

PSDF Fusion: Probabilistic Signed Distance Function for On-the-fly 3D Data Fusion and Scene Reconstruction

Wei Dong, Qiuyuan Wang, Xin Wang, and Hongbin Zha

Key Laboratory of Machine Perception (MOE), School of EECS, Peking University
Cooperative Medianet Innovation Center, Shanghai Jiao Tong University
{w.dong,qiuyuanwang,xinwang.cis}@pku.edu.cn, zha@cis.pku.edu.cn

Abstract. We propose a novel 3D spatial representation for data fusion and scene reconstruction. *Probabilistic Signed Distance Function* (Probabilistic SDF, PSDF) is proposed to depict uncertainties in the 3D space. It is modeled by a joint distribution describing SDF value and its inlier probability, reflecting input data quality and surface geometry. A *hybrid data structure* involving voxel, surfel, and mesh is designed to fully exploit the advantages of various prevalent 3D representations. Connected by PSDF, these components reasonably cooperate in a consistent framework. Given sequential depth measurements, PSDF can be incrementally refined with less ad hoc parametric Bayesian updating. Supported by PSDF and the efficient 3D data representation, high-quality surfaces can be extracted on-the-fly, and in return contribute to reliable data fusion using the geometry information. Experiments demonstrate that our system reconstructs scenes with higher model quality and lower redundancy, and runs faster than existing online mesh generation systems.

Keywords: Signed Distance Function, Bayesian Updating

1 Introduction

In recent years, we have witnessed the appearance of consumer-level depth sensors and the increasing demand of real-time 3D geometry information in next-generation applications. Therefore, online dense scene reconstruction has been a popular research topic. The essence of the problem is to fuse noisy depth data stream into a reliable 3D representation where clean models can be extracted. It is necessary to consider uncertainty in terms of sampling density, measurement accuracy, and surface complexity so as to better understand the 3D space.

Many representations built upon appropriate mathematical models are designed for robust data fusion in such a context. To handle uncertainties, *surfel* and *point* based approaches [29,13,15] adopt filtering-based probabilistic models that explicitly manipulate input data. *Volume* based methods [20,12,5,27], on the other hand, maximize spatial probabilistic distributions and output discretized 3D properties such as SDF and occupancy state. With fixed topologies, *mesh* based methods [35] may also involve parametric minimization of error functions.

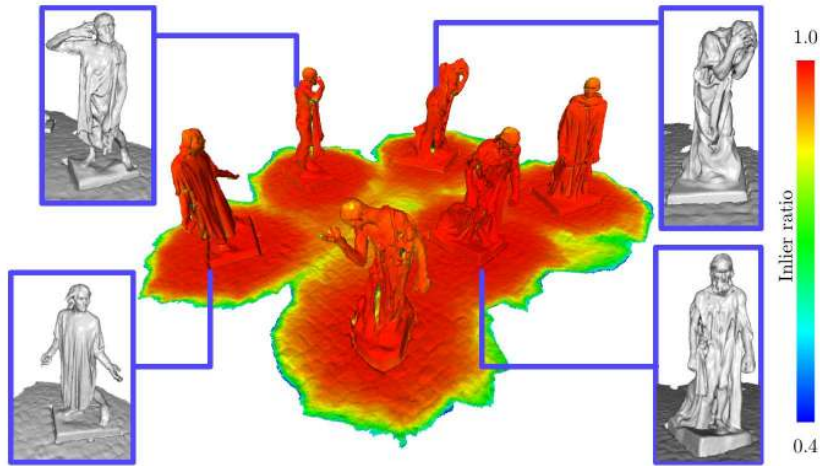


Fig. 1. Reconstructed mesh of *burghers*. The heatmap denotes the SDF inlier ratio (SDF confidence). Note details are preserved and outliers have been all removed without any post-processing. Inlier-ratio heatmaps in Fig.4,5,6 also conform to this colorbar.

