

Multi-sensor large-scale dataset for multi-view 3D reconstruction

Oleg Voynov^{1,2} Gleb Bobrovskikh¹ Pavel Karpyshev¹ Saveliy Galochkin¹
 Andrei-Timotei Ardelean¹ Arseniy Bozhenko¹ Ekaterina Karmanova¹ Pavel Kopanov¹
 Yaroslav Labutin-Rymsho³ Ruslan Rakhimov¹ Aleksandr Safin¹ Valerii Serpiva¹
 Alexey Artemov^{4*} Evgeny Burnaev^{1,2} Dzmitry Tsetserukou¹ Denis Zorin⁵

¹Skolkovo Institute of Science and Technology ²Artificial Intelligence Research Institute
³Moscow Engineering Physics Institute ⁴Technical University of Munich ⁵New York University

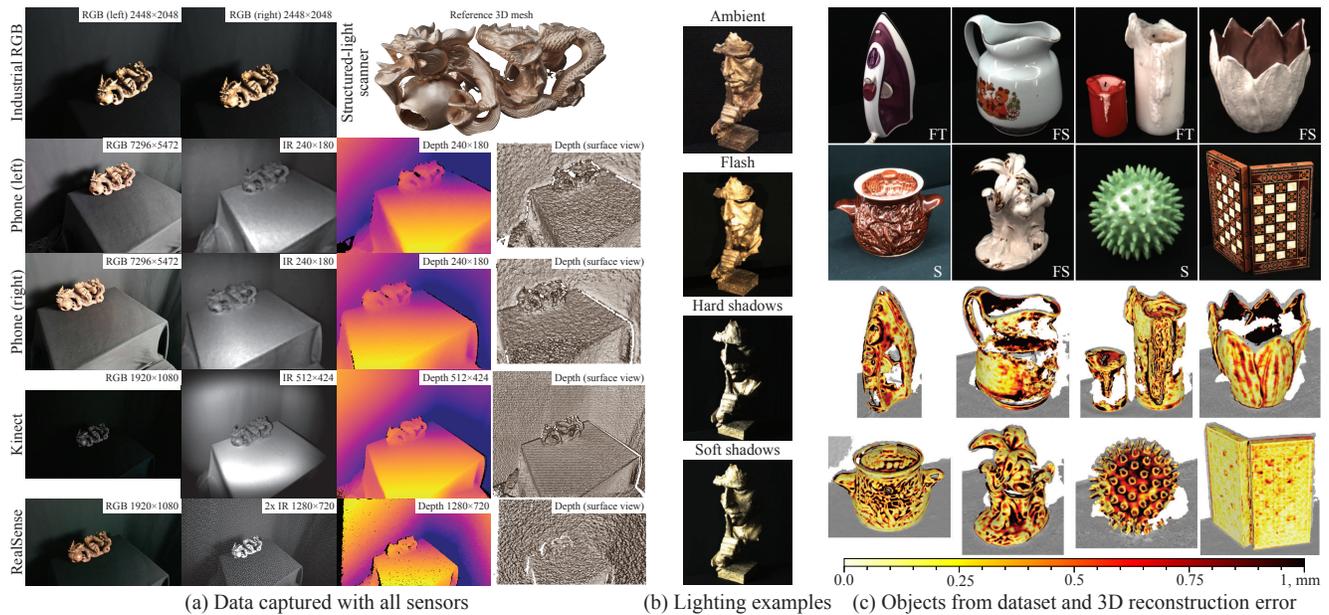


Figure 1. Our dataset contains RGB-D data captured (a) with 7 different devices, (b) under various lighting conditions. (c) We focus on materials challenging for 3D reconstruction algorithms, such as featureless (F), highly specular with sharp reflections (S), or translucent (T), as illustrated with reconstructions produced by state-of-the-art algorithms (compare with an “easy” object on the bottom right).

Abstract

We present a new multi-sensor dataset for multi-view 3D surface reconstruction. It includes registered RGB and depth data from sensors of different resolutions and modalities: smartphones, Intel RealSense, Microsoft Kinect, industrial cameras, and structured-light scanner. The scenes are selected to emphasize a diverse set of material properties challenging for existing algorithms. We provide around 1.4 million images of 107 different scenes acquired from 100 viewing directions under 14 lighting conditions. We expect our dataset will be useful for evaluation and training of 3D reconstruction algorithms and for related tasks. The dataset is available at skoltech3d.appliedai.tech.

*Work partially performed while at Skolkovo Institute of Science and Technology.

1. Introduction

Sensor data used in 3D reconstruction range from highly specialized and expensive CT, laser, and structured-light scanners to video from commodity cameras and depth sensors; computational 3D reconstruction methods are typically tailored to a specific type of sensors. Yet, even commodity hardware increasingly provides multi-sensor data: for example, many recent phones have multiple RGB cameras as well as lower resolution depth sensors. Using data from different sensors, RGB-D data in particular, has the potential to considerably improve the quality of 3D reconstruction. For example, multi-view stereo algorithms (e.g., [63, 96]) produce high-quality 3D geometry from RGB data, but may miss featureless surfaces; supplementing RGB images with depth sensor data makes it possible to have more complete reconstructions. Conversely, commodity depth sensors often

lack resolution provided by RGB cameras.

Learning-based techniques substantially simplify the challenging task of combining data from multiple sensors. However, learning methods require suitable data for training. Our dataset aims to complement existing ones (e.g., [13, 29, 33, 67, 72, 97]), as discussed in Section 2, most importantly, by providing multi-sensor data and high-accuracy ground truth for objects with challenging reflection properties.

The structure of our dataset is expected to benefit research on 3D reconstruction in several ways.

- *Multi-sensor data.* We provide data from *seven* different devices with high-quality alignment, including low-resolution depth data from commodity sensors (phones, Kinect, RealSense), high-resolution geometry data from a structured-light scanner, and RGB data at different resolutions and from different cameras. This enables supervised learning for reconstruction methods relying on different combinations of sensor data, in particular, increasingly common combination of high-resolution RGB with low-resolution depth data. In addition, multi-sensor data simplifies comparison of methods relying on different sensor types (RGB, depth, and RGB-D).
- *Lighting and pose variability.* We chose to focus on a setup with controlled (but variable) lighting and a fixed set of camera positions for all scenes, to enable high-quality alignment of data from multiple sensors, and systematic comparisons of algorithm sensitivity to various factors. We aimed to make the dataset large enough (1.39 M images of different modalities in total) to support training machine learning algorithms, and provide systematic variability in camera poses (100 per object), lighting (14 lighting setups, including “hard” and diffuse lighting, flash, and backlighting, as illustrated in Figure 1b), and reflection properties these algorithms need.
- *Object selection.* Among 107 objects in our dataset, we include primarily objects that may present challenges to existing algorithms mentioned above (see examples in Figure 1c); we made special effort to improve quality of 3D high-resolution structured-light data for these objects.

Our focus is on RGB and depth data for individual objects in laboratory setting, similar to the DTU dataset [29], rather than on complex scenes with natural lighting, as in Tanks and Temples [33] or ETH3D [67]. This provides means for systematic exploration and isolation of different factors contributing to strengths and weaknesses of different algorithms, and complements more holistic evaluation and training data provided by datasets with complex scenes.

2. Related work

Many datasets for tasks related to 3D reconstruction were developed (see, for example, a survey of datasets related to simultaneous localization and mapping (SLAM) [40]); we only discuss datasets most closely related to ours. A summary of comparisons to key datasets from previous work

Dataset	Sensor types	RGB, MPix	Depth, MPix	HR geo.	Poses /scene	Lighting	# Scenes	# Frames
DTU [29]	RGB (2) SLS	2	—	✓	49/64	8	80	27K
ETH3D [67]	RGB TLS	24	—	✓	10–70	U	24	11K
TnT [33]	RGB TLS	8	—	✓	150–300	U	21	148K
BlendedMVG [97]	unknown	3/0.4	—	—	20–1000	U	502	110K
BigBIRD [72]	RGB (5) RGB-D (5)	12	0.3	—	600	1	120	144K
ScanNet [13]	RGB-D	1,3	0,3	—	NA	U	1513	2.5M
Ours	RGB (2) RGB-D 1 (2) RGB-D 2 RGB-D 3 SLS	5 40	— 0.04 0.2 0.9	✓	100	14	107	877K

Table 1. **Comparison of our dataset to the most widely used related datasets.** U indicates uncontrolled lighting; frames are counted per sensor, i.e., all data from an RGB-D sensor are counted as a single frame. The number of separate images acquired may be considerably larger (1.4 M for our dataset). All scenes, from both training and testing sets, were counted.

is shown in Table 1.

RGB datasets with high-resolution 3D ground truth. A number of datasets are designed for multi-view stereo (MVS) methods, such as PatchMatch-based [4, 63, 93, 94], learning-based [9, 22, 36, 42, 80, 88, 96, 102], or hybrid methods [19, 35]. These datasets are also used for evaluation of methods reconstructing an implicit surface representation encoded by a neural network from a set of RGB images [47, 50, 82, 98], and in the novel view synthesis task [2, 44, 56, 58, 83].

Datasets from this category typically include high-resolution RGB, either photo [1, 29, 67, 97] or video [33], and high-resolution 3D ground truth obtained with a structured-light scanner (SLS) [1] or a terrestrial laser scanner (TLS) [33, 67]. The Middlebury datasets [59, 60, 68] focus on two-frame stereo, providing accurate ground truth for disparity in addition to RGB.

In this group, the DTU dataset [29] with controlled lighting and high-resolution ground truth is most often used for training learning-based MVS methods. Most other MVS datasets, while containing some images of isolated objects, focus on complete scenes, often collected with hand-held, freely positioned cameras [33, 67, 75].

Compared to previously developed datasets with high-resolution 3D scanner data, we provide the largest number of sensors, objects, and lighting conditions, and the most challenging objects.

Datasets with low-resolution depth. Datasets designed for tasks like SLAM, object recognition and segmentation are often collected using low-resolution depth sensors like Microsoft Kinect or Intel RealSense [7, 13, 46, 51, 72, 73]; some of these datasets combine high-resolution RGB with low-resolution depth (e.g., [13, 72]) but do not provide high-

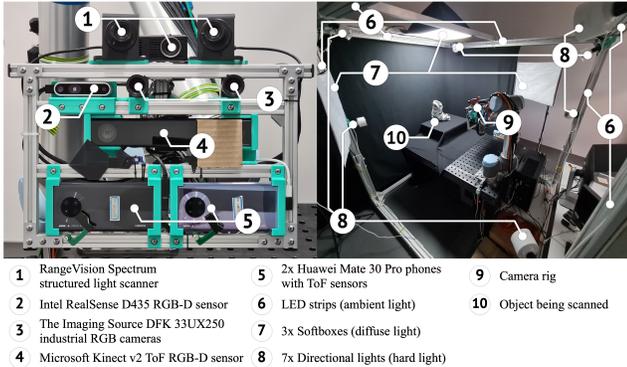


Figure 2. **Our acquisition setup** (view in zoom). We included a diverse set of seven commonly used RGB and RGB-D sensors, mounting them on a shared metal rig to aid data alignment. We constructed a metal frame surrounding the scanning area and installed various light sources to provide 14 lighting setups.

resolution ground truth for depth. A notable exception is CoRBS dataset [84], but it contains only four scenes. The degree of control over camera position and lighting in these datasets varies from complete [11, 72] to none (e.g., [51]).

These datasets are often used for qualitative evaluation of non-learning-based depth fusion methods [15, 28, 66, 89] which produce voxel- or surfel-based surface representations from depth maps; they are also used to train learning-based methods which produce voxel-based surface representation from RGB [45, 76], depth [14, 16], or RGB-D [17] data. These methods are likely to benefit from including our high-resolution depth data in training.

Synthetic datasets ShapeNet [8] and ModelNet [90] are often used to train learning-based depth fusion methods, e.g., [26, 43, 86, 87]. Synthetic ICL-NUIM SLAM benchmark [23] is used for evaluation, for example, in [15, 66, 89]. Such approach is limited by the differences between real-world and synthetic data.

Synthetic benchmarks [6, 34, 37] allow to generate large training sets easily via simulation of the acquisition process. However, real-world data is still required to faithfully model sensors, train generators, and test trained algorithms. Our dataset can be used for these purposes.

Datasets with multiple depth resolutions. Recently proposed RGB-D-D [24] and ARKitScenes [5] datasets make a step in a similar direction to ours, pairing low-resolution depth data acquired by smartphones with medium resolution (0.3 MPix) time-of-flight data and high-resolution laser scans, respectively.

