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4645–4653, 2018. 7, 12

[14] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014. 2

[15] Jingwei Huang, Angela Dai, Leonidas J Guibas, and Matthias Nießner. 3dlite: towards commodity 3d scanning for content creation. *ACM Trans. Graph.*, 36(6):203–1, 2017. 1, 2, 5, 6, 8

[16] Jingwei Huang, Hao Su, and Leonidas Guibas. Robust watertight manifold surface generation method for shapenet models. *arXiv preprint arXiv:1802.01698*, 2018. 7

[17] Jingwei Huang, Yichao Zhou, Thomas Funkhouser, and Leonidas J Guibas. Framenet: Learning local canonical frames of 3d surfaces from a single rgb image. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 8638–8647, 2019. 4

[18] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017. 3, 4, 6

[19] Shahram Izadi, David Kim, Otmar Hilliges, David Molyneaux, Richard Newcombe, Pushmeet Kohli, Jamie Shotton, Steve Hodges, Dustin Freeman, Andrew Davison, et al. Kinectfusion: real-time 3d reconstruction and interaction using a moving depth camera. In *Proceedings of the 24th annual ACM symposium on User interface software and technology*, pages 559–568. ACM, 2011. 1, 2

[20] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pages 694–711. Springer, 2016. 5, 6, 8

[21] Victor Lempitsky and Denis Ivanov. Seamless mosaicing of image-based texture maps. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–6. IEEE, 2007. 2

[22] Richard A Newcombe, Shahram Izadi, Otmar Hilliges, David Molyneaux, David Kim, Andrew J Davison, Pushmeet Kohi, Jamie Shotton, Steve Hodges, and Andrew Fitzgibbon. Kinectfusion: Real-time dense surface mapping and tracking. In *Mixed and augmented reality (ISMAR), 2011 10th IEEE international symposium on*, pages 127–136. IEEE, 2011. 1

[23] Eyal Ofek, Erez Shilat, Ari Rappoport, and Michael Werman. Multiresolution textures from image sequences. *IEEE Computer Graphics and Applications*, (2):18–29, 1997. 2

[24] Frédéric Pighin, Jamie Hecker, Dani Lischinski, Richard Szeliski, and David H Salesin. Synthesizing realistic facial expressions from photographs. In *ACM SIGGRAPH 2006 Courses*, page 19. ACM, 2006. 2

[25] Kari Pulli and Linda G Shapiro. Surface reconstruction and display from range and color data. *Graphical Models*, 62(3):165–201, 2000. 2

[26] Denis Simakov, Yaron Caspi, Eli Shechtman, and Michal Irani. Summarizing visual data using bidirectional similarity. In *2008 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2008. 2

[27] Sudipta N Sinha, Drew Steedly, Richard Szeliski, Maneesh Agrawala, and Marc Pollefeys. Interactive 3d architectural

modeling from unordered photo collections. In *ACM Transactions on Graphics (TOG)*, volume 27, page 159. ACM, 2008. 2

[28] Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhofer. Deepvoxels: Learning persistent 3d feature embeddings. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2437–2446, 2019. 2

[29] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: Image synthesis using neural textures. *arXiv preprint arXiv:1904.12356*, 2019. 2

[30] Luiz Velho and Jonas Sossai Jr. Projective texture atlas construction for 3d photography. *The Visual Computer*, 23(9-11):621–629, 2007. 2

[31] Cuong T Vu, Thien D Phan, and Damon M Chandler. S3: A spectral and spatial measure of local perceived sharpness in natural images. *IEEE transactions on image processing*, 21(3):934–945, 2011. 5

[32] Michael Waechter, Nils Moehrle, and Michael Goesele. Let there be color! large-scale texturing of 3d reconstructions. In *European Conference on Computer Vision*, pages 836–850. Springer, 2014. 2, 7, 12

[33] Thomas Whelan, Stefan Leutenegger, Renato F Salas-Moreno, Ben Glocker, and Andrew J Davison. Elasticfusion: Dense slam without a pose graph. *Proc. Robotics: Science and Systems, Rome, Italy*, 2015. 1

[34] Zexiang Xu, Sai Bi, Kalyan Sunkavalli, Sunil Hadap, Hao Su, and Ravi Ramamoorthi. Deep view synthesis from sparse photometric images. *ACM Transactions on Graphics (TOG)*, 38(4):76, 2019. 2

[35] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 586–595, 2018. 4, 5, 6, 7