While such representations have been proven effective by various applications, their underlying data structures endure more or less drawbacks. Surfels and points are often loosely managed without topology connections, requiring additional modules for efficient indexing and rendering [1], and is relatively prone to noisy input. Volumetric grids lack flexibility to some extent, hence corresponding data fusion can be either oversimplified using weighted average [20], or much time-consuming in order to maximize joint distributions [27]. In addition, ray-casting based volume rendering is also non-trivial. Usually storing vertices with strong topological constraints, mesh is similarly hard to manipulate and is less applicable to many situations. There have been studies incorporating aforementioned data structures [18,24,14,23], yet most of these pipelines are loosely organized without fully taking the advantages of each representation.

In this paper, we design a novel framework to fully exploit the power of existing 3D representations, supported by PSDF-based probabilistic computations. Our framework is able to perform reliable depth data fusion and reconstruct high-quality surfaces in real-time with more details and less noise, as depicted in Fig.1. Our contributions can be concluded as:

1. We present a novel hybrid data structure integrating voxel, surfel, and mesh;
2. The involved 3D representations are systematically incorporated in the consistent probabilistic framework linked by the proposed PSDF;
3. Incremental 3D data fusion is built upon less ad-hoc probabilistic computations in a parametric Bayesian updating fashion, contributes to online surface reconstruction, and benefits from iteratively recovered geometry in return.

2 Related Work

Dense reconstruction from depth images. There is a plethora of successful off-the-shelf systems that reconstruct 3D scenes from depth scanners’ data. [4] presents the volumetric 3D grids to fuse multi-view range images. [20] extends [4] and proposes KinectFusion, a real-time tracking and dense mapping system fed by depth stream from a consumer-level Kinect. KinectFusion has been improved by VoxelHashing [22], InfiniTAM [12], and BundleFusion [5] that leverage more memory-efficient volumetric representations. While these systems perform online depth fusion, they usually require offline MarchingCubes [17] to output final mesh models; [14,24,6] incorporates online meshing modules in such systems. Instead of utilizing volumetric spatial representations, [13] proposes point-based fusion that maintains light-weight dense point cloud or surfels as 3D maps. ElasticFusion [31] and InfiniTAM_v3 [23] are efficient implementations of [13] with several extensions; [16] further improves [13] by introducing surface curvatures. The point-based methods are unable to output mesh online, hence may not be suitable for physics-based applications. [34,2] split scenes into fragments and register partially reconstructed mesh to build comprehensive models, but their offline property limits their usages.

3D scene representations. Dense 3D map requires efficient data structures to support high resolution reconstructions. For volumetric systems, plain 3D arrays [20] are unsuitable for large scale scenes due to spatial redundancies. In view of this, moving volumes method [30] is introduced to maintain spatial properties only in active areas. Octree is used to ensure a complete yet adaptive coverage of model points [10,33,24]. As the tree might be unbalanced causing long traversing time, hierarchical spatial hashing is utilized [22,12] supporting $O(1)$ 3D indexing, and is further extended to be adaptive to local surface complexities [11].

There are also studies that directly represent scenes as point clouds or mesh during reconstruction. In [13,31] point clouds or surfels are simply arranged in an 1D array. Considering topologies in mesh, [18] manages point clouds with inefficient KD-Tress for spatial resampling. [35] maintains a 2.5D map with fixed structured triangles which will fail to capture occlusions. Hybrid data structures are also used to combine volumes and mesh. [24] builds an octree-based structure where boundary conditions have to be carefully considered in term of mesh triangles. [14] uses spatial hashed blocks and stores mesh triangles in the block level, but ignores vertex sharing between triangles. [6] reveals the correspondences between mesh vertices and voxel edges, reducing the redundancy in the aspect of data structure. Yet improvement is required to remove false surfaces generated from noisy input data.

Uncertainty-aware data fusion. Uncertainty is one of the core problems remain in 3D reconstruction, which may come from imperfect inputs or complex environments. In volumetric representations that split the space into grids, probability distributions are usually utilized to model spatial properties. Binary occupancy is an intuitive variable configuration, denoting whether a voxel is physically occupied. [32] proposes a joint distribution of point occupancy state and inlier ratio over the entire volume with visual constraints and achieves compet-

itive results. [28,27] similarly emphasize ray-consistency to reconstruct global-consistent surfaces, whose inferences are too sophisticated to run in real-time. Although surface extraction can be performed on occupancy grids via thresholding or ray-casting, it is usually very sensitive to parameters. Instead of maintaining a $\{0, 1\}$ field, [4] introduces SDF which holds signed projective distances from voxels to their closest surfaces sampled by input depth measurements. [20,19] use weight average of Truncated SDF in the data fusion stage by considering a per-voxel Gaussian distribution regarding SDF as a random variable. While Gaussian noise can be smoothed by weight average, outliers have to be carefully filtered out with ad hoc operations. [22] uses a temporal recycling strategy by periodically subtracting weight in volumes; [14] directly carves out noisy inputs; [5] proposes a weighted subtraction to de-integrate data which are assumed to be incorrectly registered. As a non-local prior, [7] refines SDF value on-the-go using plane fitting, which performs well mainly in flat scenes with relatively low voxel resolutions. We find a lack of systematic probabilistic solution for SDF dealing both Gaussian noise and possible outliers.

For point-based representations, [29] proposes an elegant math model treating inlier ratio and depth of a point as random variables subject to a special distribution. The parameters of such distributions are updated in a Bayesian fashion. This approach is adopted by [8] in SLAM systems for inverse depths, achieving competitive results. An alternative is to select ad hoc weights involving geometry and photometric properties of estimated points and computing weighted average [26,15]. This simple strategy shares some similarity to the fusion of SDF values, and is also used in RGB-D systems where depths are more reliable [13,23]. Despite the solid math formulations, point-based methods are comparatively prone to noise due to their discrete representations.

3 Overview

Our framework is based on the *hybrid data structure* involving three *3D representations* linked by *PSDF*. The pipeline consists of iterative operations of data fusion and surface generation.

3.1 Hybrid Data Structure

We follow [22,14,6] and use a spatial hashing based structure to efficiently manage the space. A hash entry would point to a block, which is the smallest unit to allocate and free. A block is further divided into $8 \times 8 \times 8$ small voxels. Following [6] we consider a voxel as a 3-edge structure instead of merely a cube, as depicted in Fig.2(a), which will avoid ambiguity when we refer to shared edges. PSDF values are stored at the corners of these structures. In addition, we maintain surfels on the volumetric grids by limiting their degree of freedom on the edges of voxels; within a voxel at most 3 surfels on edges could be allocated. This constraint would regularize the distribution of surfels, guarantee easier access, and avoid duplicate allocation. Triangles are loosely organized in the level of blocks,

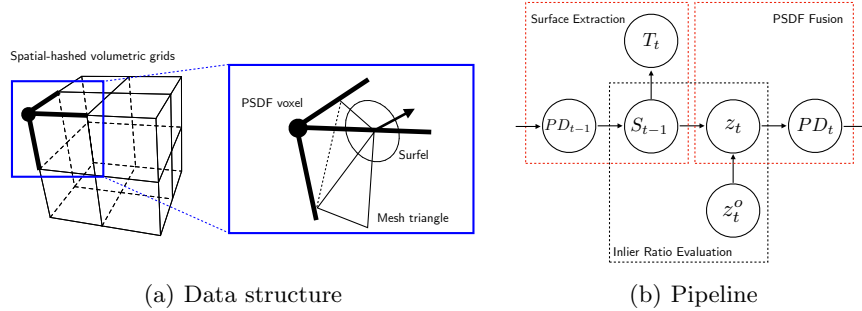


Fig. 2. A basic hybrid unit, and the system pipeline.

linking adjacent surfels. In the context of mesh, a surfel could be also interpreted as a triangle vertex.