Our dataset contains inputs from multiple depth sensors of low and high resolutions along with associated registered higher resolution RGB images, providing a framework for evaluating and training both depth fusion and RGB-D fusion algorithms previously trained on synthetic data, as well as developing new ones. Having multiple depth resolutions also supports applications such as depth map super-resolution [27, 74, 79, 91] and depth map completion [30, 103].

Device	#	RGB	Depth	IR	Intr.	Extr.	Rec.
DFK 33UX250	2	✓*	—	—	✓	✓	—
Mate 30 Pro	2	✓*	✓	✓	✓	✓	—
RealSense D435	1	✓*	✓*	✓✓*	✓	✓	—
Kinect v2	1	✓*	✓	✓	✓	✓	—
Spectrum	1	—	✓	—	—	—	✓

Table 2. **Composition of our dataset.** We provide RGB, depth, and IR images, intrinsic (Intr.) and extrinsic (Extr.) calibration parameters, and a reference mesh reconstruction (Rec.). The data marked with * is captured per lighting setup.

Similarly to our dataset, a concurrent work [71] complements RGB-D sequences from Intel RealSense with registered ground-truth captured by structured-light scanning. We provide registered depth images from common sensors at three levels of accuracy, and a controlled lighting variation.

3. Dataset

3.1. Overview

Our dataset consists of 107 scenes with a single everyday object or a small group of objects on a black background, see examples in Figure 1 and a complete list of scenes in the supplementary material.

We collected the dataset using multiple sensors mounted on Universal Robots UR10 robotic arm with 6 degrees of freedom and sub-millimeter position repeatability. We used the sensors shown in Figure 2 on the left: (1) RangeVision Spectrum structured-light scanner (SLS), (2) two The Imaging Source DFK 33UX250 industrial RGB cameras, (3) two Huawei Mate 30 Pro phones with time-of-flight (ToF) depth sensors, (4) Intel RealSense D435 stereo RGB-D camera, and (5) Microsoft Kinect v2 ToF RGB-D camera.

We surrounded the scanning area with a metal frame to which we attached the light sources, shown in Figure 2 on the right: seven directional lights, three diffuse soft-boxes, and LED strips which imitate ambient light. We also used flashlights on the phones as the light source moving with the camera. To prevent cross-talk between depth sensors we added external shutters that close the infrared (IR) projector of one sensor while the others are imaging.

For each scene, we moved the camera rig through 100 positions on a sphere with a radius of 70 cm around the object, using the same trajectory for all scenes, and collected the data using 14 lighting setups. For each device, except the SLS, we collected raw RGB, depth, and IR images, including both left and right IR images for RealSense. In total, we collected 15 raw images per scene, camera position, and lighting setup: 6 RGB, 5 IR, and 4 depth images, as illustrated in Table 2 and Figure 1a. As the IR and depth data from ToF sensors of the phones and Kinect is unaffected by the lighting conditions, we captured this data once per camera position. For the SLS we collected partial scans from 27 positions and combined them into a single scan, as we

describe further.

In our dataset, we included a large number of objects with challenging and varied surface material properties, as shown in Table 3. A set of qualitative labels corresponding to the key surface reflection parameters were assigned to each object based on visual estimation of these parameters. For example, *Specularity* represents the ratio of specular to diffuse reflectance for one of the dominant materials of the object, and *Reflection sharpness* characterizes how sharp the reflectance function peak is. These labels and their relation to performance of 3D reconstruction methods are discussed in the supplementary material.

3.2. Data acquisition and post-processing

Here, we summarize the main steps of our data acquisition, including selection of camera settings, camera calibration, and preparation of objects for scanning, and then describe our post-processing pipeline, which employs new methods for high-accuracy data registration and automatic generation of occlusion-valid reference depth images from SL scans.

The data acquisition procedure for each scene consisted of the following steps:

1. We automatically selected optimal camera exposure and gain settings for the scene.
2. We performed a preliminary SL scan, and if it was incomplete or low quality due to surface reflection properties, we applied a vanishing opaque matte coating to the object.
3. We scanned the object with the SLS from 27 viewpoints.
4. For coated objects, we accelerated the coating sublimation with hot air while keeping the object stationary. We then collected additional validation SL scans from five viewpoints to verify that the object was not deformed during coating removal.
5. We acquired RGB and low-resolution depth sensor data.

Sensor exposure/gain selection is critical for obtaining useful data: uniform settings result in frequent over- or underexposure due to variations of object surface properties. Hardware auto-exposure, being biased by the black background, proved to be inadequate for our setup and produced too high exposure/gain. Instead, we designed an auto-exposure algorithm inspired by [69, 70], which we used to find optimal camera settings for each scene, lighting, and sensor individually.

To find the *minimal noise* camera settings, we extracted the foreground mask of the scene from the images with and without the object using the method of [38], and then maximized the Shannon entropy of the foreground image w.r.t. the exposure value while keeping the gain at minimum. We then additionally maximized the entropy w.r.t. the gain value while keeping the found exposure value, to prevent the underexposure of very dark objects. To find the *real-time / high-noise* camera settings we swapped exposure and gain, optimizing the gain first while keeping the exposure at 30 FPS, and then optimizing the exposure while keeping the

Feature	Possible values	Dataset statistics
Specularity	diffuse/low/medium/high	40/15/35/24
Refl. sharpness	no refl./low/medium/high/very high	26/20/23/29/10
Translucency	none/low/medium/high/transparent	93/7/7/5/2
Mirror-like	yes/no	11/96
Metallic	yes/no	9/99
Texture type	color/3D/imperfections	9/24/43
Geom. features	small/medium/no dominant/flat	19/22/47/30

Table 3. **Surface material properties in our dataset.**

found gain.

We obtained *minimal noise* camera settings for each lighting variant, and *real-time* settings only for the dim flash and ambient lighting. We picked the settings for a single position of the rig per scene and used these settings for all positions. For RealSense, we also picked the optimal power of the IR projector out of 12 options, selecting the one leading to the lowest number of depth pixels missing a value.

Preparation of objects for 3D scanning.

For our dataset we picked the objects with surface material properties that challenge common sensors and reconstruction algorithms. However, these properties often challenge the SLS too, making it hard to obtain reliable high-resolution reference data. To get the highest-quality SL scans we applied a temporary coating to about half of the scanned objects (see Figure 3).

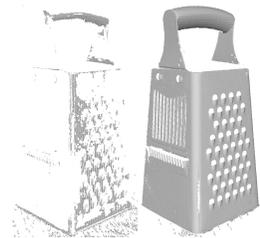


Figure 3. A partial scan without (left) and with coating (right).

Post-processing of structured-light scans. To simplify the use of the SL data for evaluation and training we merged raw partial scans into a single clean surface reconstruction. For this, we globally aligned partial scans using a variant of the iterative closest point algorithm, initialized with calibrated scanner poses. We assessed the alignment quality by comparing the inter-scan distance to the scanner resolution and in case of poor alignment, due to a lack of geometrical features on the object, re-scanned the object with attached 3D markers. From the aligned scans, we reconstructed a surface mesh using screened Poisson Surface Reconstruction [31] with the cell size of 0.3 mm, corresponding to a conservative estimate of the scanner resolution. Finally, we cleaned the surface keeping only the vertices close to the raw scans and manually removing the scanning and reconstruction artifacts and the supporting stand. We refer to the resulting surface as *the SL scan*.

For the temporary coating of the scanned objects we used Aesub Blue scanning spray. It sublimates from the surface at room temperature in a few hours, which we reduced to 5–15 minutes by slightly heating the object with a heat gun. To detect potential object deformations, we aligned the five validation scans to the SL scan and visually checked for

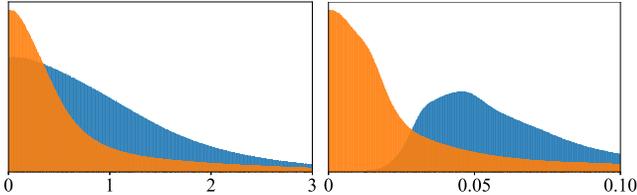


Figure 4. **Calibration refinement results.** Left: distribution of COLMAP reconstruction errors in mm before camera pose refinement (blue) and after (orange). Right: distribution of relative Kinect depth errors w.r.t. SLS before correction (blue) and after (orange).

distance variation indicating deformation. If a deformation was localized only to an isolated part of the object (*e.g.*, a power cord) we removed this part from the scan, otherwise, we excluded the whole object from the dataset.

Camera calibration. To measure the trajectories of the sensors and their intrinsic camera parameters we used the calibration pipeline of [65]. It is based on generic camera models known to result in a more accurate calibration than parametric camera models, commonly used in benchmarks related to ours, including DTU and ETH3D.

Application of this calibration pipeline to our camera rig as is was numerically unstable due to the large variability of sensor properties, such as field of view and resolution. To eliminate instability we split the calibration procedure into several steps, as we describe in the supplementary material.

For all RGB and depth sensors and the SLS, we obtained central generic camera models with several thousands parameters and the trajectory of the sensor in the global coordinate system. The mean calibration error for different sensors at different steps of the procedure was in the range 0.04–0.4 px, or approximately in the range 0.024–0.15 mm.

We observed thermal effects to contribute noticeably to calibration results, in particular, calibration errors drop by factors of 2–10 if the devices warm up after power-on to stable operating temperature. To reduce thermal drift in our data we pre-warmed all devices for 1 hour before calibration and at the beginning of each day of scanning.

Refinement of camera poses. Limited rig stiffness and variations in positioning of the robotic arm lead to a noticeable deviation of the camera trajectory during scanning from the calibrated trajectory, up to several pixels in image space. We additionally refined individual sensor poses per scene w.r.t. SL scan using a new approach inspired by [10, 18].

The idea of these works is to optimize a photometric discrepancy between a sensor image and a render of a reference texture on the scan to the image plane, with respect to the pose of the sensor. In [10], geometric features of the scan, such as normals, are used as the reference texture, while in [18], poses for multiple viewpoints are optimized simultaneously and for each viewpoint the projections from images from all other viewpoints are used as the reference. In both cases, the mutual information between intensity distributions of the sensor image and the render is used as the discrepancy

measure, and it is optimized using derivative-free NEWUOA method [55]. We used the same overall idea but changed both key components of the method: the reference texture and the discrepancy measure.

The modification of the original method was required as the previously proposed variants did not yield sufficiently accurate results for our data. The variant of [10] assumes the presence of photometric similarity between the render of the scan surface and the photo of the real surface, which is often not present, *e.g.*, for a smooth surface with color texture, and, importantly, relies on the match between the silhouettes of the rendered scan and the real object in the photo, which is often violated as not all parts of the object visible in the photo may be scanned. Simultaneous optimization of multiple viewpoints in [18] often results in pose drift.

Our approach is to use the raw SLS camera images as the reference texture, since the scan is constructed from these images and they are perfectly aligned. Instead of mutual information between intensity distributions, which we found to be often unstable, we optimized the smooth L_1 distance in deep CNN feature space. Inspired by [39], for feature extraction we used a model pre-trained specifically for feature matching [21].

For some camera poses and sensors, the direct alignment to the SLS images was unstable, as these images were captured under SLS illumination, which differs significantly from the illumination used for the other cameras. To further stabilize the alignment procedure for the whole dataset, we split it into three steps. First, we aligned, directly to the SLS reference texture, a single image from an industrial RGB camera, for a viewpoint selected for each scene individually to maximize the alignment tightness. Then, we aligned the images from the industrial RGB camera for the remaining viewpoints successively, using the already aligned images as the reference. Finally, for all the other sensors, we aligned each image individually to the image from the industrial RGB camera captured at the same position of the rig.

Using this sequential sensor alignment procedure we obtained subpixel alignment accuracy in most cases. In Figure 4 left, we show that refinement of camera poses reduces the error of a 3D reconstruction method.

Refinement of depth camera calibration. The lower quality depth sensors were calibrated by the manufacturers using calibration systems different from the one we used for the SLS and the camera trajectories. This leads to a misalignment between the lower quality depth data and the SL scan, which can be represented as a composition of a 3D rigid transform, scaling and a non-linear warping [25, 77, 100, 104].

To reduce the misalignment, we used a correction model $d_{\text{undist}} = S(u, v)S(d_{\text{raw}})$, which transforms the raw depth measurement d_{raw} at pixel (u, v) in the depth image to the depth value in the 3D space of the SL scan, using the cubic basis splines $S(u, v)$ and $S(d)$ defined on regular grids. We

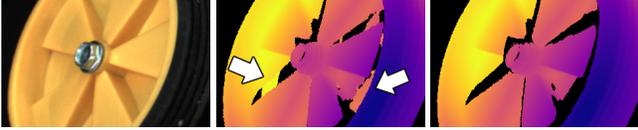


Figure 5. **Rendering SL depth with occluding surface.** Although the real surface of the object is solid (left), its SL scan is incomplete, so the background parts of the object bleed through in the rendered depth map (middle). Occluding surface helps to discard incorrect depth values (right).

trained the model on depth images of a calibration pattern, and as the reference used the calibration data from the same system that we used for camera calibration.

In Figure 4 right, we show that refinement of depth camera calibration reduces the error of a low quality sensor depth w.r.t. the SL scan.

Rendering SL scans. Occluding surface. In many tasks that we target, the ground-truth 3D data is needed in the form of depth maps. At the same time, the SL scanning often does not capture the complete object due to its geometric limitations, so a direct rendering of a depth map from the SL scan is likely to produce incorrect large depth values at some points because closer parts of the surface are absent from the scan, as illustrated in Figure 5 middle.

Previous works have addressed this problem by discarding depth values either inconsistent with an MVS reconstruction from RGB images [41], or occluded by a complete surface reconstructed from the scan and adjusted manually [67]. We follow the latter, more reliable approach, but propose a fully automatic method for building the occluding surface.

An *occluding surface* for an object is any surface fully enclosing the object. Our goal is to construct an occluding surface tight to the real surface of the object, using the information about visibility of the real surface from the viewpoint of the SLS camera during scanning. The depth values rendered from the SL scan that differ from the values rendered from such an occluding surface by more than a certain threshold can then be identified as spurious and removed from the generated depth maps.

To find the occluding surface we assume that the convex hull of all partial SL scans C encloses the part of the object that we are going to render. Conceptually, we find the occluding surface as $C \setminus \cup_i \text{Cone}(v_i, S_i)$, where the cone $\text{Cone}(v_i, S_i)$ encloses all the points between the surface of the partial scan S_i and the respective viewpoint of the SLS camera v_i . Intuitively, we carve the free space between the scanner and the object out of the convex hull, for all positions of the scanner.

The direct computation of such an occluding surface by a sequence of boolean operations on closed meshes is extremely expensive as the number of triangles for each mesh can be very large. Instead, we compute an accurate approximation by sampling all surfaces involved (we used the sampling distance of 0.1 mm) and keeping only the points

that are in the interior of the convex hull and in the exterior of all cones, which we check using depth testing. We then reconstruct the occluding surface from these samples using screened Poisson Surface Reconstruction.

The occluding surface obtained with this method can be used to generate accurate, occlusion-valid depth images from SL data as shown in Figure 5 right.

4. Experimental evaluation

We demonstrate the challenges of our dataset for RGB-based 3D reconstruction methods, and additionally, demonstrate one of the uses of our dataset as a benchmark to compare methods of 3D reconstruction from different modalities. For this, we tested 5 methods which reconstruct the surface only from RGB data, 3 methods which use only depth data, and one method which uses both modalities.

Methods. COLMAP [62, 63] is an RGB-to-3D reconstruction pipeline. It implements a non-learning-based multi-view stereo (MVS) method: from multi-view RGB images, the depth maps are estimated per-view and then fused into a point cloud representing the reconstructed 3D surface. COLMAP is commonly evaluated on benchmarks similar to ours, so we include it as a baseline. ACMP [94] is a non-learning, PatchMatch-based MVS method with planar prior, which shows a strong performance on benchmarks such as Tanks and Temples (TnT). VisMVSNet [102] and UniMVSNet [54] are learning, plane-sweeping-based MVS methods with top performance on TnT benchmark. NeuS [82] is a method which reconstructs the surface as a signed distance function (SDF) represented with a neural network which is fitted directly to RGB images via differentiable rendering. TSDF Fusion [12, 28] is a classical non-learning-based depth-fusion approach. It reconstructs a truncated SDF (TSDF) of a surface on a voxel grid via iterative integration of depth maps. RoutedFusion [86] is a learning-based extension, which performs depth map integration using neural networks. SurfelMeshing [66] is a non-learning-based depth-fusion method that reconstructs the surface using a surfel cloud as an intermediate representation. Neural RGB-D surface reconstruction [3] is a recent method which, similarly to NeuS, reconstructs the TSDF of a surface represented with a neural network, but in contrast to NeuS, fits the network to RGB and depth images.

Data. We tested the methods under the most favorable conditions available in our dataset. We obtained the reconstructions using full-resolution images from all 100 viewpoints for each scene, and the ground-truth camera poses and intrinsic camera models. For each data modality we used the highest-quality option: the photos from one of the industrial RGB cameras with the minimal noise settings, taken under ambient light to rule out the effects of lighting, and the depth maps from the ToF sensor of Kinect, of a higher resolution than of the phones, and more stable than RealSense.

For the learning-based methods we tested the models trained on prior datasets: DTU [29] and BlendedMVS [97] for MVSNet, and ModelNet [90] for RoutedFusion. We provide more details in the supplementary material.

Measures. We evaluated the reconstructions w.r.t. the full SL scans using Precision, Recall, and F-score quality measures, similarly to the prior benchmarks [33, 67]. Precision measures the reconstruction accuracy: we calculated it as the percentage of the reconstruction points closer to the reference than a certain distance threshold. Recall measures the reconstruction completeness: we calculated it as the percentage of the *reference* points closer to the reconstruction than a threshold. F-score is the harmonic mean of these two numbers, which, to be high, requires the reconstruction to be both accurate and complete.

For a careful calculation of Precision and Recall two problems have to be considered. First, both the reconstruction and the reference may have varying point densities, which will cause uneven contribution of different parts of the surface to the value of the measure. Second, the reference SL scans are incomplete, so the distance from some reconstruction points to the SL scan does not represent the distance to the real surface of the object, specifically, if the point lies near the missing part of the surface. We addressed these problems similarly to [67], extending the approach using our new occluding surface. Specifically, we calculated each measure in small cells of 3D space and then took the average as the final value, and used only the points for which the distance to the real surface of the object is certain. We describe this in more detail in the supplementary material.

Results. In Figure 6 we show the average performance on our dataset for the tested methods. Each curve shows the number of scenes that the respective method reconstructs with a better value of the measure than the value on the X axis. We tested Neural RGB-D surface reconstruction only on a fraction of scenes, due to its long running time, so we exclude the curve for this method.

In Figure 7 we show reconstructions for several scenes. The top part of the figure shows the results produced by methods of different types: an RGB-only UniMVSNet, an RGB-only NeuS with neural representation, Neural RGB-D surface reconstruction, and a depth-only TSDF Fusion. The middle part shows the best result per scene w.r.t. Recall. The bottom part shows the best result w.r.t. Precision.

In these figures we show the measures calculated with the threshold of 0.5 mm. In the supplementary material we show distributions of the measures for other values of the threshold, distributions of additional quality measures, and additional reconstruction visualizations.

Discussion. The best method w.r.t. Recall (ACMP) reconstructs all scenes in our dataset on at least 53%, with the distance threshold of 0.5 mm, but only half of the scenes on 80% or more. The best method w.r.t. Precision and the

overall quality represented with F-score (VisMVSNet) reconstructs all scenes with at least 32% accuracy and only 14 scenes with accuracy higher than 80%. This demonstrates that our dataset contains plenty of challenges for state-of-the-art 3D reconstruction methods. In particular, featureless parts of the surface, especially with sharp reflections, are often missing in the reconstruction, as illustrated in Figure 7. At the same time, VisMVSNet significantly outperforms UniMVSNet w.r.t. Precision and F-score and performs almost as well w.r.t. Recall. This is opposite to their relative performance on prior benchmarks (TnT and DTU), which indicates that our dataset poses a different set of challenges.

Comparison of the methods of different types shows that reconstructions from depth maps are significantly less accurate than reconstructions from RGB images, although sometimes may be more complete, as illustrated at the top of Figure 7. This demonstrates that these modalities can complement each other. Similarly, NeuS, which uses a neural surface representation, fills-in the areas of the surface challenging for RGB-based methods, but often inaccurately. Remarkably, Neural RGB-D surface reconstruction produces the surface of a significantly lower quality in comparison to VisMVSNet and NeuS, while using the same input RGB images and additionally the depth maps. This illustrates that there is a room for development of 3D reconstruction methods that effectively use both modalities, and we believe that our dataset will facilitate the development of such methods.

5. Conclusion

We presented a new dataset for evaluation and training of 3D reconstruction algorithms. Compared to prior datasets the distinguishing features of ours include a large number of sensors of different modalities and resolutions, depth sensors in particular, selection of scenes presenting difficulties for many existing algorithms, and high-quality reference data for these scenes. Our dataset can support training and evaluation of methods for many variations of 3D reconstruction tasks, in particular, learning-based 3D surface reconstruction from multi-view RGB-D data.

The main (intentional) limitation of our dataset is the use of the laboratory setting: the focus on static isolated objects with easy-to-separate background, the same camera trajectory for all scenes, the laboratory lighting. Another possible limitation is a small range of object sizes, limited by the physical size of the setup.

Acknowledgements. The authors acknowledge the use of Skoltech supercomputer Zhores [99] for obtaining the results presented in this paper. E. Burnaev and O. Voynov were supported by the Analytical center under the RF Government (subsidy agreement 000000D730321P5Q0002, Grant No. 70-2021-00145 02.11.2021).

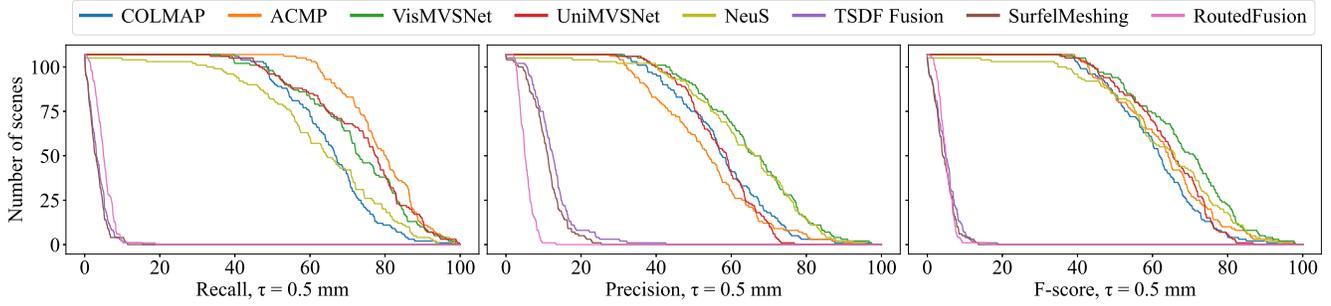


Figure 6. **Average performance per method** as the number of scenes with reconstruction quality better than the value on the X axis.

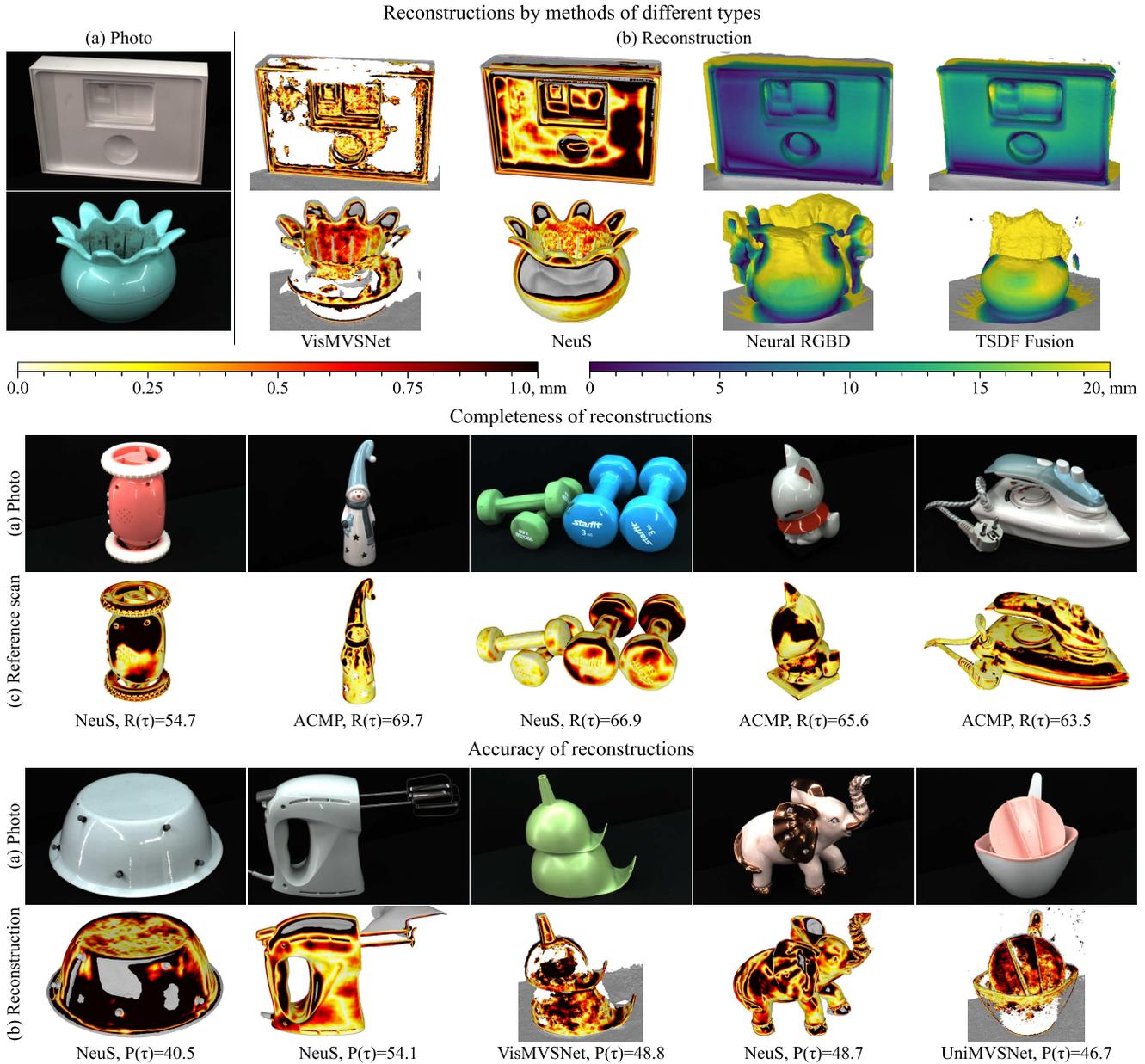


Figure 7. **Qualitative results.** (a) A photo of the scene. (b) Reconstruction with color-coded distance to the SL scan. (c) The SL scan with color-coded distance to the reconstruction. The two bottom parts show only the best result per scene and the respective value of Recall or Precision for $\tau = 0.5$ mm. Grey color represents undefined distance to the real surface, as explained in the supplementary material.

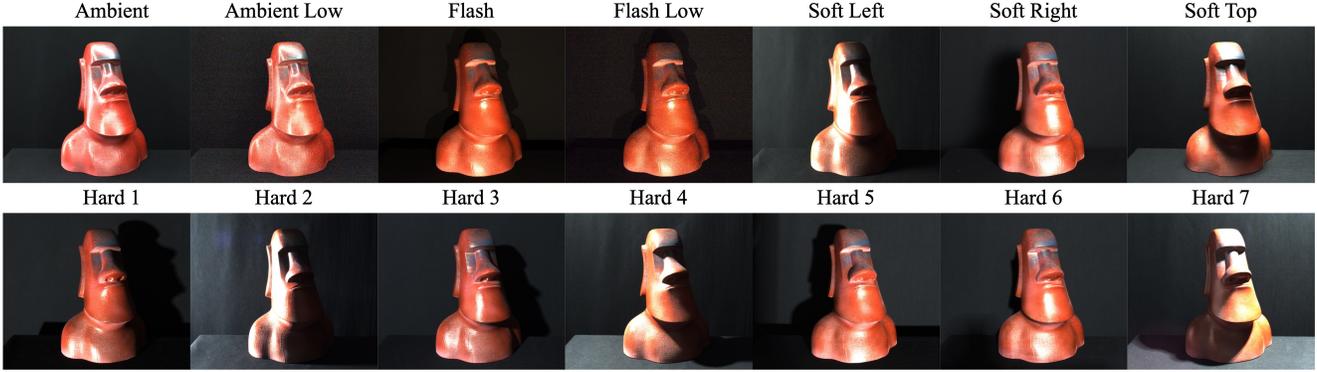


Figure 8. **Lighting setups in our dataset:** ambient lighting, a flashlight attached to the camera, three variants of soft light coming from different directions, and seven variants of hard light. For the ambient and the flash lighting there is an additional low exposure variant: notice a higher noise level in Ambient Low and Flash Low.

Supplementary material

In Table 4 we show an extended comparison of our dataset to a number of relevant datasets. In Appendices A and B we discuss the choice of lighting variations in our dataset, and the choice of sensors. In Appendices C and D we provide details of data acquisition and camera calibration processes. In Appendices E and F we provide details of testing and evaluation of 3D reconstruction methods on our dataset. In Appendix G we show complete evaluation results. Finally, in Appendix H we describe and discuss the variability of key surface reflection parameters in our dataset.

A. Lighting variability

In Figure 8 we illustrate all 14 lighting setups in our dataset. *Ambient Low* and *Flash Low* correspond to *real-time / high-noise* camera settings for the ambient diffuse lighting and the phone flashlight. We aimed to provide a broad range of realistic lighting conditions: directional light sources and flashlights of the phones provide eight samples of “hard” light, typical, for example, for streetlight; soft-boxes provide three samples of diffuse light, typical for indoor illumination; LED strips imitate ambient light, typical for cloudy weather.

B. Sensors

We aimed to include commodity RGB-D sensors with different properties. Smartphones are ubiquitous, and are increasingly commonly augmented with a depth sensor, while Kinect and RealSense devices represent dedicated RGB-D cameras. We included smartphones that capture depth with a time-of-flight sensor, Kinect v2 also uses a time-of-flight sensor but with a higher resolution and accuracy, and the RealSense device uses stereo-matching of infra-red images (a different technology). These devices capture depth maps with different resolution, level of noise, and different artefacts, as briefly illustrated in Figure 4. The structured-light scanner provides the reference 3D data for these devices.

These devices include commodity RGB sensors with dif-

ferent resolutions; we supplemented them with industrial RGB cameras with high-quality optics and low-noise sensors. The pair of industrial RGB cameras can serve as yet another source of depth maps based on stereo-matching of RGB images, in contrast to IR images in RealSense.

The data captured in the same environment with different sensors can be used to test generalization ability of computer vision methods, or to train a generator of synthetic data to reproduce a specific sensor, etc.

C. Data acquisition details

Lens settings. We set the focal distance for the industrial RGB cameras and for the cameras of the phones to the expected average distance from the cameras to the surface of the scanned object, namely 62.5 cm. To extend the depth of field for the industrial cameras to the whole scanning area, we set the aperture to the minimal value that does not cause any visual blur due to diffraction. We kept the lens parameters for Kinect and RealSense at their factory settings.

White balance. We set the white balance for all RGB sensors fixed at the start of each scanning session (*i.e.*, daily), using a black-and-white calibration pattern. Most of our light sources have the same light temperature so that their light appears white under this setting, except for the ambient illumination which has a somewhat higher green component, and for the flashlight of the phones which has a yellow tint. We did not fix the white balance setting for Kinect, for which this control is not available.

Exposure and gain. For Kinect, it is not possible to control camera exposure and gain directly. Instead, the built-in auto-exposure function is permanently turned on, which sets the exposure and gain so that the mean pixel value over the whole image is 50% gray. Since in our setup a large portion of the image is the black background, the built-in algorithm produces over-exposed images. To minimize this effect, we added a dimming light filter mounted on a servomotor to the rig, which is placed in front of the RGB camera of Kinect automatically when the bright light sources are activated,

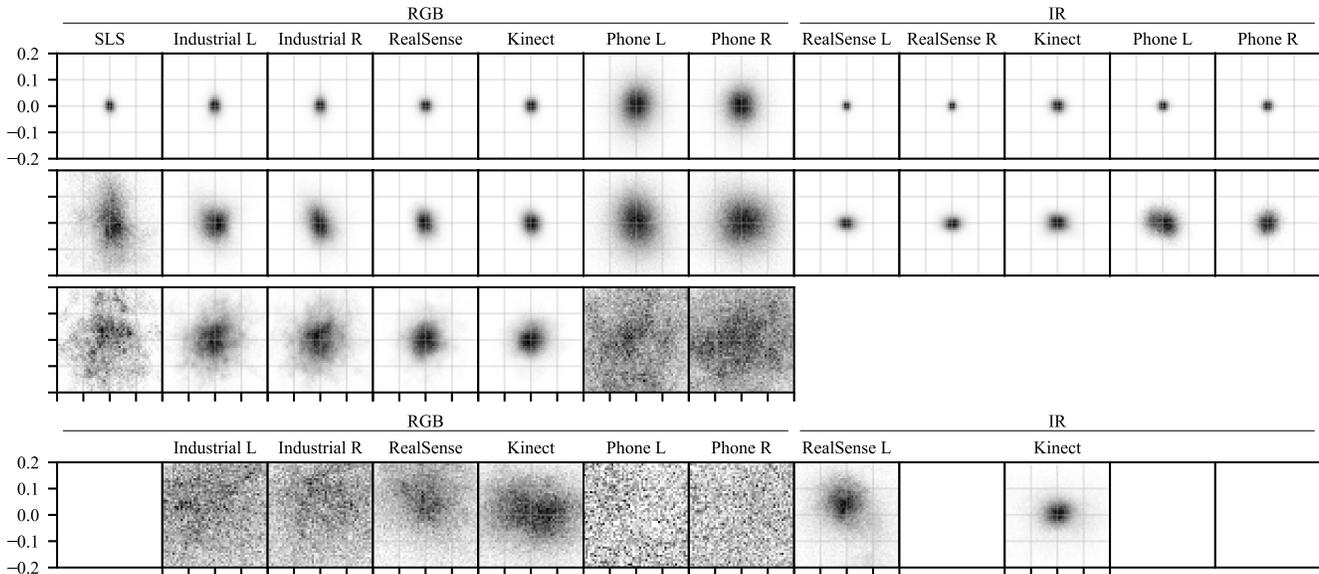


Figure 9. **Error distributions for camera calibration**, as histograms of 2D differences between the detection of a feature on the calibration pattern and its projection from 3D to the image. The range of each histogram is $[-0.2, 0.2]$ pixels in each dimension. Each histogram represents one sensor at one of the three steps of the calibration procedure: estimation of intrinsic camera models (first row), estimation of relative position of the sensor within the rig (second row), estimation of position of the sensor on the scanning trajectory (third row). The histograms for the IR sensors in the third row are not shown, since the trajectory of the IR sensors is calculated directly from the trajectory of their RGB companions. The fourth row shows the errors for the baseline approach of estimation of relative positions of the sensors within the rig (compare with the second row).

and removed for the dim ones.

D. Camera calibration details

The original implementation of the calibration pipeline of [65] supports calibration of rigid camera rigs, however, a straightforward application of this pipeline for our camera rig in its entirety proved to be numerically unstable due to the properties of the setup. Firstly, the included sensors have a large variation in the field of view (from 30° for the SLS camera to 90° for the IR sensors of RealSense) and resolution (from 0.04 MPix for the IR sensors of the phones to 40 MPix for their RGB sensors). Secondly, the focus of the cameras of the phones, being fixed programmatically, fluctuates slightly over time, which we relate to thermal deformations of the device (see [20] for a study of such an effect). Finally, the camera rig deforms slightly depending on its tilt in different scanning positions. To avoid the loss of accuracy we split the calibration procedure into several steps.

First, we obtained intrinsic camera models for each sensor independently. Then, for the sensors with a relatively high resolution and stable focus, namely the SLS and all RGB sensors except the sensors of the phones, we estimated their relative position within the rig for its vertical orientation. Next, we estimated the relative position of RGB sensors of the phones within the rig, w.r.t. the other sensors. After that, we estimated the position of each IR (depth) sensor w.r.t. the RGB sensor of the respective device, assuming

that the sensors are fixed rigidly to the frame of the device. Finally, we estimated the position of each RGB sensor individually for different scanning positions of the robotic arm, and then linked the positions of all sensors together through a scanning position with the vertical orientation of the rig.

This whole procedure required capturing thousands of images of the calibration pattern, which we almost fully automated with the use of the robotic arm, except for several manual reorientations of the pattern.

In Figure 9 we show distributions of calibration error for each sensor at different stages of the calibration procedure. Note that the resolution of the RGB sensors of the phones is relatively high compared to the other RGB sensors so a higher value of the calibration error measured in pixels is expected.

At the bottom of Figure 9 we show calibration errors for the straightforward application of the calibration pipeline to our camera rig. Here, we estimated the relative positions for a subset of sensors within the rig simultaneously. The mean error of the straightforward approach is 1.8-5.7 times higher depending on the sensor compared to ours.

E. Parameters of tested methods

We tested all methods using ground-truth camera poses and intrinsic camera models, after the refinements of camera poses and depth camera calibration. We used RGB images from the right industrial camera and the depth maps from Kinect at full resolution after removing distortion (which

resulted in a small amount of cropping at the boundaries), specifically, 2368×1952 for RGB and 496×400 for depth.

For each method, we tried to pick the values of the method parameters which resulted in the best reconstructions on our dataset, based on 5-10 typical scenes. We describe the parameters using the original notations from the respective works.

We tested **COLMAP** [61–63] with parameters set to their “performance” values recommended in the software documentation. Specifically, for feature extraction, we enabled estimation of affine shape of SIFT features, and enabled the more discriminative DSP-SIFT features instead of plain SIFT; for feature matching, we disabled estimation of multiple geometric models per image pair; for patch-match stereo, we enabled the regularized geometric consistency term; finally, for stereo fusion we set the minimum number of fused pixels to produce a point to 3, the maximum relative difference between measured and projected pixels to 1, and the maximum depth error to 1 mm.

We tested **ACMP** [92, 94] with the default parameters used in the source code. For this MVS method and the two methods below we sampled the depth hypotheses uniformly from 473 mm to 983 mm with a 2 mm resolution, and for view selection used the strategy proposed in [96], applied to the reference SL scan instead of a sparse reconstruction of the scene.

To test **VisMVSNet** [101, 102] we trained it from scratch on BlendedMVG dataset [95], which is an extended version of the BlendedMVS dataset [97], originally used by the authors of the method. We used the original training parameters and trained the network on a single Nvidia GTX 1080Ti GPU for 2 epochs with a batch size of 2 (342K iterations), using Adam optimizer [32] with a learning rate of 10^{-4} .

We tested VisMVSNet with the number of neighboring source images $N_v = 7$, the numbers of depth hypotheses $N_{d,1}, N_{d,2}, N_{d,3} = 64, 32, 16$, the minimal number of consistent views during fusion $N_f = 4$, and the fusion probability thresholds $p_{t,1}, p_{t,2}, p_{t,3} = 0.8, 0.7, 0.8$.

We tested **UniMVSNet** [53, 54] using the published model trained on DTU dataset [29] and fine-tuned on BlendedMVS dataset. We used the number of neighboring source images $N = 11$, the numbers of depth hypotheses $M_1, M_2, M_3 = 64, 32, 16$, and the fusion probability thresholds $\phi_1, \phi_2, \phi_3 = 0.1, 0.15, 0.9$.

We tested **NeuS** [81, 82] with the original hyperparameters, optimizing the network for 300K iterations. To accelerate convergence and prevent oversmoothing of reconstruction we cropped the input images for this method to the bounding box of the object expanded by ~ 10 pixels.

To test **TSDF Fusion** [12, 28] we used an implementation of this algorithm from Open3D library [105], with a voxel size of 3 mm, and a TSDF truncation distance of 2 cm. We additionally tested an implementation of this algorithm

from [48, 49] and obtained very similar results; we do not report them.

We tested **SurfelMeshing** [64, 66] with the number of inliers for depth filtering set to 1, the depth map erosion radius set to 0.1, and with 1 iteration of median filtering.

To test **RoutedFusion** [85, 86] we trained it from scratch on ModelNet dataset [90]. We used the original training parameters for ShapeNet dataset [8], with an increased level of synthetic noise $\sigma = 0.01$, as suggested by the authors. We tested RoutedFusion with a grid resolution of 384^3 , which corresponds to a voxel size of 2.6 mm.

We tested **Neural RGB-D Reconstruction** [3, 57] with the original hyperparameters, optimizing the network until convergence for 200K iterations, and sampling the coarse points on the ray for every 2 mm.

F. Evaluation details

We evaluated 3D reconstruction methods using *precision* $P(\tau)$ defined as the percentage of reconstruction points which are closer to the reference surface than a distance threshold τ ; *recall* $R(\tau)$ defined as the percentage of reference points which are closer to the reconstruction than the distance threshold; and *F-score* $F(\tau)$ defined as the harmonic mean of precision and recall. Additionally, we calculated the mean distance from the reconstruction to the reference, and the mean distance from reference to reconstruction, which we report in Appendix G.

To calculate the measures based on the distance from the reference to the reconstruction, we used vertices of the full SL scan; their average distance from the nearest neighbor is around 0.15 mm. To evaluate the methods which reconstruct the surface in the form of a triangular mesh, we sampled points from the mesh uniformly at a sampling distance of 0.1 mm.

The reference SL scans are incomplete, so the distance from some reconstruction points to the SL scan does not represent the distance to the real surface of the object, specifically, if the point lies near the missing part of the surface. To calculate the measures using only the points for which the distance to the real surface is reliable we used the approach of [67], and extended it using our new occluding surface.

First, we only kept the points which lie in the free space between the SL scanner and the object or near the surface of the object: we checked if the depth of the point w.r.t. SL scanner for any of its scanning positions is less than the depth given by the SL scan, plus a small tolerance $t_{\text{subsurf}} = 3$ mm to keep the points just below the surface for evaluation.

Next, to evaluate the precision metric, we checked if a point is closer to the real surface of the object than a distance threshold. For every reconstructed point, we calculated the distance to the SL scan and the distance to the occluding surface. If the distance to the SL scan was below the threshold, we considered the point to be closer to the real surface than

the threshold. If the distance to the occluding surface was above the threshold, we considered the point to be farther from the real surface than the threshold. In any other case (in which the point can only be closer to the occluding surface and farther from the SL scan than the threshold, since the occluding surface encloses the SL scan by definition), we considered the distance from the point to the real surface to be unknown and excluded the point from the calculation of the measure.

To visualize the distance from the reconstruction to the reference and to calculate the mean distance we used a similar strategy and considered the distance to the real surface to be unknown whenever the point was closer to the occluding surface than to the SL scan. To account for the approximate nature of our occluding surface and to prevent computational instabilities at points where it must coincide with the SL scan exactly in case of exact calculations, we replaced the distance to the occluding surface in all calculations by its value increased by $\varepsilon_{occ} = 0.1$ mm.

G. Complete evaluation results

In Figures in a separate PDF file¹ we show qualitative evaluation of 3D reconstruction methods. For Neural RGB-D surface reconstruction we show the results only for four scenes: two relatively easy ones, based on performance of all the methods, `dragon` and `small_wooden_chessboard`, and two relatively hard ones `green_flower_pot` and `white_box`.

In each figure in the first column titled *Reconstruction*, we show the reconstruction produced by each method, and at the bottom of this column we show the reference surface from the SL scanner. In the column *Accuracy on reconstruction*, we show the reconstructed surface with color-coded distance to the reference surface; at the bottom of this column we show a photo of the scene. In the column *Accuracy on reference*, we show the reference surface, with the color showing the distance from the reconstruction to the reference. For each vertex of the reference surface, the distance is averaged over the reconstructed points for which this vertex is the closest one. Thanks to this projection of the distance values from the reconstructed points to the reference surface, these images show the accuracy of all points not just the ones closest to the camera. In the column *Completeness on reference*, we show the reference surface with color-coded distance to the reconstructed surface.

We use two colormaps for two different scales of error. In the last three columns, the points with no definite value of the distance are grey.

Since the methods TSDF Fusion, RoutedFusion, SurfelMeshing, and Neural RGB-D surface reconstruction produce the surface with a significant error, and in particular deep below the reference surface, we increased the value of

sub-surface tolerance t_{subsurf} for these methods from 3 mm to 20 mm.

In Figures in a separate PDF file¹ we show the recall, precision, and F-score curves for reconstructed surfaces produced by the methods for each scene. Each curve represents the measure in percent for different values of the distance threshold τ in millimeters. Additionally, we mark the mean distance from the reference to the reconstruction for each method with a vertical dashed line on the recall plot, and the mean distance from the reconstruction to the reference with a vertical dashed line on the precision plot. Note that for some methods these lines may be out of the plot range.

Finally, in Figures 10 and 11 we show the values of the recall and precision metrics calculated with the distance threshold $\tau = 0.5$ mm for the reconstructions produced by the RGB-based methods COLMAP, ACMP, VisMVSNet, UniMVSNet, and NeuS. The value of the measure for each method is represented by the right edge of the bar with the respective color. The scenes in each figure are sorted by the best result on the scene.

H. Material properties in our dataset

We provide more information on the qualitative descriptors of surface reflection parameters assigned to each object and their relation to performance indicators for various reconstruction methods.

Identification of surface properties. As outlined in the main text, all labels were assigned to objects by visually inspecting photos of each object, to provide a preliminary assessment of the importance of various factors for reconstruction quality. Multiple photos of each object were used, with backlight proving particularly useful to establish the degree of translucency. The labels refer to the *dominant* material or materials of the objects, which have most influence on integral metrics of reconstruction quality. While this classification is imprecise and not a substitute for photometric measurements, as we use only rough estimation for each property, typically 3 buckets, we expect that visually assigned labels are highly correlated with the actual physical parameters. Where necessary, more than one label was assigned (*e.g.*, for objects consisting of multiple parts with distinct reflectivity properties, with no single dominant material); in this case several labels for some properties are used simultaneously.

Most of our labels are aligned with typical parameters of simple reflectance models: diffuse and specular reflection coefficient, shininess/roughness, measuring the reflection peak width, and translucency.

- *Translucency.* We assess the degree of translucency for the materials of the large parts of the object, from *none* (the default), through *low*, *medium*, *high*, all the way to completely *transparent*.
- *Reflection sharpness.* We characterize how sharp the

¹See at skoltech3d.appliedai.tech/data/skoltech3d_supp_results.pdf

reflectance function peak is, if highlights are present on the object (i.e., it is not purely diffuse), grouping the objects into four categories, from low to very high.

- *Specularity.* We visually estimate the ratio of specular to diffuse reflection for the dominant object materials, resulting in a degree of manifestation of view-dependent highlights, complementing *mirror-like* tag with a weaker label. We distinguish between no view dependent highlights (*diffuse*), largely diffuse highlights (*low*), somewhat diffuse light reflections and wide highlights (*medium*), and narrow highlights for partially mirror-like surfaces where one can only see sharp reflections of lights (*high*).
- *Mirror-like.* A binary tag for surfaces for which, in addition to a near-perfect reflection peak, the diffusive component is low relative to specular. This label overlaps very high reflection sharpness labels.
- *Metallic.* A binary tag for surfaces made of metal; compared to dielectric surfaces, these surfaces are always completely non-transparent, and the specular component of reflected light tend to have the same color as diffuse, while for dielectric materials it is closer to the color of the incident light.
- *Geometric features.* We visually assess the dominant qualitative scale of geometric structures of the surface, if any. We seek to distinguish between fine 3D structure with characteristic scale close but above SLS 3D resolution (*small*), a larger, discernable geometry with feature size equal to a small fraction of the object size (*medium*). Additionally, among surfaces lacking a dominant feature scale, we identified those with flat areas of sufficiently large size to make an impact on overall reconstruction quality (*flat* label), in the absence of 3D texture.
- *Texture type.* We differentiate between several types of textured surfaces: most of the surface covered by high-frequency image texture (*color*), sufficiently small scale (below or close to what the scanner can resolve), high-frequency displacement variation (*3D*), or other texture-like imperfections such as dirt, speckles, or visible roughness (*imperfections*).

Correlating properties and reconstruction performance.

Most surface reflectance features mentioned above are expected to influence performance of common reconstruction methods.

We use data-driven approach to model the variability in a given measure of reconstruction performance caused by different reflectance properties. Viewing all measures described in Appendix F as target variables, we constructed a numerical feature representation for each scene in our dataset; and used sparse L_1 regularised (Lasso) regression [78] to determine how each surface reflectance feature modulates reconstruction performance in each scene.

As we wanted to determine which features had statistically the most influence on reconstruction performance, we used the automatic feature selection mechanism built into the Lasso algorithm.

Our target variables were *precision* P (0.5 mm), *recall* R (0.5 mm), and *F-score* F (0.5 mm) for the image-based methods ACMP, COLMAP, VisMVSNet, and UniMVSNet (Appendix E). For constructing the feature representation of each scene, we formed 25 binary features from the 7 properties mentioned above by encoding the presence of each surface property as 1, and its absence as 0. For fitting the model, we used the implementation of Lasso available in `scikit-learn` [52]; we z-scored all features, and, for each target variable, sought to find an optimal value of the regularisation parameter α (a minimizer of RMSE measure) using leave-one-out cross-validation by varying it over a range [0.01, 2.0], and record the values of the R^2 statistic, the regularization weight α as well as the regression weights of the final fit.

Table 5 displays (normalized) values of coefficients weighting the (normalized) value of each feature in each constructed linear model. Most regression problems have demonstrated a consistent value of the regularization weight $\alpha \in [0.18, 0.57]$ with an average of 0.36; the values of coefficient of determination R^2 vary from 0.24 to 0.54 with a mean of 0.43, indicating that a reasonable performance has been achieved.

Factors affecting algorithms performance. Using the process described above, we have identified five labels having the most pronounced *negative effect* on reconstruction quality for image-based MVS methods, according to a mean F-score across methods: very high, high, or medium reflection sharpness, high specularity, medium-scale geometric features; additionally, a negative effect from high translucency, low specularity is also expected, unlike moderate but notable contribution of non-metallic and diffuse materials. Conversely, having any type of texture (either color, 3d, or texture-like surface imperfections) *contributes positively* to reconstruction performance as one expect for MVS-type algorithms; the same can be said about non-translucent objects.

We consider this study very preliminary, given the approximate nature of our labels. Its results are encouraging as these show clear correlation between surface reflectance properties and algorithm performance.

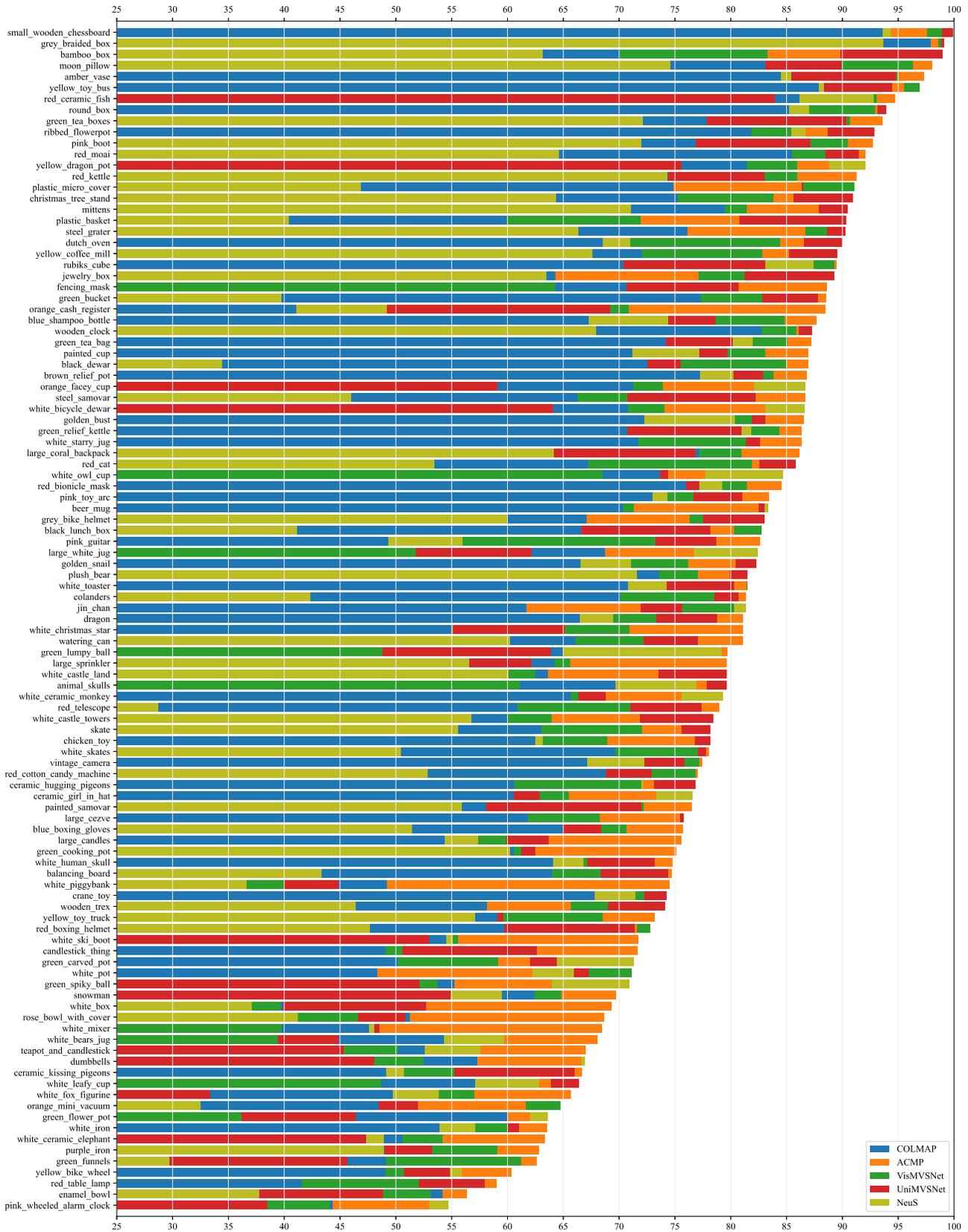


Figure 10. Recall for the RGB-based methods for all scenes with the distance threshold $\tau = 0.5$ mm.

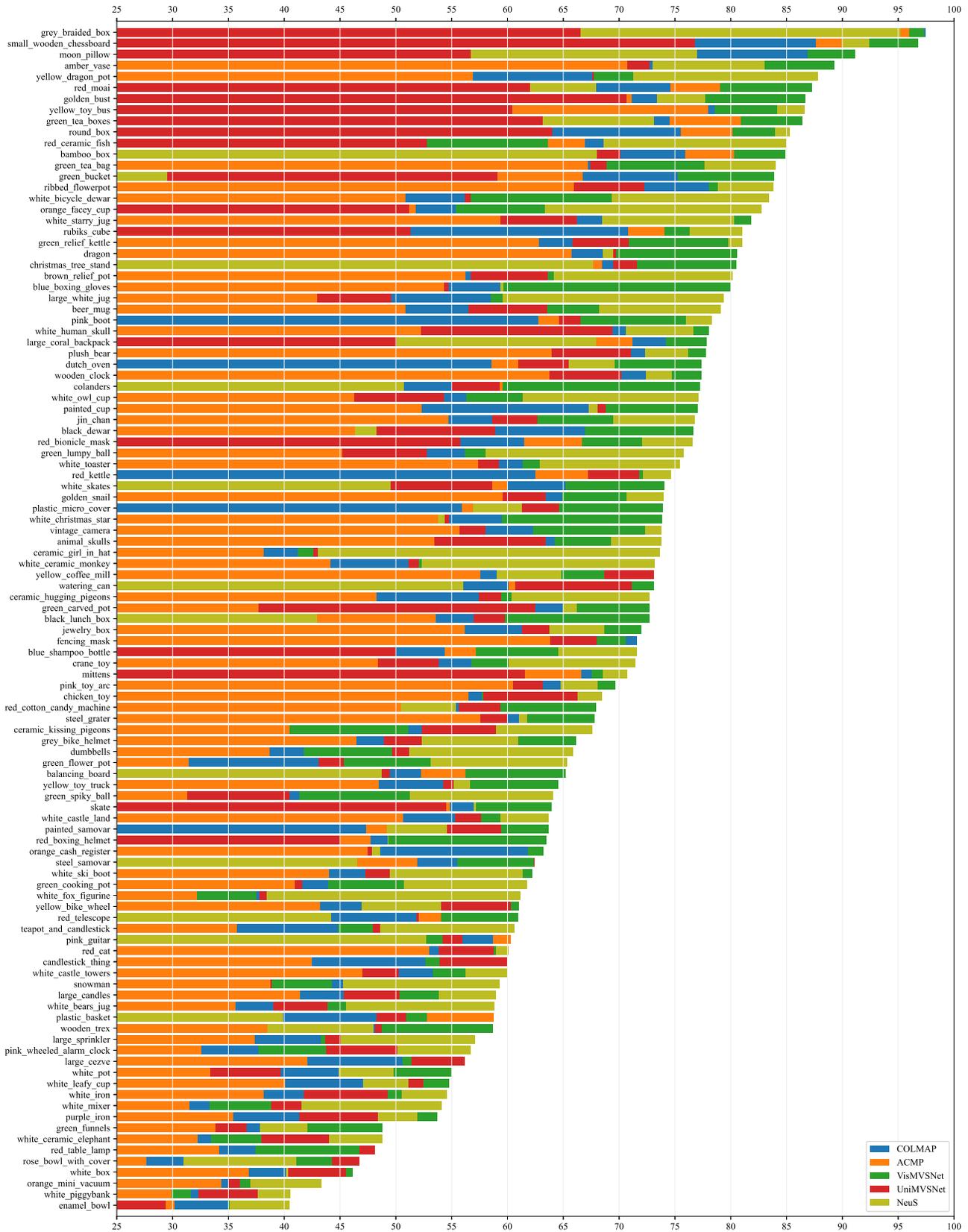


Figure 11. Precision for the RGB-based methods for all scenes with the distance threshold $\tau = 0.5$ mm.

Dataset	Devices	RGB res.	Depth res.	Hi-res. geom.	Video	Camera positioning	Poses/scene	Lighting	# Scenes	#points in high-res pointclouds	#Frames	Scene type	Data provided	Target tasks
DTU	2 RGB unknown model	1200×1600	—	✓	—	6DOF robotic	49 or 64	8	80	13.4M	27K	small objects	<ul style="list-style-type: none"> • intrinsics, extrinsics • raw RGB • undistorted RGB • surfel clouds • voxels where GT is known 	MVS
ETH3D (hi)	Nikon D3X, Faro Focus X 330	6048×4032	—	✓	—	tripod	10-70	U	13 train 12 test	28M	11K	in/outdoor large scenes	<ul style="list-style-type: none"> • intrinsics, extrinsics • raw RGB • undistorted RGB • raw 3D scans • scans with outliers removed • rendered depth maps 	MVS
ETH3D (lo)	4 global shutter, Faro Focus X 330	752×480	—	✓	Mono	handheld	160-300	U	5 train 5 test	28M	10K	in/outdoor large scenes	same as ETH3D hi-res	MVS, stereo
TnT	DJI Zenmuse X5R, Sony a7S II, Faro Focus X 330	3840×2160	—	✓	RGB	handheld + gimbal	150-500	U	7 train 14 test	5.5M-53.4M	148K	large objects, in/outdoor large scenes	<ul style="list-style-type: none"> • raw videos • sets of extracted RGB • raw 3D scans • point clouds • camera params (extracted with COLMAP) 	SfM and MVS pipelines, view synthesis
BlendedMVS / MVG	unknown RGB cameras	1536×2048 (rescaled)	—	—	—	unknown	20-1000	U	113 in MVS 502 in MVG	—	18K/110K	MVS: large outdoor, small scale objects, sculptures MVG: mostly aerial imaging	<ul style="list-style-type: none"> • MVS and MVG: <ul style="list-style-type: none"> • intrinsics and extrinsics • 768x576 rendered RGBD • masked rendered RGB • MVS only: <ul style="list-style-type: none"> • 2048x1536 rendered RGBD • reconstructed meshes 	MVS (including training)
Redwood	PrimeSense Carmine	1280×1024	640×480	—	—	handheld	—	U	10K	—	144K	medium objects	<ul style="list-style-type: none"> • RGBD • for 398 scenes, meshes 	—
SUN RGBD	RealSense, Xtion, Kinect v1 Kinect v2	1920×1080 640×480	628×428 512×424 640×480	—	—	handheld	—	U	unknown	—	10K	indoor scenes	<ul style="list-style-type: none"> • 2D polygon object segmentation; • 3D object bounding boxes; • 3D polygon room layouts. 	classification, segmentation
BigBIRD	5x Canon Rebel T3, 5x PrimeSense Carmine 1.08	4272×2848	640×480	—	—	fixed 5 cameras+ turntable	600	1	120	—	75K	small objects	<ul style="list-style-type: none"> • hi-res RGB • low-res RGBD • point clouds for each scan • merged point clouds 	instance recognition, category recognition, 3D reconstruction
CORBS	Kinect, 3Digify	1920×1080	512×424	✓	—	RGBD (tracked)	5 trajectories, 700-7000 frames each	U	4	unknown	47K	indoor scenes	<ul style="list-style-type: none"> • RGBD • IR • camera trajectories • reconstructed mesh 	SLAM
ScanNet	Structure	1296×968	640×480	—	—	handheld	—	U	1513	—	2.5M	indoor scenes	<ul style="list-style-type: none"> • intrinsics, camera trajectories • reconstructed meshes (semantically labeled) • aligned CAD objects 	classification, retrieval, voxel level annotation
ARKitScenes	Faro Focus S70 Apple iPad Pro	4032×3024 3680×2760	256×192	✓	—	handheld	—	U	5047	—	450K	indoor scenes	<ul style="list-style-type: none"> • intrinsics, camera trajectories • reconstructed meshes (semantically labeled) 	classification, segmentation, depth upsampling
RGB-D-D	LUCID Helios ToF Huawei P30 Pro	3648×2736	640×480 240×180	—	—	unknown	—	U	4811	—	4811	portraits, models, plants, lights	<ul style="list-style-type: none"> • intrinsics, camera trajectories 	depth upsampling
Ours	2x DFK 33ux250 cameras 2x Huawei Mate 30 Pro RealSense D435 Kinect V2 RangeVision Spectrum	2448×2048 7296×5472 1920×1080	240×180 512×424 1280×720	✓	—	6DOF robotic	100	14	107	9.5M-66.5M	877K	small objects	<ul style="list-style-type: none"> • intrinsics, extrinsics • RGB: cameras • undistorted RGB • RGBD: phones, Kinect, RealSense • undistorted RGBD • raw IR: phones, Kinect, RealSense • raw 3D scans • reconstructed mesh • occluding mesh 	MVS, 3D reconstruction, depth fusion, depth upsampling

Table 4. Comparison of our dataset to the most widely used related datasets. U indicates uncontrolled lighting; frames are counted per sensor, *i.e.*, all data from an RGB-D sensor are counted as a single frame. The number of separate images acquired may be considerably larger (1.4 M for our dataset). All scenes, from both training and testing sets, were counted.

Parameter	Value	Recall (0.5 mm)				Precision (0.5 mm)				F-score, thres=0.5mm				
		ACMP	COL.	Vis.	Uni.	ACMP	COL.	Vis.	Uni.	ACMP	COL.	Vis.	Uni.	Mean
Translucency	none	0	0	0.91	1.45	0	1.3	0.48	1.57	0	1.04	0	1.85	0.72
	low	-0.23	0	0.3	0	0	0	0	0	0	0	0	0	0
	medium	0.67	0	0.62	0.83	0	0	0	0	0	0	0	0.38	0.09
	high	-1.72	-1.15	-1.98	-1.81	-1.93	-1.36	0	-0.56	-1.88	-1.53	-0.48	-1.19	-1.27
Reflection sharpness	no refl.	0	0	-0.17	0	0	0	0	0	0	0	0	0	0
	low	-0.83	-0.8	-0.64	-1.82	-0.56	0	1.13	0	-0.52	-0.57	0	-0.62	-0.43
	medium	-2.05	-0.58	-3.87	-3.37	-2.6	-2.45	-1.57	-3.26	-2.3	-2.18	-2.14	-3.32	-2.49
	high	-1.59	-0.12	-3.22	-1.46	-2.41	-2.93	-1.03	-3.66	-2.04	-2.34	-1.34	-2.71	-2.11
Specularity	very high	-2.24	-1.47	-4.8	-3.47	-3.99	-3.46	-0.91	-4.21	-3.4	-3.37	-1.51	-4.11	-3.1
	diffuse	-1.81	0	-3.66	-2.15	-1.15	0	0	0	-1.21	-0.5	0	-1.32	-0.76
	low	-0.75	-0.28	-4	-2.32	-1.39	-0.97	0	-1.02	-1.28	-1.32	0	-2.08	-1.17
	medium	0	0	-2.07	0	0	0	0	0.76	0	0	0	0	0
Geometric features	high	-1.66	-0.64	-4.33	-2.79	-2.02	-1.69	-0.91	-1.08	-1.91	-1.35	-1.56	-2.29	-1.77
	small	0	0.06	2.49	1.05	0	0	0.47	0	0.06	0.29	1.15	0.59	0.52
	medium	-1.58	-1.1	-0.79	-2.31	-1.43	-1.25	0	-2.08	-1.5	-1.42	0	-2.42	-1.34
	flat	0.23	0	1.15	0.59	2.08	0.27	0	0	1.58	0	0.3	0.32	0.55
Mirror-like	NA	0	0	0	0	0	0	0	0	0	0	0	0	
	no	0	0	0	0	0	0	0	0	0	0	0	0	
Metallic	yes	0.42	0.72	0.85	0.23	0	1.13	1	1.38	0	1.3	0.8	0.83	0.73
	no	-0.62	-0.21	-0.23	-1.3	-1.25	-0.34	-1.24	-1.16	-1.14	-0.24	-1.21	-1.27	-0.97
Texture type	yes	0	0	0.67	0	0	0	0	0	0	0	0	0	0
	color	2.61	2.41	3.9	3.41	3.65	2.48	2.31	2.97	3.36	2.59	2.73	3.34	3.01
	imperf.	2.06	1.96	5.1	4.14	2.61	2.01	1.24	3.4	2.59	2.43	2.32	3.92	2.81
Coef. of det.	3d	1.23	1.89	2.42	1.51	3.43	4.07	3.21	4.14	2.8	3.13	2.88	3.01	2.96
Reg. coef.	R^2	0.39	0.24	0.54	0.41	0.46	0.47	0.39	0.49	0.44	0.4	0.44	0.48	
	α	0.26	0.57	0.18	0.3	0.37	0.4	0.39	0.38	0.37	0.28	0.54	0.27	

Table 5. **Data-driven estimates of surface reflection features in our datasets.** Note the *negative effect* on reconstruction quality (rightmost column, negative numbers highlighted in **red**) of very high, high, or medium reflection sharpness, high specularity, medium-scale geometric features. Conversely, note the *positive contribution* to reconstruction performance of color textures, 3d textures, or texture-like surface imperfections (rightmost column, large positive values highlighted in **blue**). *COL.*, *Vis.*, and *Uni.* denote COLMAP, VisMVSNet, and UniMVSNet respectively.

References

- [1] Henrik Aanæs, Rasmus Ramsbøl Jensen, George Vogiatzis, Engin Tola, and Anders Bjorholm Dahl. Large-scale data for multiple-view stereopsis. *International Journal of Computer Vision*, 120(2):153–168, 2016. [2](#)
- [2] Kara-Ali Aliev, Artem Sevastopolsky, Maria Kolos, Dmitry Ulyanov, and Victor Lempitsky. Neural point-based graphics. In *European Conference on Computer Vision*, pages 696–712. Springer, 2020. [2](#)
- [3] Dejan Azinović, Ricardo Martin-Brualla, Dan B Goldman, Matthias Nießner, and Justus Thies. Neural rgb-d surface reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6290–6301, 2022. [6](#), [11](#)
- [4] Connelly Barnes, Eli Shechtman, Adam Finkelstein, and Dan B Goldman. PatchMatch: A randomized correspondence algorithm for structural image editing. *ACM Transactions on Graphics (Proc. SIGGRAPH)*, 28(3), Aug. 2009. [2](#)
- [5] Gilad Baruch, Zhuoyuan Chen, Afshin Dehghan, Tal Dimry, Yuri Feigin, Peter Fu, Thomas Gebauer, Brandon Joffe, Daniel Kurz, Arik Schwartz, and Elad Shulman. ARK-itscenes - a diverse real-world dataset for 3d indoor scene understanding using mobile RGB-d data. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*, 2021. [3](#)
- [6] Matthew Berger, Joshua A Levine, Luis Gustavo Nonato, Gabriel Taubin, and Claudio T Silva. A benchmark for surface reconstruction. *ACM Transactions on Graphics (TOG)*, 32(2):1–17, 2013. [3](#)
- [7] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. *International Conference on 3D Vision (3DV)*, 2017. [2](#)
- [8] Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiang Xiao, Li Yi, and Fisher Yu. ShapeNet: An Information-Rich 3D Model Repository. Technical Report arXiv:1512.03012 [cs.GR], Stanford University — Princeton University — Toyota Technological Institute at Chicago, 2015. [3](#), [11](#)
- [9] Jaesung Choe, Sunghoon Im, Francois Rameau, Minjun Kang, and In So Kweon. Volumefusion: Deep depth fusion for 3d scene reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 16086–16095, October 2021. [2](#)
- [10] Massimiliano Corsini, Matteo Dellepiane, Federico Ponchio, and Roberto Scopigno. Image-to-geometry registration: a mutual information method exploiting illumination-related geometric properties. In *Computer Graphics Forum*, volume 28, pages 1755–1764. Wiley Online Library, 2009. [5](#)
- [11] D. Cremers and K. Kolev. Multiview stereo and silhouette consistency via convex functionals over convex domains. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(6):1161–1174, 2011. [3](#)
- [12] Brian Curless and Marc Levoy. A volumetric method for building complex models from range images. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques - SIGGRAPH '96*, pages 303–312, Not Known, 1996. ACM Press. [6](#), [11](#)
- [13] 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. [2](#)
- [14] Angela Dai, Christian Diller, and Matthias Niessner. Sg-nn: Sparse generative neural networks for self-supervised scene completion of rgb-d scans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. [3](#)
- [15] 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*, 36(3):1–18, July 2017. [3](#)
- [16] Angela Dai, Daniel Ritchie, Martin Bokeloh, Scott Reed, Jürgen Sturm, and Matthias Nießner. Scancomplete: Large-scale scene completion and semantic segmentation for 3d scans. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4578–4587, 2018. [3](#)
- [17] Angela Dai, Yawar Siddiqui, Justus Thies, Julien Valentin, and Matthias Niessner. Spsg: Self-supervised photometric scene generation from rgb-d scans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1747–1756, June 2021. [3](#)
- [18] Matteo Dellepiane and Roberto Scopigno. Global refinement of image-to-geometry registration for color projection. In *2013 Digital Heritage International Congress (Digital Heritage)*, volume 1, pages 39–46. IEEE, 2013. [5](#)
- [19] Simon Donne and Andreas Geiger. Learning non-volumetric depth fusion using successive reprojections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. [2](#)
- [20] Melanie Elias, Anette Eltner, Frank Liebold, and Hans-Gerd Maas. Assessing the influence of temperature changes on the geometric stability of smartphone-and raspberry pi cameras. *Sensors*, 20(3):643, 2020. [10](#)
- [21] Hugo Germain, Guillaume Bourmaud, and Vincent Lepetit. S2dnet: Learning accurate correspondences for sparse-to-dense feature matching. *arXiv preprint arXiv:2004.01673*, 2020. [5](#)
- [22] Xiaodong Gu, Zhiwen Fan, Siyu Zhu, Zuozhuo Dai, Feitong Tan, and Ping Tan. Cascade cost volume for high-resolution multi-view stereo and stereo matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. [2](#)
- [23] A. Handa, T. Whelan, J.B. McDonald, and A.J. Davison. A benchmark for RGB-D visual odometry, 3D reconstruction and SLAM. In *IEEE Intl. Conf. on Robotics and Automation, ICRA*, Hong Kong, China, May 2014. [3](#)
- [24] Lingzhi He, Hongguang Zhu, Feng Li, Huihui Bai, Runmin Cong, Chunjie Zhang, Chunyu Lin, Meiqin Liu, and Yao Zhao. Towards fast and accurate real-world depth super-resolution: Benchmark dataset and baseline. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9229–9238, June 2021. [3](#)

- [25] Daniel Herrera, Juho Kannala, and Janne Heikkilä. Joint depth and color camera calibration with distortion correction. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(10):2058–2064, 2012. 5
- [26] Jiahui Huang, Shi-Sheng Huang, Haoxuan Song, and Shi-Min Hu. Di-fusion: Online implicit 3d reconstruction with deep priors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8932–8941, June 2021. 3
- [27] Tak-Wai Hui, Chen Change Loy, , and Xiaoou Tang. Depth map super-resolution by deep multi-scale guidance. In *Proceedings of European Conference on Computer Vision (ECCV)*, pages 353–369, 2016. 3
- [28] Shahram Izadi, Andrew Davison, Andrew Fitzgibbon, David Kim, Otmar Hilliges, David Molyneaux, Richard Newcombe, Pushmeet Kohli, Jamie Shotton, Steve Hodges, and Dustin Freeman. 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 - UIST '11*, page 559, Santa Barbara, California, USA, 2011. ACM Press. 3, 6, 11
- [29] Rasmus Jensen, Anders Dahl, George Vogiatzis, Engil Tola, and Henrik Aanæs. Large scale multi-view stereopsis evaluation. In *2014 IEEE Conference on Computer Vision and Pattern Recognition*, pages 406–413. IEEE, 2014. 2, 7, 11
- [30] Junho Jeon and Seungyong Lee. Reconstruction-based pairwise depth dataset for depth image enhancement using cnn. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018. 3
- [31] Michael Kazhdan and Hugues Hoppe. Screened poisson surface reconstruction. *ACM Transactions on Graphics (ToG)*, 32(3):1–13, 2013. 4
- [32] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. 11
- [33] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics (ToG)*, 36(4):1–13, 2017. 2, 7
- [34] Sebastian Koch, Yurii Piadyk, Markus Worchel, Marc Alexa, Cláudio Silva, Denis Zorin, and Daniele Panozzo. Hardware design and accurate simulation for benchmarking of 3D reconstruction algorithms. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021. 3
- [35] Andreas Kuhn, Christian Sormann, Mattia Rossi, Oliver Erdler, and Friedrich Fraundorfer. Deepc-mvs: Deep confidence prediction for multi-view stereo reconstruction. In *2020 International Conference on 3D Vision (3DV)*, pages 404–413, 2020. 2
- [36] Vincent Leroy, Jean-Sebastien Franco, and Edmond Boyer. Shape reconstruction using volume sweeping and learned photoconsistency. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018. 2
- [37] Andreas Ley, Ronny Hänsch, and Olaf Hellwich. Syb3r: A realistic synthetic benchmark for 3d reconstruction from images. In *European Conference on Computer Vision*, pages 236–251. Springer, 2016. 3
- [38] Shanchuan Lin, Andrey Ryabtsev, Soumyadip Sengupta, Brian L Curless, Steven M Seitz, and Ira Kemelmacher-Shlizerman. Real-time high-resolution background matting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8762–8771, 2021. 4
- [39] Philipp Lindenberger, Paul-Edouard Sarlin, Viktor Larsson, and Marc Pollefeys. Pixel-perfect structure-from-motion with featuremetric refinement. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5987–5997, 2021. 5
- [40] Yuanzhi Liu, Yujia Fu, Fengdong Chen, Bart Goossens, Wei Tao, and Hui Zhao. Simultaneous localization and mapping related datasets: A comprehensive survey. *arXiv preprint arXiv:2102.04036*, 2021. 2
- [41] Keyang Luo, Tao Guan, Lili Ju, Yuesong Wang, Zhuo Chen, and Yawei Luo. Attention-aware multi-view stereo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1590–1599, 2020. 6
- [42] Xinjun Ma, Yue Gong, Qirui Wang, Jingwei Huang, Lei Chen, and Fan Yu. Epp-mvsnet: Epipolar-assembling based depth prediction for multi-view stereo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5732–5740, October 2021. 2
- [43] Marko Mihajlovic, Silvan Weder, Marc Pollefeys, and Martin R. Oswald. Deepsurfels: Learning online appearance fusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14524–14535, June 2021. 3
- [44] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European Conference on Computer Vision*, pages 405–421. Springer, 2020. 2
- [45] Zak Murez, Tarrence van As, James Bartolozzi, Ayan Sinha, Vijay Badrinarayanan, and Andrew Rabinovich. Atlas: End-to-end 3d scene reconstruction from posed images. In *European Conference on Computer Vision (ECCV)*, 2020. 3
- [46] Pushmeet Kohli Nathan Silberman, Derek Hoiem and Rob Fergus. Indoor segmentation and support inference from rgb-d images. In *ECCV*, 2012. 2
- [47] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3504–3515, 2020. 2
- [48] Matthias Nießner. Voxelhashing. <https://github.com/niessner/VoxelHashing>. 11
- [49] Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Marc Stamminger. Real-time 3d reconstruction at scale using voxel hashing. *ACM Transactions on Graphics (ToG)*, 32(6):1–11, 2013. 11
- [50] Michael Oechsle, Songyou Peng, and Andreas Geiger. Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. *arXiv preprint arXiv:2104.10078*, 2021. 2
- [51] Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Colored point cloud registration revisited. In *ICCV*, 2017. 2, 3

- [52] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830, 2011. **13**
- [53] Rui Peng. Unimvsnet. <https://github.com/prstrive/UniMVSNet>. **11**
- [54] Rui Peng, Rongjie Wang, Zhenyu Wang, Yawen Lai, and Ronggang Wang. Rethinking depth estimation for multi-view stereo: A unified representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8645–8654, 2022. **6, 11**
- [55] Michael JD Powell. The newuoa software for unconstrained optimization without derivatives. In *Large-scale nonlinear optimization*, pages 255–297. Springer, 2006. **5**
- [56] Ruslan Rakhimov, Andrei-Timotei Ardelean, Victor Lempitsky, and Evgeny Burnaev. Npbg++: Accelerating neural point-based graphics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15969–15979, 2022. **2**
- [57] Neural RGB-D Surface Reconstruction. Neural rgbd surface reconstruction. <https://github.com/dazinovic/neural-rgbd-surface-reconstruction>. **11**
- [58] Gernot Riegler and Vladlen Koltun. Stable view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12216–12225, 2021. **2**
- [59] Daniel Scharstein, Heiko Hirschmüller, York Kitajima, Greg Krathwohl, Nera Nešić, Xi Wang, and Porter Westling. High-resolution stereo datasets with subpixel-accurate ground truth. In *German Conference on Pattern Recognition*, pages 31–42. Springer, 2014. **2**
- [60] Daniel Scharstein, Heiko Hirschmüller, York Kitajima, Greg Krathwohl, Nera Nešić, Xi Wang, and Porter Westling. High-resolution stereo datasets with subpixel-accurate ground truth. In *German conference on pattern recognition*, pages 31–42. Springer, 2014. **2**
- [61] Johannes Lutz Schönberger. Colmap. <https://github.com/colmap/colmap>. **11**
- [62] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-Motion Revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. **6, 11**
- [63] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise View Selection for Unstructured Multi-View Stereo. In *European Conference on Computer Vision (ECCV)*, 2016. **1, 2, 6, 11**
- [64] Thomas Schops. Surfelmeshing. <https://github.com/puzzlepaint/surfelmeshing>. **11**
- [65] Thomas Schops, Viktor Larsson, Marc Pollefeys, and Torsten Sattler. Why having 10,000 parameters in your camera model is better than twelve. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2535–2544, 2020. **5, 10**
- [66] Thomas Schops, Torsten Sattler, and Marc Pollefeys. SurfMeshing: Online Surfel-Based Mesh Reconstruction. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(10):2494–2507, Oct. 2020. **3, 6, 11**
- [67] Thomas Schops, Johannes L Schönberger, Silvano Galliani, Torsten Sattler, Konrad Schindler, Marc Pollefeys, and Andreas Geiger. A multi-view stereo benchmark with high-resolution images and multi-camera videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3260–3269, 2017. **2, 6, 7, 11**
- [68] Steven M Seitz, Brian Curless, James Diebel, Daniel Scharstein, and Richard Szeliski. A comparison and evaluation of multi-view stereo reconstruction algorithms. In *2006 IEEE computer society conference on computer vision and pattern recognition (CVPR'06)*, volume 1, pages 519–528. IEEE, 2006. **2**
- [69] Inwook Shim, Tae-Hyun Oh, Joon-Young Lee, Jinwook Choi, Dong-Geol Choi, and In So Kweon. Gradient-based camera exposure control for outdoor mobile platforms. *IEEE Transactions on Circuits and Systems for Video Technology*, 29(6):1569–1583, 2018. **4**
- [70] Ukcheol Shin, Jinsun Park, Gyumin Shim, Francois Rameau, and In So Kweon. Camera exposure control for robust robot vision with noise-aware image quality assessment. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1165–1172. IEEE, 2019. **4**
- [71] Rakesh Shrestha, Siqi Hu, Minghao Gou, Ziyuan Liu, and Ping Tan. A real world dataset for multi-view 3d reconstruction. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part VIII*, pages 56–73. Springer, 2022. **3**
- [72] Arjun Singh, James Sha, Karthik S Narayan, Tudor Achim, and Pieter Abbeel. BigBIRD: A large-scale 3D database of object instances. In *2014 IEEE international conference on robotics and automation (ICRA)*, pages 509–516. IEEE, 2014. **2, 3**
- [73] Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. SUN RGB-D: A RGB-D scene understanding benchmark suite. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 567–576, 2015. **2**
- [74] Xibin Song, Yuchao Dai, Dingfu Zhou, Liu Liu, Wei Li, Hongdong Li, and Ruigang Yang. Channel attention based iterative residual learning for depth map super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. **3**
- [75] Christoph Strecha, Wolfgang Von Hansen, Luc Van Gool, Pascal Fua, and Ulrich Thoennessen. On benchmarking camera calibration and multi-view stereo for high resolution imagery. In *2008 IEEE conference on computer vision and pattern recognition*, pages 1–8. Ieee, 2008. **2**
- [76] Jiaming Sun, Yiming Xie, Linghao Chen, Xiaowei Zhou, and Hujun Bao. Neuralrecon: Real-time coherent 3d reconstruction from monocular video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15598–15607, June 2021. **3**
- [77] Alex Teichman, Stephen Miller, and Sebastian Thrun. Unsupervised intrinsic calibration of depth sensors via slam. In *Robotics: Science and systems*, volume 248, page 3, 2013. **5**
- [78] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996. **13**

- [79] Oleg Voynov, Alexey Artemov, Vage Egiazarian, Alexander Notchenko, Gleb Bobrovskikh, Evgeny Burnaev, and Denis Zorin. Perceptual deep depth super-resolution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5653–5663, 2019. 3
- [80] Fangjinhua Wang, Silvano Galliani, Christoph Vogel, Pablo Speciale, and Marc Pollefeys. Patchmatchnet: Learned multi-view patchmatch stereo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [81] Peng Wang. Neus. <https://github.com/Totoro97/NeuS>. 11
- [82] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *NeurIPS*, 2021. 2, 6, 11
- [83] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4690–4699, 2021. 2
- [84] Oliver Wasenmüller, Marcel Meyer, and Didier Stricker. CoRBS: Comprehensive rgb-d benchmark for slam using kinect v2. In *IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, March 2016. 3
- [85] Silvan Weder. Routedfusion. <https://github.com/weders/RoutedFusion>. 11
- [86] Silvan Weder, Johannes Schonberger, Marc Pollefeys, and Martin R. Oswald. RoutedFusion: Learning Real-Time Depth Map Fusion. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4886–4896, Seattle, WA, USA, June 2020. IEEE. 3, 6, 11
- [87] Silvan Weder, Johannes L. Schonberger, Marc Pollefeys, and Martin R. Oswald. Neurfusion: Online depth fusion in latent space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3162–3172, June 2021. 3
- [88] Zizhuang Wei, Qingtian Zhu, Chen Min, Yisong Chen, and Guoping Wang. Aa-rmsnet: Adaptive aggregation recurrent multi-view stereo network. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6187–6196, 2021. 2
- [89] Thomas Whelan, Renato F Salas-Moreno, Ben Glocker, Andrew J Davison, and Stefan Leutenegger. ElasticFusion: Real-time dense SLAM and light source estimation. *The International Journal of Robotics Research*, 35(14):1697–1716, Dec. 2016. 3
- [90] Zhirong Wu, Shuran Song, Aditya Khosla, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shape modeling. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Boston, USA, June 2015. 3, 7, 11
- [91] Chuhua Xian, Kun Qian, Zitian Zhang, and Charlie C. L. Wang. Multi-Scale Progressive Fusion Learning for Depth Map Super-Resolution. *arXiv e-prints*, page arXiv:2011.11865, Nov. 2020. 3
- [92] Qingshan Xu. Acmp. <https://github.com/GhiXu/ACMP>. 11
- [93] Qingshan Xu and Wenbing Tao. Multi-scale geometric consistency guided multi-view stereo. *Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [94] Qingshan Xu and Wenbing Tao. Planar prior assisted patch-match multi-view stereo. *AAAI Conference on Artificial Intelligence (AAAI)*, 2020. 2, 6, 11
- [95] Yao Yao. Blendedmvg. <https://github.com/YoYo000/BlendedMVS#upgrade-to-blendedmvg>. 11
- [96] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. *European Conference on Computer Vision (ECCV)*, 2018. 1, 2, 11
- [97] Yao Yao, Zixin Luo, Shiwei Li, Jingyang Zhang, Yufan Ren, Lei Zhou, Tian Fang, and Long Quan. Blendedmvs: A large-scale dataset for generalized multi-view stereo networks. *Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 7, 11
- [98] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Basri Ronen, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. *Advances in Neural Information Processing Systems*, 33, 2020. 2
- [99] Igor Zacharov, Rinat Arslanov, Maksim Gunin, Daniil Stefonishin, Andrey Bykov, Sergey Pavlov, Oleg Panarin, Anton Maliutin, Sergey Rykovanov, and Maxim Fedorov. “zhores”-petaflops supercomputer for data-driven modeling, machine learning and artificial intelligence installed in skolkovo institute of science and technology. *Open Engineering*, 9(1):512–520, 2019. 7
- [100] Bernhard Zeisl and Marc Pollefeys. Structure-based auto-calibration of rgb-d sensors. In *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5076–5083. IEEE, 2016. 5
- [101] Jingyang Zhang. Vismvsnet. <https://github.com/jzhangbs/Vis-MVSNet>. 11
- [102] Jingyang Zhang, Yao Yao, Shiwei Li, Zixin Luo, and Tian Fang. Visibility-aware multi-view stereo network. *British Machine Vision Conference (BMVC)*, 2020. 2, 6, 11
- [103] Yinda Zhang and Thomas Funkhouser. Deep Depth Completion of a Single RGB-D Image. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 175–185, Salt Lake City, UT, USA, June 2018. IEEE. 3
- [104] Qian-Yi Zhou and Vladlen Koltun. Simultaneous localization and calibration: Self-calibration of consumer depth cameras. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 454–460, 2014. 5
- [105] Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Open3D: A modern library for 3D data processing. *arXiv:1801.09847*, 2018. 11