[36] Qian-Yi Zhou and Vladlen Koltun. Color map optimization for 3d reconstruction with consumer depth cameras. *ACM Transactions on Graphics (TOG)*, 33(4):155, 2014. 1, 2, 5, 6, 8
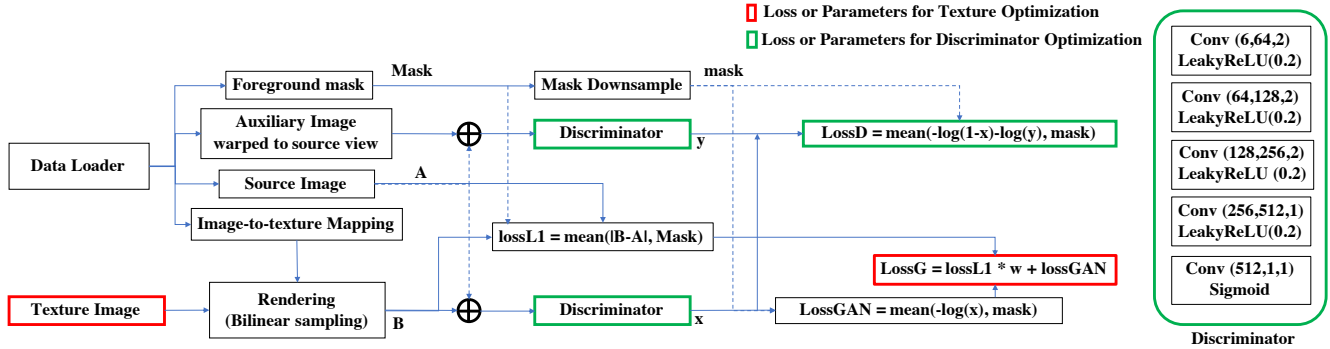
Figure 13. Detailed description of the optimization process.

# Supplemental

# A. Implementation details

## A.1. Data Loader During Optimization

Generating real and fake examples are two critical steps in our optimization process.

To achieve this, we provide three images for each input view as shown in Figure 12. For each view, we provide the color and the depth image from the device, as shown in the first two rows. The background pixels can be removed simply by deciding whether the corresponding rays are intersection the reconstructed mesh. Additionally, we provide the view-to-texture mapping for the differentiable rendering of the texture image.
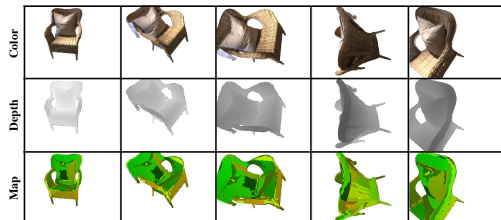


Figure 12. For each view of the scan, we cache the corresponding color and depth images. In settings where depth data is unavailable, we render depth maps from the target geometry. Additionally, we pre-compute the image-to-texture map for differentiable rendering.

Figure 13 shows the details of our optimization process. For each iteration, the data loader randomly selects a pair of related views as the source and the auxiliary, and we feed the network with a foreground mask containing the 3D scan, the real example synthesized by reprojecting the auxiliary image to the source view, the source image and the image-to-texture mapping. In the texture optimization stage, we render the target texture to source view based on image-to-mapping with a differentiable bilinear sampling as B. It is combined with the source image A as condition, sent to the discriminator to derive the prediction $x$. $x$ is used to compute the adversarial loss as "lossGAN". We additionally compute the L1 loss between A and B as "L1", and linearly combine "lossGAN" and "L1" with an exponentially decayed $w$ as "lossG", which is the objective function for optimizing the texture image. For every 1000 steps, we exponentially decay the $w$ by a factor of 0.8. In the discriminator optimization stage, the real example is combined with source image A and sent to discriminator to derive the prediction $y$. $x$ from the fake example and $y$ from the real example are combined to compute the object adversarial loss "lossD" for discriminator optimization.

The discriminator architecture consists of 5 convolutional blocks (figure 13). Conv(x,y,z) represents a 4x4 convolution with padding as 1, input channel as $x$, output channel as $y$ and stride as $z$. Each convolution block is followed with a gate function where the first four are leakyReLUs and the last is sigmoid.

## A.2. Details for Synthetic 3D Data Generation

We studied the behavior of our approach given the inaccurate camera pose or geometry. Camera perturbation is achieved by adding uniformly distributed noises to each dimension of the translation ranging from $[-e_t, e_t]$ and rotation as euler angle ranging from $[-e_a, e_a]$. To simulate geometry errors, we randomly generate a scalar for each vertex following the uniform distribution ranging from $[-e_g, e_g]$. Then, we apply 3 steps of Laplacian smooth to the scalars. We move the vertices along their normal directions with the distance specified by these scalars. We compare our method with different approaches for all selected ShapeNet objects with different amount errors. For camera errors, we set $e_t = 0.01 * 1.5^n$ ($n \in \{1.5, 2, 2.5, 3, 3.5, 4, 4.5\}$) and $e_a = 5°$. For geometry errors, we set $e_g = 0.02 * 1.5^n$ with $n \in \{1.5, 2, 2.5, 3, 3.5, 4, 4.5\}$.
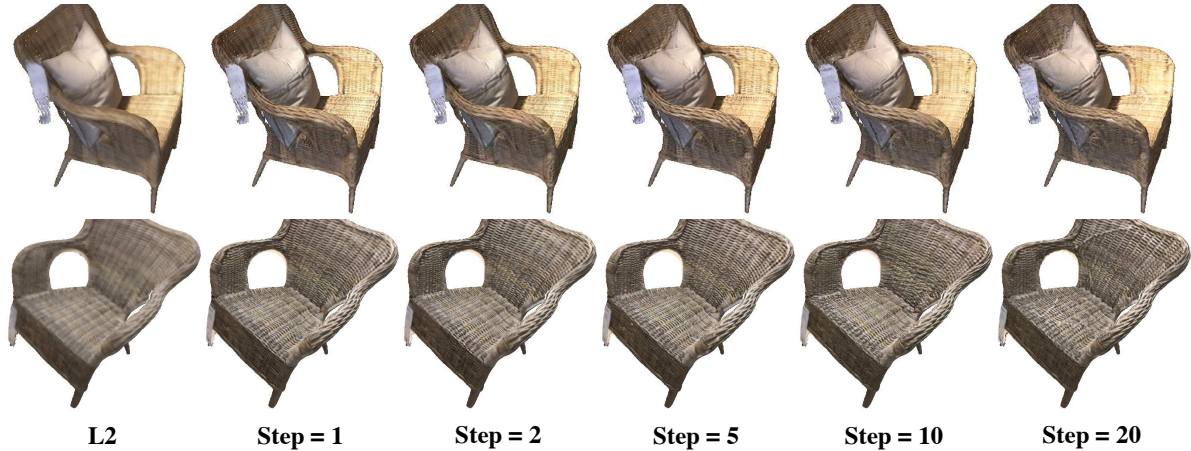
Figure 14. Ablation study on view density required for training. The two video sequences are with 426 and 244 frames respectively. The first column shows the average of frames under L2 loss. We sample the sequence by uniformly pick frames every k steps, with $k \in \{1, 2, 5, 10, 20\}$ in the remaining five columns.
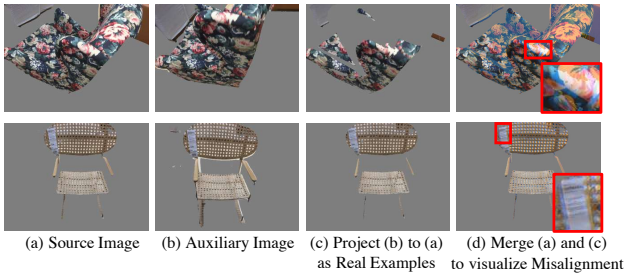


Figure 15. Real examples with Misalignment. (a) is the source image. (b) is an auxiliary image. In (c), we project (b) to the geometry and render to (a), which is a misaligned version of (a). (d) visualizes the misalignment by overlaying (a) and (c).

## B. Misalignment in the Data

Figure 15 shows an real example of the misaligned version of source views that we created.

## C. Sparsity of Views

We run our algorithm on scans with different number of frames to study the behavior and the robustness of our algorithm under different level of view sparsity, as shown in Figure 14. The two video sequences are with 426 and 244 frames respectively. The first column shows the average of frames under L2 loss. We sample the sequence by uniformly pick frames every k steps, with $k \in \{1, 2, 5, 10, 20\}$ in the remaining five columns. We notice that our algorithm produces appealing results if number of frames is larger than 25, but starts to show artifact patterns under this number. We believe the reason is that with very sparse views, there are not enough patches for learning a good misalignment tolerant metric. We believe this can be addressed by data augmentation with virtual camera perturbations.
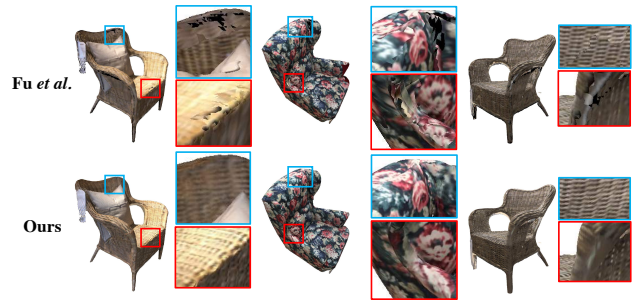


Figure 16. Comparison with [13] that develops based on [32].

## D. Additional Comparisons



Our method can deal with standard texture optimization dataset well as shown in the example of the fountain scan. Our experiments are aimed to showcase more challenging scenarios with complex scene geometry and lighting, as well as approximate surface reconstruction and alignment.

We provide additional comparisons between our method and Fu *et al.* [13] that develops based on [32]. View selection-based method yields inconsistent boundaries, while our method generates consistent texture.

## E. Additional Results

We provide a video called "supp/video.mp4" that contains the explanation of our approach and part of video and image results. We additionally provide the visualization of the full real dataset in our experiments with 200 renderings comparing to different methods for objects, scenes and CAD models. Please check "supp/*.html" for details.