3.2 3D Representations

Voxel and PSDF. In most volumetric reconstruction systems SDF or truncated SDF (TSDF) (denoted by D) of a voxel is updated when observed 3D points fall in its neighbor region. Projective signed distances from measurements, which could be explained as SDF observations, are integrated by computing weight average. Newcombe [19] suggests that it can be regarded as the solution of a maximum likelihood estimate of a joint Gaussian distribution taking SDF as a random variable. While Gaussian distribution could depict the uncertainty of data noise, it might fail to handle outlier inputs which are common in reconstruction tasks using consumer-level sensors. Moreover, SDF should depict the projective distance from a voxel to its *closest* surface. During integration, however, it is likely that *non-closest* surface points are taken into account, which should also be regarded as outlier SDF observations. In view of this, we introduce another random variable π to denote the inlier ratio of SDF, initially used in [29] to model the inlier ratio of 3D points:

$$p(D_i^o | D, \tau_i, \pi) = \pi \mathcal{N}(D_i^o; D, \tau_i^2) + (1 - \pi) \mathcal{U}(D_i^o; D_{min}, D_{max}), \quad (1)$$

where D_i^o reads an SDF observation computed with depth measurements, τ_i is the variance of the SDF observation, \mathcal{N} and \mathcal{U} are Gaussian and Uniform distributions.

Following [29], the posterior of PSDF can be parameterized by a Beta distribution multiplying a Gaussian distribution $\mathcal{B}(a, b) \mathcal{N}(\mu; \sigma^2)$, given a series of observed input SDF measurements. The details will be discussed in §4.1. The parameters a, b, μ, σ of the parameterized distribution are maintained per voxel. **Surfel.** A surfel in our pipeline is formally defined by a position \mathbf{x} , a normal \mathbf{n} , and a radius r . Since a certain surfel is constrained on an edge in the volume, \mathbf{x} is generally an interpolation of 2 adjacent voxel corners.

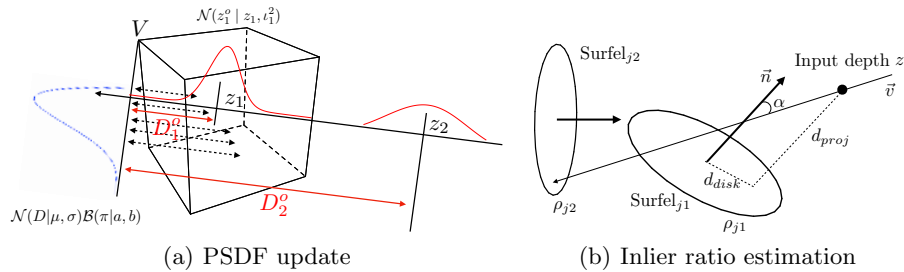


Fig. 3. Left, illustration of PSDF distribution with multiple SDF observations (D_1^o is likely to be an inlier observation, D_2^o is possibly an outlier). The red curves show Gaussian distributions per observation. The blue curve depicts a Beta distribution which intuitively suggests that inlier observations should be around D_1^o . Right, estimation of SDF inlier ratio ρ involving observed input 3D point and already known surfels.

Triangle. A triangle consists of 3 edges, each linking two adjacent surfels. These surfels can be located in different voxels, even different blocks. In our framework triangles are mainly extracted for rendering; the contained topology information may be further utilized in extensions.

Depth input. We receive depth measurements from sensors as input, while sensor poses are assumed known. Each observed input depth z^o is modeled as a random variable subject to a simple Gaussian distribution:

$$p(z^o | z, \iota) = \mathcal{N}(z^o; z, \iota^2), \quad (2)$$

where ι can be estimated from a precomputed sensor error model.

3.3 Pipeline

In general, our system will first generate a set of surfels S_{t-1} in the volumetric PSDF field PD_{t-1} . Meanwhile, mesh triangle set T_t is also determined by linking reliable surfels in S_{t-1} . S_{t-1} explicitly defines the surfaces of the scene, hence can be treated as a trustworthy geometry cue to estimate outlier ratio of the input depth data z_t^o . PD_{t-1} is then updated to PD_t by fusing evaluated depth data distribution z_t via Bayesian updating. The process will be performed iteratively every time input data come, as depicted in Fig.2(b). We assume the poses of the sensors are known and all the computations are in the world coordinate system.

4 PSDF Fusion and Surface Reconstruction

4.1 PSDF Fusion

Similar to [20], in order to get SDF observations of a voxel V given input 3D points from depth images, we first project V to the depth image to find the

projective closest depth measurement z_i . Signed distance from V to the input 3D data is defined by

$$D_i = z_i - z^V, \quad (3)$$

where z^V is a constant value, the projective depth of V along the scanning ray.

The observed D_i^o is affected by the variance ι_i of z_i in Eq.2 contributing to the Gaussian distribution component in Eq.1, provided D_i is an inlier. Otherwise, D_i would be counted in the uniform distribution part. Fig.3(a) illustrates the possible observations of SDF in one voxel.

Variance ι_i can be directly estimated by pre-measured sensor priors such as proposed in [21]. In this case, due to the simple linear form in Eq.3, we can directly set $D_i^o = z_i^o - z^V$ and $\tau_i = \iota_i$ in Eq.1.

Given a series of independent observations D_i^o , we can derive the posterior

$$p(D, \pi \mid D_1^o, \tau_1 \cdots D_n^o, \tau_n) \propto p(D, \pi) \prod_i p(D_i^o \mid D, \tau_i^2, \pi), \quad (4)$$

where $p(D, \pi)$ is a prior and $p(D_i^o \mid D, \tau_i^2, \pi)$ is defined by Eq.1. It would be intractable to evaluate the production of such distributions with additions. Fortunately, [29] proved that the posterior could be approximated by a parametric joint distribution:

$$p(D, \pi \mid D_1^o, \tau_1 \cdots D_n^o, \tau_n) \propto \mathcal{B}(\pi \mid a_n, b_n) \mathcal{N}(D \mid \mu_n, \sigma_n^2), \quad (5)$$

therefore the problem could be simplified as a parameter estimation in an incremental fashion:

$$\mathcal{B}(\pi \mid a_n, b_n) \mathcal{N}(D \mid \mu_n, \sigma_n^2) \propto p(D_n^o \mid D, \tau_n^2, \pi) \mathcal{B}(\pi \mid a_{n-1}, b_{n-1}) \mathcal{N}(D \mid \mu_{n-1}, \sigma_{n-1}^2). \quad (6)$$

In [29] by equating first and second moments of the random variables π and D , the parameters could be easily updated, in our case evoking the change of SDF distribution and its inlier probability:

$$\mu_n = \frac{C_1}{C_1 + C_2} m + \frac{C_2}{C_1 + C_2} \mu_{n-1}, \quad (7)$$

$$\mu_n^2 + \sigma_n^2 = \frac{C_1}{C_1 + C_2} (s^2 + m^2) + \frac{C_2}{C_1 + C_2} (\sigma_{n-1}^2 + \mu_{n-1}^2), \quad (8)$$

$$s^2 = 1 / \left(\frac{1}{\sigma_{n-1}^2} + \frac{1}{\tau_n^2} \right), \quad (9)$$

$$m = s^2 \left(\frac{\mu_{n-1}}{\sigma_{n-1}^2} + \frac{D_i^o}{\tau_n^2} \right), \quad (10)$$

$$C_1 = \frac{a_{n-1}}{a_{n-1} + b_{n-1}} \mathcal{N}(D_i^o; \mu_{n-1}, \sigma_{n-1}^2 + \tau_n^2), \quad (11)$$

$$C_2 = \frac{b_{n-1}}{a_{n-1} + b_{n-1}} \mathcal{U}(D_i^o; D_{min}, D_{max}), \quad (12)$$

the computation of a and b are the same as [29] hence ignored here. In our experiments we find that a truncated D_i^o leads to better results, as it directly rejects distant outliers. SDF observations from non-closest surfaces are left to be handled by PSDF.

4.2 Inlier Ratio Evaluation

In Eq.11-12, the expectation of $\mathcal{B}(\pi | a, b)$ is used to update the coefficients, failing to make full use of known geometry properties in scenes. In our pipeline, available surface geometry is considered to evaluate the inlier ratio ρ_n of D_n^o , replacing the simple $\frac{a_{n-1}}{a_{n-1}+b_{n-1}}$. Note ρ_n is computed per-frame in order to update C_1, C_2 ; π is still parameterized by a and b .

ρ_n can be determined by whether an input point z is near the closest surface of a voxel and results in an inlier SDF observation. We first cast the scanning ray into the volume and collect the surfels maintained on the voxels hit by the ray. Given the surfels, 3 heuristics are used, as illustrated in Fig.3(b).

Projective distance. This factor is used to measure whether a sampled point is close enough to a surfel which is assumed the nearest surface to the voxel:

$$w_{dist} = \exp\left(\frac{-|\mathbf{n}^T(\mathbf{x} - z\mathbf{v})|^2}{2\theta^2}\right), \quad (13)$$

where \mathbf{v} is the normalized direction of the ray in world coordinate system and θ is a preset parameter proportional to the voxel resolution.

Angle. Apart from projective distance, we consider angle as another factor, delineating the possibility that a scanning ray will hit a surfel. We use the empirical angle weight in [15]:

$$w_{angle} = \begin{cases} \frac{\cos(\alpha) - \cos(\alpha_{max})}{1 - \cos(\alpha_{max})}, & \text{if } \alpha < \alpha_{max}, \\ w_{angle}^0, & \text{else,} \end{cases} \quad (14)$$

where $\alpha = \langle \mathbf{n}, \mathbf{v} \rangle$, α_{max} is set to 80 deg and w_{angle}^0 assigned to 0.1.

Radius. The area that surfels could influence vary, due to the local shape of the surface. The further a point is away from the center of a surfel, the less possible it would be supported. A sigmoid-like function is used to encourage a smooth transition of the weight:

$$w_{radius} = \gamma + \frac{2(1 - \gamma)}{1 + \exp\left(\frac{-d_{disk}}{r}\right)}, \quad (15)$$

$$d_{disk} = \sqrt{(z\mathbf{v} - \mathbf{x})^T(I - \mathbf{n}\mathbf{n}^T)(z\mathbf{v} - \mathbf{x})}, \quad (16)$$

where parameter $\gamma \in [0, 1)$ and is set to 0.5 in our case.

Putting all the factors together, we now have

$$\rho = w_{dist} \cdot w_{radius} \cdot w_{angle}. \quad (17)$$

To compute the ρ predicted by all the surfels, one may consider either summations or multiplications. However, we choose the highest ρ instead – intuitively a depth measurement is a sample on a surface, corresponding to exactly one surfel. A more sophisticated selection might include a ray consistency evaluation [32,27,28] where occlusion is handled. When a new area is explored where no surfels have been extracted, we use a constant value $\rho_{pr} = 0.1$ to represent a simple occupancy prior in space, hence we have

$$\rho = \max_j \{\rho_{pr}, \rho_1, \dots, \rho_j, \dots\}. \quad (18)$$

4.3 Surface Extraction

PSDF implicitly defines zero crossing surfaces and decides whether they are true surfaces. The surface extraction is divided into two steps.

Surfel generation. In this stage we enumerate zero-crossing points upon 3 edges of each voxel and generate surfels when condition

$$\begin{aligned} \mu^{i1} \cdot \mu^{i2} < 0, \\ \frac{a^{i1}}{a^{i1} + b^{i1}} > \pi_{thr} \text{ and } \frac{a^{i2}}{a^{i2} + b^{i2}} > \pi_{thr}, \end{aligned} \quad (19)$$

are satisfied, where $i1$ and $i2$ are indices of adjacent voxels and π_{thr} is a confidence threshold. Supported by the reliable update of the PSDF, false surfaces could be rejected and duplicates could be removed. According to our experiments, our framework is not sensitive to π_{thr} ; 0.4 would work for all the testing scenes. A surfel’s position \mathbf{x} would be the linear interpolation of corresponding voxels’ positions \mathbf{x}^i indexed by i , and the radius would be determined by σ of adjacent voxels, simulating its affecting area. Normal is set to normalized gradient of the SDF field, as mentioned in [20].

$$\mathbf{x} = \frac{|\mu^{i2}|}{|\mu^{i1}| + |\mu^{i2}|} \mathbf{x}^{i1} + \frac{|\mu^{i1}|}{|\mu^{i1}| + |\mu^{i2}|} \mathbf{x}^{i2}, \quad (20)$$

$$r = \frac{|\mu^{i2}|}{|\mu^{i1}| + |\mu^{i2}|} \sigma^{i1} + \frac{|\mu^{i1}|}{|\mu^{i1}| + |\mu^{i2}|} \sigma^{i2}, \quad (21)$$

$$\mathbf{n} = \nabla \mu / \|\nabla \mu\|. \quad (22)$$

Triangle generation. Having sufficient geometry information within surfels, there is only one more step to go for rendering-ready mesh. The connections between adjacent surfels are determined by the classical MarchingCubes [17] method. As a simple modification, we reject edges in the voxel whose σ is larger than a preset parameter σ_{thr} . This operation will improve the visual quality of reconstructed model while preserving surfels for the prediction stage.

5 Experiments

We test our framework (denoted by *PSDF*) on three RGB-D datasets: TUM [25], ICL-NUIM [9], and dataset from Zhou and Koltun [34]. Our method is

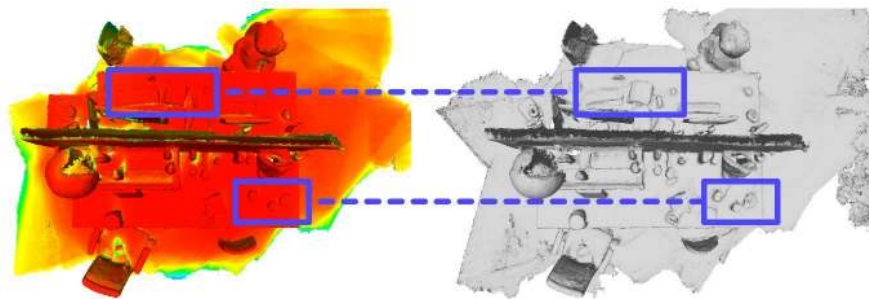


Fig. 4. Comparison of output mesh of *frei3.long_office* scene. Left, *PSDF*. Right, *TSDF*. Our method generates smooth surfaces and clean edges of objects, especially in blue boxes.

compared against [6] (denoted by *TSDF*) which incrementally extracts mesh in spatial-hashed TSDF volumes. The sensors’ poses are assumed known for these datasets, therefore the results of *TSDF* should be similar to other state-of-the-art methods such as [22,23] where TSDF integration strategies are the same. We demonstrate that our method reconstructs high quality surfaces by both qualitative and quantitative results. Details are preserved while noise is removed in the output models. The running speed for online mesh extraction is also improved by avoiding computations on false surfel candidates.

For [25,34] we choose a voxel size of 8mm and $\sigma_{thr} = 16\text{mm}$; for [9] voxel size is set to 12mm and $\sigma_{thr} = 48\text{mm}$. The truncation distance is set to $3 \times$ voxel size plus $3 \times \tau$; with a smaller truncation distance we found strides and holes in meshes. Kinect’s error model [21] was adopted to get ι where the factor of angle was removed, which we think might cause double counting considering w_{angle} in the inlier prediction stage. The program is written in C++/CUDA 8.0 and runs on a laptop with an Intel i7-6700 CPU and an NVIDIA 1070 graphics card.

5.1 Qualitative Results

We first show that *PSDF* accompanied by the related mesh extraction algorithm produces higher quality surfaces than *TSDF*. Our results are displayed with shaded heatmap whose color indicates the inlier ratio of related SDF. Both geometry and probability properties can be viewed in such a representation.

Fig.4 shows that *PSDF* outperforms *TSDF* by generating clean boundaries of small objects and rejecting noisy areas on the ground. In Fig.5, in addition to the results of *TSDF*, we also display the reconstructed mesh from offline methods provided by [34] as references. It appears that our method produces results very similar to [34]. While guaranteeing well-covered reconstruction of scenes, we filter outliers and preserve details. In *copyroom*, the wires are completely reconstructed, one of which above PC is smoothed out in [34] and only partially recovered by *TSDF*. In *lounge*, we can observe a complete shape of table, and de-noised details in the clothes.

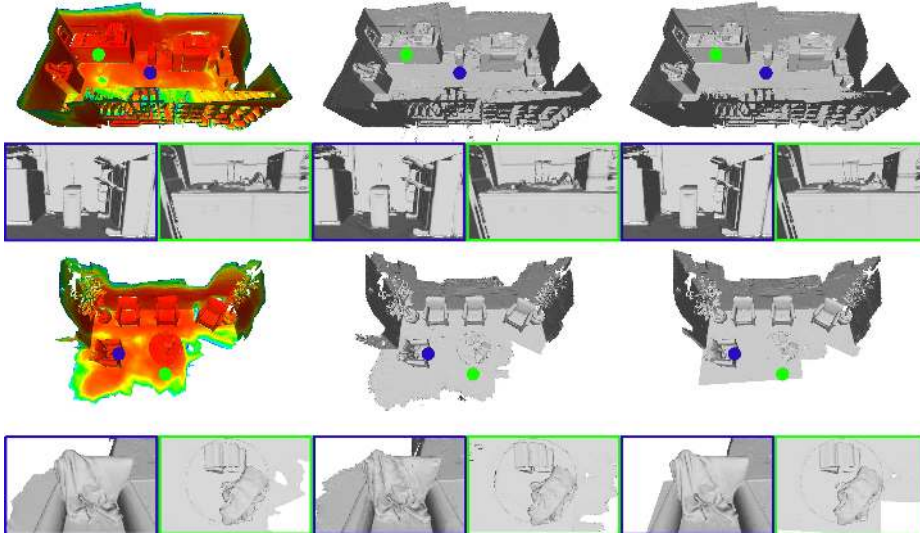


Fig. 5. Output mesh of the *copyroom* and *lounge* scenes. From left to right, *PSDF*, *TSDF*, mesh provided by [34]. Zoomed in regions show that our method is able to filter outliers while maintaining complete models and preserving details. Best viewed in color.

We also visualize the incremental update of π by rendering $\mathbb{E}(\beta(\pi|a, b)) = a/(a + b)$ as the inlier ratio of reconstructed mesh in sequence using colored heatmap. Fig.6 shows the fluctuation of confidence around surfaces. The complex regions such as fingers and wrinkles on statues are more prone to noise, therefore apparent change of shape along with color can be observed.

5.2 Quantitative Results

Reconstruction Accuracy. We reconstruct mesh of the synthetic dataset *livingroom2* with added noise whose error model is presented in [9]. Gaussian noise on inverse depth plus local offsets is too complicated for our error model, therefore we simplify it by assigning inverse sigma at certain inverse depths to ι_i . The mesh vertices are compared with the ground truth point cloud using the free software *CloudCompare* [3].

Table.1 indicates that *PSDF* reconstructs better models than *TSDF*. Further details in Fig.7 suggest that less outliers appear in the model reconstructed by *PSDF*, leading to cleaner surfaces.

Mesh size. Our method maintains the simplicity of mesh by reducing false surface candidates caused by noise and outliers. As shown in Fig.8 and Table 2, *PSDF* in most cases generates less vertices (% 20) than *TSDF*, most of which are outliers and boundaries with low confidence. Fig.8 shows that the vertex count remains approximately constant when a loop closure occurred in *frei3.long-office*,

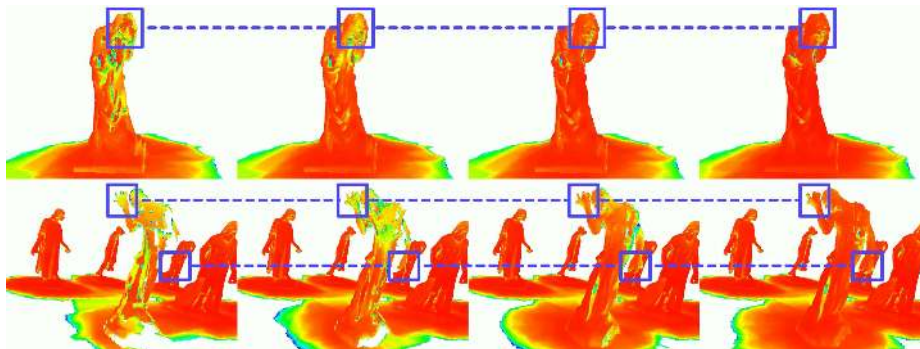


Fig. 6. Incremental reconstruction of the *burghers* dataset. