

CAD Priors for Accurate and Flexible Instance Reconstruction

Tolga Birdal^{1,2}

Slobodan Ilic^{1,2}

¹ Technische Universität München, Germany

² Siemens AG, Munich, Germany

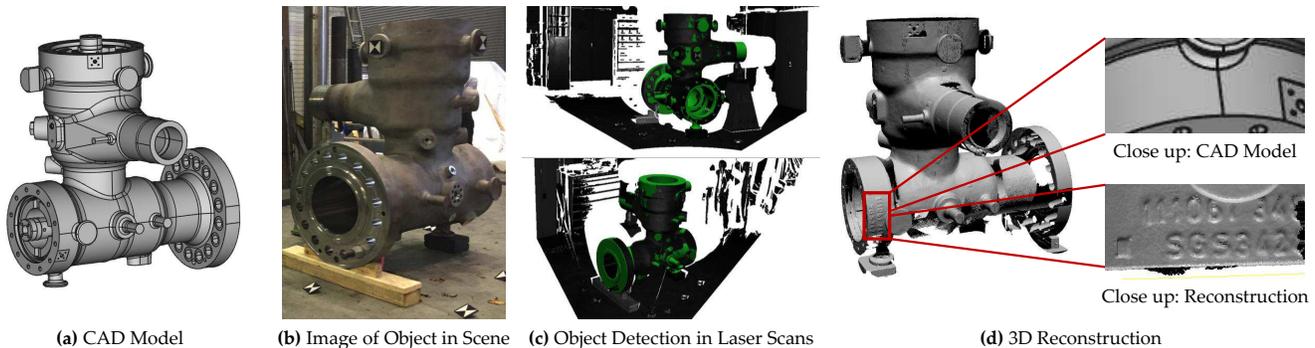


Figure 1: Our 3D reconstruction method. (a) Input 3D CAD model. (b) Image of the instance to reconstruct. (c) Detection of 3D model in point clouds. (d) Final reconstruction we obtain, with close-up comparisons to the nominal CAD prior.

Abstract

We present an efficient and automatic approach for accurate instance reconstruction of big 3D objects from multiple, unorganized and unstructured point clouds, in presence of dynamic clutter and occlusions. In contrast to conventional scanning, where the background is assumed to be rather static, we aim at handling dynamic clutter where the background drastically changes during object scanning. Currently, it is tedious to solve this problem with available methods unless the object of interest is first segmented out from the rest of the scene. We address the problem by assuming the availability of a prior CAD model, roughly resembling the object to be reconstructed. This assumption almost always holds in applications such as industrial inspection or reverse engineering. With aid of this prior acting as a proxy, we propose a fully enhanced pipeline, capable of automatically detecting and segmenting the object of interest from scenes and creating a pose graph, online, with linear complexity. This allows initial scan alignment to the CAD model space, which is then refined without the CAD constraint to fully recover a high fidelity 3D reconstruction, accurate up to the sensor noise level. We also contribute a novel object detection method, local implicit shape models (LISM) and give a fast verification scheme. We evaluate our method on multiple datasets, demonstrating the ability to accurately reconstruct objects from small sizes up to $125m^3$.

1. Introduction

3D reconstruction involves the recovery of digitized object or scene using a 2D/3D sensor, typically through multiple acquisition steps. From reverse engineering to industrial inspection, its applications are plentiful. Due to such wide use, from the early days of vision, it attracted significant attention of both research community and industry [24, 7].

Despite the huge demand, many marker-free approaches based solely on 3D data either involve acquisition of ordered scans [31, 21], or follow the de-facto standard pipeline [19] in case of unordered scans. The former suffers from the requirements of redundant depth capture with large overlap and scenes with very little clutter or occlusions. Due to the volumetric nature of scan fusion, such techniques also do not scale well to large objects while retaining high precision. The latter exploits 3D keypoint matching of all scans to one another, alleviating the order constraint. Thanks to 3D descriptors, it could well operate in full 3D. Yet, matching of scans to each other is an $O(N^2)$ problem and prevents the methods from scaling to an arbitrary number of scans. In addition, neither of those can handle scenes with extensive dynamic clutter or occlusions.

Nowadays, with the capability of collecting high quality, large scale and big data, it is critical to offer automated solutions for providing highly accurate reconstructions regardless of the acquisition scenarios. In this paper, we tackle the problem of 3D instance reconstruction from a handful of unorganized point clouds, where the object of interest

is large in size, texture-less, surrounded by significant dynamic background clutter and is viewed under occlusions. Our method can handle scenes between which no single transformation exists, i.e. the same objects appear in different locations, such as the one in Fig. 2. We also do not impose any constraints on the order of acquisition. To solve all of these problems simultaneously, we make use of the reasonable assumption that a rough, nominal 3D CAD model prior of the object-to-reconstruct is available beforehand and propose a novel reconstruction pipeline. Such assumption of a nominal prior is valid for many applications especially in industry, where the objective is to compare the reconstruction to the designed CAD model. Even for the cases where this model does not exist, one could always generate a rough, inaccurate mesh model with existing methods, e.g. KinectFusion [31], to act as a prior. Note that, due to manufacturing errors, sensor noise, damages and environmental factors, physical instances deviate significantly from the CAD models and the end-goal is an automatic algorithm to accurately recover the particular instance of the model. With the introduction of the prior model, we refactor the standard 3D reconstruction procedures via multi-fold contributions. We replace the scan-to-scan matching with model-to-scan matching resulting in absolute poses for each camera. Unlike the case in object instance detection where false positives (FP) are tolerable, object instance reconstruction is easily jeopardized by the inclusion of a single FP. Therefore, one of our goals is to suppress FP, even at the expense of some true negatives. To achieve this, we contribute a probabilistic Local Implicit Shape Model (LISM) formulation for the object instance detection and pose estimation, accompanied by a rigorous hypotheses verification to reject all wrong pose candidates. This matching is followed by automatically segmenting out the points belonging to the object and transforming them back onto the model coordinate frame. Doing this for multiple views results in roughly aligned partial scans in the CAD space. To fully recover the exact object, the CAD prior is then discarded (as it might cause undesired bias) and a global multi-view refinement is conducted only to optimize the camera poses.

In global scan registration, the creation of a *pose graph* indicating which scans can be registered to each other is required. The typical complexity of obtaining this graph is $O(N^2)$, where all scans are matched to one another. We profit from the CAD prior and contribute by automatically computing this graph, in which only cameras sharing significant view overlap are linked. This reduces the complexity to $O(N)$, and robustifies the whole pipeline.

The entire procedure is made efficient so that large scans are handled in reasonable time. Exhaustive evaluations demonstrate high accuracy regardless of the size of the objects, clutter and occlusions.¹

¹Our suppl. video is under: https://youtu.be/KPA_8BNuOvg

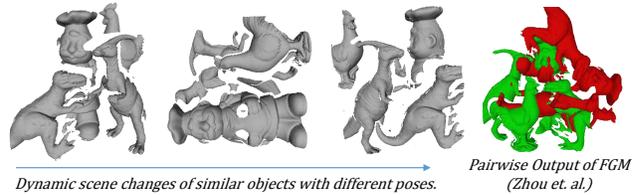


Figure 2: Dynamic Clutter and Occlusions on Mian Dataset [28]: There is change in relative locations of the objects between scenes. This easily fools the modern global registration algorithms [46]. We operate on the object level and circumvents this problem.

2. Prior Art

Arguably, the most wide-spread 3D reconstruction methods are KinectFusion[31] and its derivatives [9, 21, 37]. These methods have been successful in reconstructing small isolated objects, but their application is not immediate when the size increases, or clutter and occlusions are introduced [32]. Due to extensive usage of signed-distance fields, they are bound to depth images and a sequential acquisition.

Abundance of works exists in multi-view global alignment from multiple 3D unorganized point clouds [2, 10, 13, 12, 43, 47, 38, 17, 40, 22, 8, 12, 46]. These methods assume the scans to be roughly initialized and reasonably well-segmented. They, in general, handle slight synthetic noise well enough, but they do not deal with cluttered and occluded data. Another track of works try to overcome the aforementioned constraints by using keypoint detectors and matching descriptors in 3D scans. These methods operate on a subset of points during matching. One of the pioneering works proposing a feature agnostic, automatic and constraint-free algorithm is the graph based in-hand scanning from D. Huber and M. Hebert [19]. The authors set a baseline for this family of methods. Novatnack and Nishino [33] developed a scale dependent descriptor for better initialization and fused it with [19] to assess the power of their descriptor. Yet both of these relied on range image data. Mian et. al. [28] proposed a tensor feature and a hashing framework operating on meshes. Fantoni et al. [14] uses 3D keypoint matching as an initial stage of multiview alignment to bring the scans to a rough alignment. Zhu et. al. [47, 49] as well as Liu and Yonghuai [27] use genetic algorithms to discover the matching scans and use this in global alignment context. These stochastic schemes are correspondence free but slow. Similar to [28], Zhu et. al. [48] devises a local-to-global minimum spanning tree method to align the scans. A majority of these automatic alignment procedures suffer from increased worst case complexity of $O(N^2)$, where N is the number of views. Moreover, since there is no integrated segmentation, the registration procedure cannot handle clutter and occlusions.

Use of CAD models in reconstruction is not novel by itself. Savarese et. al. [45] enrich the multiview reconstruction from 2D images with a CAD prior. Guney and Geiger

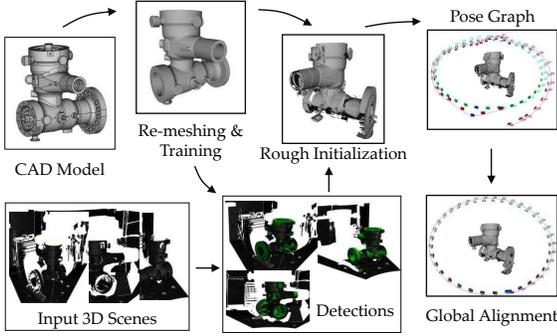


Figure 3: Proposed 3D reconstruction pipeline: Prior CAD model is trained to create the model representation. Input scenes are then parsed for the model pose. Pose estimates initialize a rough reconstruction, with segmentation and automatically computed pose graph. This is then further refined to the full reconstruction.

[18] use object knowledge to resolve the stereo ambiguities. Bidal et. al. use models in triangulation by registration [4]. Other works [26, 44, 50] use CAD prior to detect generic object classes. Our approach differs from all of those in the sense that we use proxy instance priors for initial alignment of scans and then operate directly on 3D points.

Salas et. al. [36] propose SLAM++ using object priors to constrain a SLAM system. Our work is differentiated from theirs in the sense that we perform instance reconstruction using the CAD prior, and not SLAM. In our setting, the background as well as the object is allowed to change dynamically between different acquisitions.

3. Method

Given a set of unstructured and unordered 3D scenes $\{\mathbf{S}_i\} \in \mathbf{S}$, we seek to find a set of transformations $\{\mathbf{T}_i\} \in SE(3)$ so as to stitch and reconstruct a global model \mathbf{S}_G :

$$\mathbf{S}_G = \bigcup_{i=1}^N (\mathbf{T}_i^0 \circ \mathbf{S}_i) \quad (1)$$

$\mathbf{T} \circ \mathbf{S}$ applies transformation \mathbf{T} to the scene \mathbf{S} . In this setting, both the transformations $\{\mathbf{T}_i^0\} \in SE(3)$ and the global model \mathbf{S}_G are unknown and we do not assume known initialization. Due to the lack of a common reference frame and apriori information about $\{\mathbf{S}_i\}$, obtaining the set $\{\mathbf{T}_i^0\}$ typically requires $O(N^2)$ worst case complexity, where all the scene clouds are matched to one another to obtain the relative transformations aligning them. To better condition the problem and reduce its complexity, we introduce the supervision of a CAD proxy \mathbf{M} in form of a mesh model and re-write Eq. (1):

$$\mathbf{S}_G = \bigcup_{i=1}^N (\mathbf{T}_i^M \circ (\mathbf{S}_i | \mathbf{M})) \quad (2)$$

where $\mathbf{T}_i^M \in SE(3)$ is the transformation from the scene to the model space, such that the segmented scene points

$(\mathbf{S}_i | \mathbf{M})$ come to the best agreement. To estimate $\{\mathbf{T}_i^M\}$, we follow a two stage technique. First, a rather approximate estimate $\{\tilde{\mathbf{T}}_i^M\}$ is made by matching the CAD model to a single scene. Note that this time, the set $\{\tilde{\mathbf{T}}_i^M\}$ can be computed in $O(N)$, since it only requires CAD to scan alignment. However, because scene clouds suffer from partial visibility, noise and deviations w.r.t. the CAD model, the discovery of the pose of the model in the scans provides only rough initial transformations to the model coordinate frame. For this reason, as the final stage, the CAD prior is disregarded and the scans are globally refined, simultaneously. This lets us reconstruct configurations deviating significantly from the CAD prior.

Our procedure of multi-view refinement is similar to [14], where a global scheme for scan alignment is employed. Let $\mathbf{S}_1, \dots, \mathbf{S}_M$ be the set of scans that are to be brought in alignment. To generalize and formalize the notation for registrations of all point clouds to each other, we maintain a directed pose graph in form of an adjacency matrix $A \in \{0, 1\}^{M \times M}$, such that $A(h, k) = 1$ iff cloud \mathbf{S}_h can be registered to cloud \mathbf{S}_k . Let $\theta = (\theta_1, \dots, \theta_M)$ be the absolute camera poses of each view. The alignment error between two clouds \mathbf{S}_h and \mathbf{S}_k then reads:

$$E(\theta_h, \theta_k) = A(h, k) \sum_{i=1}^{N_h} \rho \left(\|d(\theta_h \circ p_i^h, \theta_k \circ q_i^h)\|^2 \right) \quad (3)$$

where $\{p_i^h \rightarrow q_i^h\}$ are the N_h closest point correspondences obtained from the clouds \mathbf{S}_h and \mathbf{S}_k . The point-to-plane distance $d(\cdot, \cdot)$ is defined to be:

$$d(p_i, q_i) = (\mathbf{R}p_i + \mathbf{t} - q_i)^T n_i^q \quad (4)$$

with n_i^q referring to the normal associated to point q_i . \mathbf{R} and \mathbf{t} are the components of the pair transformation $\{ \mathbf{R} \in SO(3), \mathbf{t} \in \mathbb{R}^3 \}$. The overall alignment error, which we want to minimize at this stage, is obtained by summing up the contribution of every pair of overlapping views:

$$E(\theta) = \sum_{h=1}^M \sum_{k=1}^M A(h, k) \sum_{i=1}^{N_h} \rho \left(\|d(\theta_h \circ p_i^h, \theta_k \circ q_i^h)\|^2 \right) \quad (5)$$

where ρ is the robust estimator. The final absolute poses are the result of the minimization $(\theta_1, \dots, \theta_M) = \text{argmin}_{\theta} (E)$, and align the M clouds in a least squares sense. In contrast to the pairwise registration error in Eq. (3), which has closed form solution for the relative transformation θ , there are no closed form solutions in the multiview setting. Therefore, we use a non-linear optimization procedure, Levenberg Marquardt. The rotations are parameterized with angle-axis representation $(\mathbf{w} \in \mathbb{R}^3, \phi)$. We constrain the frame with the highest number of points found on the CAD model to be static and update (solve for) the rest of the poses. In practice, this leads to faster convergence.

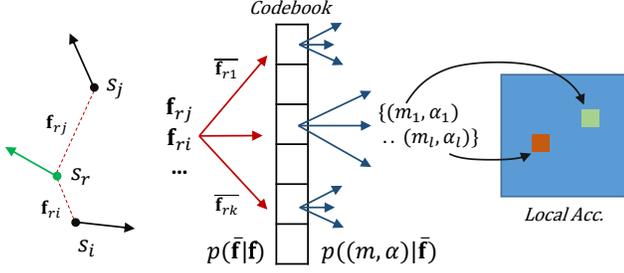


Figure 4: Local Implicit Voting: Given multiple scene point pairs, tied to a common reference s_r , we generate features \mathbf{f}_{r_i} , activating different codebook buckets (middle). Each bucket casts votes for multiple (m, α) pairs in the local voting space of s_r .

Note that, in contrast to the methods that exploit pairwise registration, our poses are absolute and do not suffer from drifts or tracking artifacts. We also do not require a conversion from relative poses to absolute ones, which are usually obtained by the computation of a minimum spanning tree or shortest paths over the pose graph [16, 19]. This property eases the implementation and reduces errors, that are to be encountered in usual heuristics.

Due to the accuracy requirements, unlike [14], we omit using distance transforms at this stage. We rather use speeded up KD-Trees to achieve exact nearest neighbors [30]. Since we optimize over the poses, and not over 3D points, the trees are built only once in the beginning and all closest point computations are done in the local coordinate frame of the view of interest. This is important for efficiency. For reasons of accuracy, we use analytical Jacobians. As cloud sizes become large, this optimization exhibits significant computational costs. This is why, a priori sampling plays a huge role, where we use $\approx 20k$ to $30k$ points per scan, distributed evenly in space.

To summarize, our key contribution lies in obtaining $\{\mathbf{T}_i^M\}$ in a robust, efficient and accurate manner. We will now show how to compute the rough alignment $\{\tilde{\mathbf{T}}_i^M\}$ and the pose graph (adjacency matrix) A .

3.1. Locally Implicit Models for Estimating $\{\tilde{\mathbf{T}}_i^M\}$

While using any method, which is capable of handling 3D points, e.g. [20, 42, 28, 41, 34], is possible, we mainly follow the Geometric Hashing of Drost et al. [11] and Birdal and Ilic [6] due to efficiency and robustness to clutter and occlusions. Yet, we introduce a more effective probabilistic formulation, inspired by the implicit shape models [23].

Model Description In the first stage, we generate a pose invariant *codebook* encoding all possible semi-global structures that could be found on the CAD model. We represent this semi-global geometry via simple point pair features (PPF) of oriented point pairs $(\mathbf{m}_i, \mathbf{m}_j)$:

$$\mathbf{f}_{ij} = (\|\mathbf{d}\|_2, \angle(\mathbf{n}_i, \mathbf{d}), \angle(\mathbf{n}_j, \mathbf{d}), \angle(\mathbf{n}_i, \mathbf{n}_j)) \quad (6)$$

where $\mathbf{d} = \mathbf{m}_i - \mathbf{m}_j$, \mathbf{n}_i and \mathbf{n}_j are the surface normals at points \mathbf{m}_i and \mathbf{m}_j . $\angle(\cdot, \cdot)$ is the angle operator. The complete set of such features $\mathbf{F} = \{\mathbf{f}_{ij}\}$ for the prior CAD model is collected and quantized to generate the codebook: $\bar{\mathbf{F}}$. We use our codebook to relate a feature \mathbf{f} (key) to a set of oriented references points $\{(\mathbf{m}_i, \mathbf{m}_j)\}$ (stored in buckets) and build the global model description as an inverted file i.e. a hashtable \mathbf{H} . Thus, each bucket in the codebook contains self-similar point pairs extracted from the CAD model. Whenever a pair from the scene is matched to one in the model, their normals at the reference points are aligned. Then, the full pose of the object can be obtained once the rotation angle α around the normal is known. This can be done by aligning Local Coordinate Frames (LCF) constructed from matched pairs. Thus, instead of storing the full PPF, we store only this local parameterization $\{\mathbf{m}_r, \alpha\}$ composed of the model reference point \mathbf{m}_r and rotation angle α . A pair correspondence resolves the full 6DOF pose and what is left is to retrieve the matching pairs $\{s_r, s_i\}$ and $\{\mathbf{m}_r, \mathbf{m}_i\}$. We now give a novel way to do this.

Probabilistic Formulation During detection, a new point cloud scene \mathbf{S} is encountered and downsampled to a set of points $\mathbf{S}_D = \{s_r\}$, some of which are assumed to lie on the object. The sampling also enforces spatial uniformity (see our suppl. material). We fix a reference point s_r and pair it with all the other samples $\{s_i\}$. Each pair makes up a PPF \mathbf{f}_{r_i} . The original method [11] associates \mathbf{f}_{r_i} to a unique key and can not account for the quantization errors that inevitably happen due to the noise. To circumvent these quantization artifacts, resulting from the hard assignment in [11], we quantize \mathbf{f}_{r_i} to K different bins ($K > 1$), activating different codebook entries as in ISM. This soft quantization results in possibly matching buckets $\bar{\mathbf{F}}_{r_i} = \{\bar{\mathbf{f}}_1 \dots \bar{\mathbf{f}}_K\}$. $\bar{\mathbf{F}}_{r_i}$ indexes the buckets of \mathbf{H} , with weights $p(\bar{\mathbf{f}}_k | \mathbf{f}_{r_i})$. For each matching bucket, we collect the valid interpretations $p(m, \alpha | \bar{\mathbf{f}})$, inversely proportional to the size of the bucket N_b , denoting the probabilities of particular pose configuration, given the quantized feature. Formally:

$$p(\mathbf{m}, \alpha | s_r, s_i) = p(\mathbf{m}, \alpha | \mathbf{f}) \quad (7)$$

$$= \sum_k p(\mathbf{m}, \alpha | \bar{\mathbf{f}}_k) p(\bar{\mathbf{f}}_k | \mathbf{f}) \quad (8)$$

with $p(\bar{\mathbf{f}}_k | \mathbf{f}) = \frac{1}{K}$ and $p(m, \alpha | \bar{\mathbf{f}}) = \frac{1}{N_b}$ being uniformly distributed. This probability is actually the prior on the PPF of the particular object and can be computed differently accounting for the nature of the object geometry using a suited distribution. At this point, the gathered pair representations for a particular scene reference point are sufficient to recover for the object pose. However, due to outliers, some of these matches will be erroneous. Therefore, a 2D voting scheme is employed, locally for each scene reference point s_r . The voting space is composed of the alignment of the

LCF α as well as the model point correspondence \mathbf{m} :

$$V(\mathbf{m}, \alpha) = \sum_i p(\mathbf{m}, \alpha | \mathbf{s}_r, \mathbf{s}_i) \quad (9)$$

For each \mathbf{s}_r , there is a voting space $V_r(\mathbf{m}, \alpha)$, from which the best alignment is extracted as:

$$(\mathbf{m}_r^*, \alpha_r^*) = \operatorname{argmax}_{\mathbf{m}, \alpha} V_r(\mathbf{m}, \alpha) \quad (10)$$

Each such $(\mathbf{m}_r^*, \alpha_r^*)$ corresponds to a pose hypothesis. This is similar to performing Generalized Hough Transform (GHT) on reference point level locally and is the reason why we attribute the name Local ISM to our method. After all pose hypotheses are extracted, as the maxima in the local spaces, the poses are clustered together to assemble the final consensus, further boosting the final confidence.

Hypotheses Verification and Rejection Devised matching theoretically generates a pose hypothesis for each scene reference point, which is assumed to be found on the model. There are typically $\sim 400 - 1000$ such points, reducing to 50 poses after the clustering, where the close-by poses are grouped together and averaged. Still, as many hypotheses as the number of clusters remain to be verified and the best pose is expected to be refined. In our problem of instance reconstruction it is critical that no false positive pose hypotheses survives. For this reason we introduce a rigorous hypothesis verification scheme. The effective verification requires fine registration, while efficient registration requires as few poses as possible. This creates a chicken and egg problem. We address this issue via a multi-level registration approach. In the first stage, sparsely sampled scan points are finely registered to the model using the efficient LM-ICP [15] variant of Iterative Closest Point (ICP) registration [3]. We also build a 3D distance transform for fast nearest neighbor access. Our sparse LM-ICP requires only 1ms per hypothesis, allowing us to verify all the hypotheses. We define the hypothesis score to be:

$$\Xi(\theta_i) = \frac{1}{N_M} \sum_j \begin{cases} 1, & \|\theta_i^{-1} \circ \mathbf{m}_j - \mathbf{s}_k(j)\| < \tau_\theta \\ 0, & \text{otherwise} \end{cases} \quad (11)$$

where θ_i is the pose hypothesis and $\mathbf{s}_k(j)$ the closest sampled scene point to transformed model point $\theta_i^{-1} \circ \mathbf{m}_j$. Intuitively, this score reflects the percentage of visible model points. The surviving poses are then sorted, taken to the next level and densely refined. This coarse to fine scheme is repeated for 3 levels of the pyramid. Finally, a dense registration is performed to accurately obtain the final pose.

Until this stage the surface normals are excluded from the fine registration process. We do this intentionally, to use them as a verification tool. Following registration, we check the surface consistency between the scene and the model.

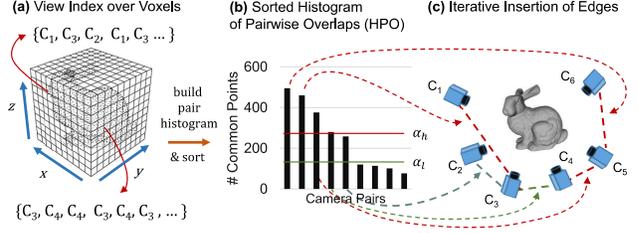


Figure 5: Pose graph computation. See text for details.

To do so for each scene point, the surface normal of the closest model point is retrieved. A scan is only accepted if a majority of the normals agree with the model. While this procedure can result in potentially good detections being removed (due to scene deviations), it does not allow false positives to survive as shown in Sec. 4.

3.2. Computing Pose Graph A and Live Feedback

Any global optimization algorithm requires an adjacency graph $G = (V, E)$, which encodes the existence of overlap between camera views. The nodes of this sparse graph contain the cameras $V = \{C_1..C_N\}$, whereas an edge E_{ij} is only created between nodes (C_i, C_j) if they share significant overlap. An absolute pose \mathbf{T}_i is associated to each node and a relative pose \mathbf{T}_{ij} is to each edge. Traditionally, this requires pair-wise overlap computation between all cameras. While a naive approach would involve linking the cameras, whose centers are found to be close, this is by no means a guarantee for shared overlap. Therefore, we present a more accurate approach, without sacrificing efficiency, thanks to the availability of the CAD model.

Consider the voxel grid index \mathbf{D} of model \mathbf{M} as in Fig. 5(a). Each segmented scene point $\mathbf{s}'_i \in \mathbf{T}_i^{-1} \circ \mathbf{S}_i$ is mapped to a voxel D_k , which stores a set of cameras $\{C_i\}$ observing it. Whenever the point \mathbf{m}_k belonging to the voxel D_k is visible in the camera C_j , this camera is added to the list of cameras seeing that model point. Each list stores unique camera indices. From that, we compute the histogram of pairwise overlaps (HPO) as shown in Fig. 5(b).

While all the possible edges are now generated (as the bins in HPO), it is not recommended to use all these in the multiview alignment i.e. the overlap might be little, causing a negative impact. Instead, we adopt an iterative algorithm, similar to hysteresis thresholding. First, HPO is sorted with decreasing overlap (Fig. 5(c)). Next, two thresholds α_l and α_h are defined. All pairs with overlap less than α_l are discarded. All cameras with overlap larger than α_h are immediately linked and edges are constructed in the graph. If, at this stage, the graph is not connected, we start inserting edges from the remaining bins of HPO into A until either the connectivity or the threshold α_l is reached. This is illustrated in Figures 5(c) and 5(d). If the final graph is still not connected, we use the largest connected sub-graph, to ensure optimize-ability. For efficient online update, a mod-

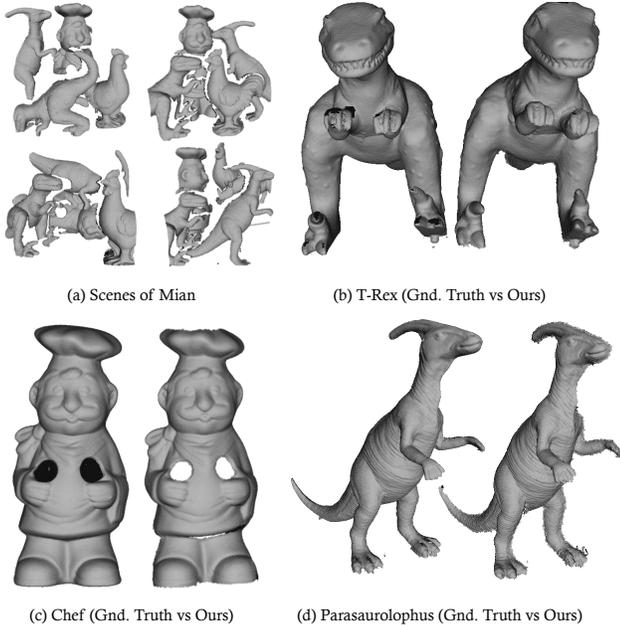


Figure 6: Results on Mian Dataset. (a) Subset of scenes from the dataset. (b, c, d) The ground truth models (left) and our reconstruction (right) for three objects.

ified union-find data structure is used to store the graph and dynamically insert edges when new views are encountered. Unlike quadratic complexity of the standard pose graph creation methods, ours has linear complexity.

Live Feedback Due to the connected-ness of pose graph, our method is able to keep track of the overlap between all the point clouds, at all times, informing the user whenever graph disconnects or overlap is small. The complement of the already reconstructed part reveals the unscanned region, which is also fed back to the operator. Incoming scans directly propagate and form links in the pose graph, allowing online response to the user’s actions.

4. Experimental Evaluation

We evaluate our method against a set of real datasets acquired by laser scanners and structured light sensors. The CAD models we work with might contain uneven distribution of vertices or inner geometry. We always eliminate the inner structure by thresholding the ambient occlusion values [29] before the models are re-meshed [25]. At detection time, a relative model and scene sampling distance of $d = \tau \text{diam}(M)$ is used, where $0.05 \geq \tau \geq 0.025$ depending on the object. We also adjust another threshold on the distance to consider a scene point to be on the model based on the sensor quality. For accurate scanners we use 1.5mm , while for less accurate ones 0.5cm . This does not affect the segmentation, but the hypothesis verification.

Table 1: Reconstruction results on Mian Dataset (in mm). Each object is compared to the model provided by [28] using [1].

Model	w/o Opt.	with Opt.	# Scans	Clutter
Chef	2.90 ± 2.40	1.07 ± 0.65	22	0.58 ± 0.11
Chicken	1.71 ± 1.60	0.33 ± 0.24	29	0.61 ± 0.12
Para.	2.52 ± 2.00	0.41 ± 0.30	12	0.24 ± 0.20
T-rex	2.36 ± 2.08	0.88 ± 0.62	27	0.14 ± 0.22

Mian Dataset We first compare the reconstruction quality on Mian Dataset [28]. This dataset includes 50 laser scanned point clouds of 4 complete 3D objects, with varying occlusion and clutter. The objects change locations from scan to scan, creating dynamic scenes. The clutter and background also varies as the objects appear together with other different ones in each scene. We quantify this dynamic clutter by relating it to the provided occlusion values:

$$\text{Clutter} = 1 - \frac{(\text{Model Surface Area}) * (1 - \text{Occlusion})}{(\text{Scene Surface Area})} \quad (12)$$

and provide it in Table 1 for each object. The models present in the scenes are provided by [28] to act as ground truth. We do not perform any prior operation to the scenes such as segmentation or post-processing except meshing via SSD [39]. For Parasaurolophus and Chicken objects, the pose graph becomes disconnected and therefore, we optimize individually the two sub-components and record the mean. We also report the number of scans in which the model is detected and verified. Not every model is visible in every scan. In the end, our mean accuracy is well below a millimeter, where the used sensor, Minolta Vivid 910 scanner, reports an ideal accuracy of $\sim 0.5\text{mm}$. We are also not aware of any other works, reporting reconstruction results on such datasets. Fig. 6 visualizes our outcome, and Table 1 shows our reconstruction accuracy both prior to and after the optimization. While our error is quantitatively small, the qualitative comparison also yields a pleasing result, sometimes being superior even to the original model.

The detection performance of a basic variant of our method has already been proven to be robust on this dataset [11]. Figures 7(a) and 7(b) provide PR-curves for LISM and the hypothesis verification. Note that although LISM already performs well, our verification clearly improves the distinction between a match vs false positive. Using a simple threshold, we could obtain 100% precision without sacrificing the recall. Thus, our score threshold, combined with the normal consistency check manages to reject all false hypotheses, at the expense of rejecting a small amount of TP.

Toy Objects Dataset Since our objective is to assess the fidelity of the reconstruction to the CAD model, we opt to use the objects from the 3D printed dataset [37]: *Leopard*, *Teddy* and *Bunny* and *Tank* (See Fig. 8). The diameters of objects vary in the range of $15 - 30\text{cm}$. The print accuracy

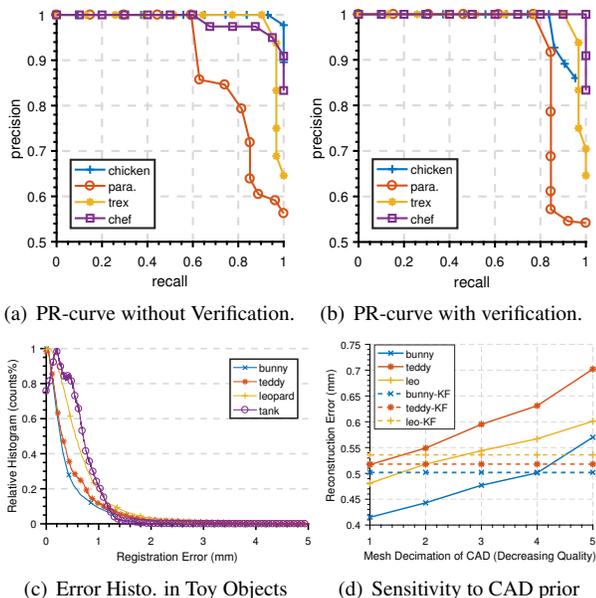


Figure 7: Performance of LISM and verification on Mian dataset.

is up to 50μ (micro-meters), well sufficient for consideration as ground truth. To capture the real scenes, a homebrew, high accuracy phase-shift sensor, delivering $<0.4\text{mm}$ point accuracy is chosen. We sample up to 10 scans per object, taken out of a 100 frame sequence. To disrupt the acquisition order, we randomly shuffle this subset and apply our reconstruction algorithm. Next, we compute the CAD-to-reconstruction distances in CloudCompare [1]. We do not explicitly register our reconstruction to CAD model, because we already end up on model coordinate frame (Having the result in the CAD space is a side benefit of our approach). Moreover, we use the original 100-frame, ordered sequence as an input to standard reconstruction pipelines such as Kinect Fusion [35], Kehl et. al. [21] (also uses color) and Slavcheva et. al. [37] all of which require a temporally ordered set of frames, with a large inter-frame overlap. All of these algorithms take depth image as input, whereas ours uses the unstructured 3D data and the model.

Our results on this dataset are shown in Tab. 2 (*Ours*) when original CAD is used. We also report the results when KinectFusion (KF) prior is used (*Ours-KF*). The individual error distribution of the objects are shown in Fig. 7(c).

Table 2: Reconstruction errors on toy objects dataset (mm).

	Leo	Teddy	Bunny	Tank
KinFU	1.785 ± 1.299	0.998 ± 0.807	0.664 ± 0.654	1.390 ± 1.315
Kehl	1.018 ± 1.378	1.028 ± 0.892	0.838 ± 0.860	1.573 ± 2.250
Sdf2Sdf	0.652 ± 0.614	0.910 ± 0.584	0.541 ± 0.436	0.466 ± 0.416
Ours	0.481 ± 0.519	0.517 ± 0.572	0.415 ± 0.501	0.451 ± 0.322
Ours-KF	0.536 ± 0.411	0.519 ± 0.582	0.502 ± 0.529	0.468 ± 0.474
Ours-CO	0.651 ± 0.628	0.544 ± 0.601	0.698 ± 0.506	0.475 ± 0.433

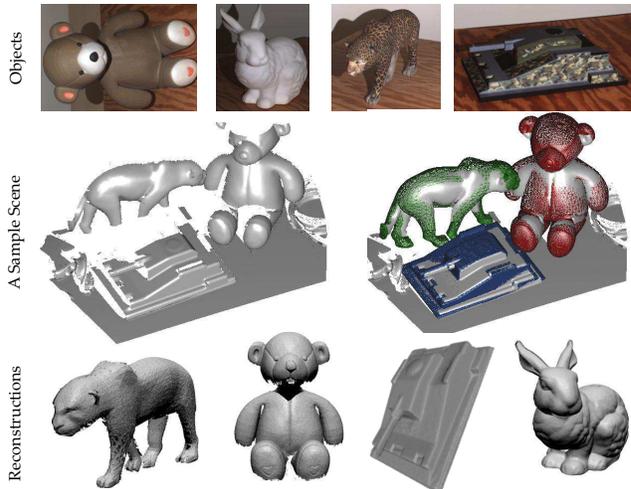


Figure 8: Qualitative results on toy objects. First row: Real images of objects; Second row: A sample scene and detections visualized; Third Row: Our results.

Because our method does not suffer from drift and computes absolute poses all the time, although we use 10 times less scans, we are still 2-4 times more accurate than conventional methods. This also shows that our method could retain the good accuracy of the sensor.

Next, we augment this dataset with further scenes of the same objects, such that clutter and occlusions are present. Some shots are shown in Fig. 8 (mid-row). Our reconstruction accuracy (*Ours-CO*) is shown in Tab. 2 for different objects. These results are still better than or close to Sdf2Sdf [37]. Due to inclusion of some outliers, our results get slightly worse than the one in no clutter, yet they are still acceptable. However, none of the other approaches can run on this new set due to the existence of significant outliers.

In a further experiment, we gradually decimate the toy models down to a mesh of ≈ 500 vertices. We exclude tank as the decimation has little effect on the planarities. As shown in Fig. 7(d), even though the CAD prior gets very crude, we are still able to achieve a reasonable reconstruction, as long as the CAD model is still detectable in the scenes. Note that, results of KF prior is plotted in dashes as it is also a form of rough mesh approximation. Furthermore, Fig. 9 visually compares our reconstruction to the state of the art on the tank object. Because we do not use smoothing voxel representations (such as SDF), our method is much better at preserving sharp features at the model edges.

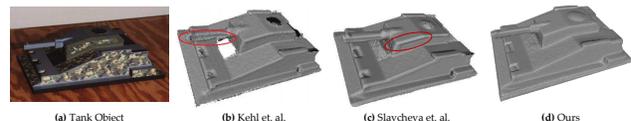


Figure 9: Visual comparisons on Tank object. Note the ability of our method in preserving sharp features.

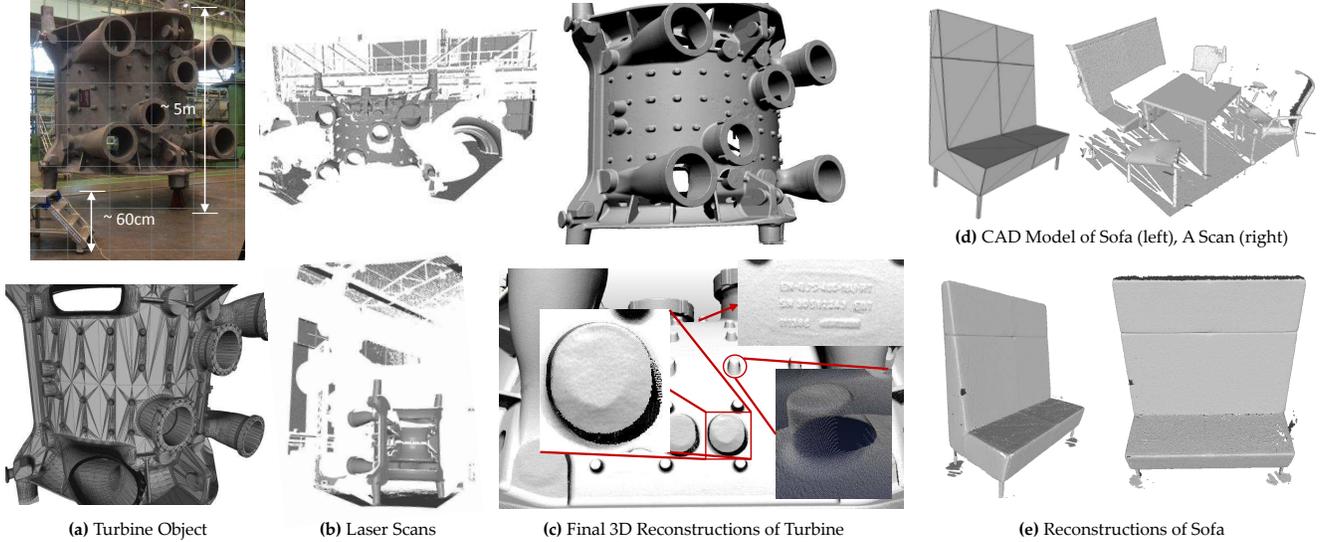


Figure 10: The reconstruction of Turbine(a) in captured cluttered scans(b) is presented in (c). Results in Sofa are shown in (d,e).

Table 3: Object Information, Average reconstruction errors w.r.t. Photogrammetry (in mm) and Timings.

Object	Scanner	Scan Res.	Obj. Size	No Scans	No PG Images	PG vs CAD	Surphaser	Our Accuracy	Detect	Verify	Refine
Ventil	Surph.	0.3 mm	$8m^3$	8	180	1.3cm	3.6 ± 3.3	2.2 ± 0.4	3.10s	0.27s	112.94s
Turbine	Surph.	0.4mm	$125m^3$	10	180	3.4cm	-	2.5 ± 1.3	3.72s	0.54s	126.13s
Sofa	Str.Light	1mm	$1.7m^3$	6	68	0.85mm	-	1.4 ± 1.2	1.44s	0.31s	68.82s

Dataset of Large Objects Finally, we apply our pipeline to quality inspection of real gas turbine casings and large objects. In this real scenario, CAD models come directly from the manufacturer. Due to space constraints, we summarize the data modality in Tab.3. The manufactured parts deviate significantly from the ideal model due to manufacturing and we scan them in the production environment within clutter and occlusions. With such large objects and little resemblance of the CAD prior, obtaining ground truth becomes a challenging task. Thus, we use a photogrammetry (PG) system [5, 4] to collect a sparse set of scene points, by attaching markers on the objects. We capture many images (see Tab.3) from different angles and run Linearis 3D software for bundle adjustment to obtain sparse ground truth. For Ventil object, we also use Surphaser software for reconstruction using external markers. Both Surphaser and our outputs are compared to the PG data in Tab. 3. The mean errors are obtained by CloudCompare [1]. We also provide running times of the individual stages. As seen, our accuracy outperforms an industry standard solution, Surphaser Software, by a margin of 38% on Ventil object. The performance in objects of varying sizes indicate that our reconstruction method is applicable from small to large scale while maintaining repeatability. Fig. 10 presents further qualitative results on our reconstruction of the Turbine and Sofa objects. Please consult the supplementary

material for more evaluations.

Limitations Due to the nature of PPF matching, our approach requires objects with rich geometry. Symmetric objects are also problematic due to ambiguity in pose estimation. Last but not least, currently, there is no mechanism to handle mis-detections. Yet, mis-detections are hardly a problem when the score threshold is reasonably high. This way, we detect in less scenes but avoid mistakes.

5. Conclusions

We proposed *reconstruction-via-detection* framework, as an alternative perspective to robust 3D instance reconstruction from unconstrained point cloud scans. Our framework integrates probabilistic object detection, hypothesis verification, pose graph construction and multi-view optimization. Such a scheme allowed us to deal with problems of dynamic clutter, occlusion and object segmentation. Moreover, the computational cost is reduced, due to model-to-scan alignment. To the best of our knowledge, this is the first method, capable of reconstructing instances within clutter and occlusions, without explicit segmentation.

As a future direction, we like to take care of confusions stemming from rotational symmetries by optimizing over possible global alignments of the scans. We also plan to extend our method to robotics via next-best view prediction.

References

- [1] Cloudcompare - 3d point cloud and mesh processing software. <http://www.danielgm.net/cc/>, 2016. 6, 7, 8
- [2] F. Arrigoni, B. Rossi, and A. Fusiello. Global registration of 3d point sets via lrs decomposition. In *European Conference on Computer Vision*, pages 489–504. Springer, 2016. 2
- [3] P. J. Besl and N. D. McKay. Method for registration of 3-d shapes. In *Robotics-DL tentative*, pages 586–606. International Society for Optics and Photonics, 1992. 5
- [4] T. Birdal, E. Bala, T. Eren, and S. Ilic. Online inspection of 3d parts via a locally overlapping camera network. In *Applications of Computer Vision (WACV), 2016 IEEE Winter Conference on*, pages 1–10. IEEE, 2016. 3, 8
- [5] T. Birdal, I. Dobryden, and S. Ilic. X-tag: A fiducial tag for flexible and accurate bundle adjustment. In *3D Vision (3DV), 2016 Fourth International Conference on*, pages 556–564. IEEE, 2016. 8
- [6] T. Birdal and S. Ilic. Point pair features based object detection and pose estimation revisited. In *3D Vision (3DV), 2015 2nd International Conference on*, pages 527–535. IEEE, 2015. 4
- [7] G. Blais and M. D. Levine. Registering multiview range data to create 3d computer objects. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 17(8):820–824, Aug 1995. 1
- [8] F. Bonarrigo and A. Signoroni. An enhanced ‘optimization-on-a-manifold’ framework for global registration of 3d range data. In *3D Imaging, Modeling, Processing, Visualization and Transmission (3DIMPVT), 2011 International Conference on*, pages 350–357. IEEE, 2011. 2
- [9] D. R. Canelhas, T. Stoyanov, and A. J. Lilienthal. Sdf tracker: A parallel algorithm for on-line pose estimation and scene reconstruction from depth images. In *2013 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 3671–3676. IEEE, 2013. 2
- [10] M. Danelljan, G. Meneghetti, F. S. Khan, and M. Felsberg. A probabilistic framework for color-based point set registration. In *CVPR*, volume 1, page 3, 2016. 2
- [11] B. Drost, M. Ulrich, N. Navab, and S. Ilic. Model globally, match locally: Efficient and robust 3d object recognition. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 998–1005. IEEE, 2010. 4, 6
- [12] B. Eckart, K. Kim, A. Troccoli, A. Kelly, and J. Kautz. Accelerated generative models for 3d point cloud data. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 5497–5505. IEEE, 2016. 2
- [13] G. D. Evangelidis and R. Horaud. Joint registration of multiple point sets. *CoRR*, abs/1609.01466, 2016. 2
- [14] S. Fantoni, U. Castellani, and A. Fusiello. Accurate and automatic alignment of range surfaces. In *2012 Second International Conference on 3D Imaging, Modeling, Processing, Visualization & Transmission*, pages 73–80. IEEE, 2012. 2, 3, 4
- [15] A. W. Fitzgibbon. Robust registration of 2d and 3d point sets. *Image and Vision Computing*, 21(13):1145–1153, 2003. 5
- [16] V. M. Govindu. Robustness in motion averaging. In *Computer Vision—ACCV 2006*, pages 457–466. Springer, 2006. 4
- [17] V. M. Govindu and A. Pooja. On averaging multiview relations for 3d scan registration. *Image Processing, IEEE Transactions on*, 23(3):1289–1302, 2014. 2
- [18] F. Guey and A. Geiger. Displets: Resolving stereo ambiguities using object knowledge. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4165–4175, 2015. 3
- [19] D. F. Huber and M. Hebert. Fully automatic registration of multiple 3d data sets. *Image and Vision Computing*, 21(7):637–650, 2003. 1, 2, 4
- [20] A. E. Johnson and M. Hebert. Using spin images for efficient object recognition in cluttered 3d scenes. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 21(5):433–449, 1999. 4
- [21] W. Kehl, N. Navab, and S. Ilic. Coloured signed distance fields for full 3d object reconstruction. In *British Machine Vision Conference*, 2014. 1, 2, 7
- [22] S. Krishnan, P. Y. Lee, J. B. Moore, S. Venkatasubramanian, et al. Global registration of multiple 3d point sets via optimization-on-a-manifold. In *Symposium on Geometry Processing*, pages 187–196, 2005. 2
- [23] B. Leibe, A. Leonardis, and B. Schiele. Combined object categorization and segmentation with an implicit shape model. In *In ECCV workshop on statistical learning in computer vision*, pages 17–32, 2004. 4
- [24] M. Levoy, K. Pulli, B. Curless, S. Rusinkiewicz, D. Koller, L. Pereira, M. Ginzton, S. Anderson, J. Davis, J. Ginsberg, J. Shade, and D. Fulk. The Digital Michelangelo Project: 3D scanning of large statues. In *Proceedings of ACM SIGGRAPH 2000*, pages 131–144, July 2000. 1
- [25] B. Lévy and Y. Liu. L p centroidal voronoi tessellation and its applications. In *ACM Transactions on Graphics (TOG)*, volume 29, page 119. ACM, 2010. 6
- [26] J. Liebelt and C. Schmid. Multi-view object class detection with a 3d geometric model. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 1688–1695. IEEE, 2010. 3
- [27] Y. Liu. Automatic 3d free form shape matching using the graduated assignment algorithm. *Pattern Recognition*, 38(10):1615–1631, 2005. 2
- [28] A. S. Mian, M. Bennamoun, and R. Owens. Three-dimensional model-based object recognition and segmentation in cluttered scenes. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 28(10):1584–1601, 2006. 2, 4, 6
- [29] G. Miller. Efficient algorithms for local and global accessibility shading. In *Proceedings of the 21st annual conference on Computer graphics and interactive techniques*, pages 319–326. ACM, 1994. 6
- [30] M. Muja and D. G. Lowe. Scalable nearest neighbor algorithms for high dimensional data. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 36(11):2227–2240, 2014. 4

- [31] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohi, J. Shotton, S. Hodges, and A. Fitzgibbon. Kinectfusion: Real-time dense surface mapping and tracking. In *Mixed and augmented reality (ISMAR), 2011 10th IEEE international symposium on*, pages 127–136. IEEE, 2011. [1](#), [2](#)
- [32] M. Nießner, M. Zollhöfer, S. Izadi, and M. Stamminger. Real-time 3d reconstruction at scale using voxel hashing. *ACM Transactions on Graphics (TOG)*, 2013. [2](#)
- [33] J. Novatnack and K. Nishino. Scale-dependent/invariant local 3d shape descriptors for fully automatic registration of multiple sets of range images. In *Computer Vision–ECCV 2008*, pages 440–453. Springer, 2008. [2](#)
- [34] C. Papazov, S. Haddadin, S. Parusel, K. Krieger, and D. Burschka. Rigid 3d geometry matching for grasping of known objects in cluttered scenes. *The International Journal of Robotics Research*, 2012. [4](#)
- [35] R. B. Rusu and S. Cousins. 3D is here: Point Cloud Library (PCL). In *IEEE International Conference on Robotics and Automation (ICRA)*, Shanghai, China, May 9-13 2011. [7](#)
- [36] R. F. Salas-Moreno, R. A. Newcombe, H. Strasdat, P. H. Kelly, and A. J. Davison. Slam++: Simultaneous localisation and mapping at the level of objects. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1352–1359, 2013. [3](#)
- [37] M. Slavcheva, W. Kehl, N. Navab, and S. Ilic. SDF-2-SDF: Highly Accurate 3D Object Reconstruction. In *European Conference on Computer Vision (ECCV)*, 2016. [2](#), [6](#), [7](#)
- [38] Y. Tang and J. Feng. Hierarchical Multiview Rigid Registration. *Computer Graphics Forum*, 2015. [2](#)
- [39] G. Taubin. Smooth signed distance surface reconstruction and applications. In *Iberoamerican Congress on Pattern Recognition*, pages 38–45. Springer, 2012. [6](#)
- [40] R. Toldo, A. Beinat, and F. Crosilla. Global registration of multiple point clouds embedding the generalized procrustes analysis into an icp framework. In *3DPVT 2010 Conference*, 2010. [2](#)
- [41] F. Tombari, S. Salti, and L. Di Stefano. Unique signatures of histograms for local surface description. In *Computer Vision–ECCV 2010*, pages 356–369. Springer, 2010. [4](#)
- [42] F. Tombari and L. D. Stefano. Hough voting for 3d object recognition under occlusion and clutter. *IPSJ Transactions on Computer Vision and Applications*, 4(0):20–29, 2012. [4](#)
- [43] A. Torsello, E. Rodola, and A. Albarelli. Multiview registration via graph diffusion of dual quaternions. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 2441–2448. IEEE, 2011. [2](#)
- [44] P. Yan, S. M. Khan, and M. Shah. 3d model based object class detection in an arbitrary view. In *2007 IEEE 11th International Conference on Computer Vision*, pages 1–6. IEEE, 2007. [3](#)
- [45] S. Yingze Bao, M. Chandraker, Y. Lin, and S. Savarese. Dense object reconstruction with semantic priors. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1264–1271, 2013. [2](#)
- [46] Q.-Y. Zhou, J. Park, and V. Koltun. Fast global registration. In *European Conference on Computer Vision*, pages 766–782. Springer International Publishing, 2016. [2](#)
- [47] J. Zhu, D. Meng, Z. Li, S. Du, and Z. Yuan. Robust registration of partially overlapping point sets via genetic algorithm with growth operator. *Image Processing, IET*, 8(10):582–590, 2014. [2](#)
- [48] J. Zhu, L. Zhu, Z. Jiang, X. Bai, Z. Li, and L. Wang. Local to global registration of multi-view range scans using spanning tree. *Computers & Electrical Engineering*, 2016. [2](#)
- [49] J. Zhu, L. Zhu, Z. Li, C. Li, and J. Cui. Automatic multi-view registration of unordered range scans without feature extraction. *Neurocomputing*, 171:1444–1453, 2016. [2](#)
- [50] M. Z. Zia, M. Stark, B. Schiele, and K. Schindler. Detailed 3d representations for object recognition and modeling. *IEEE transactions on pattern analysis and machine intelligence*, 35(11):2608–2623, 2013. [3](#)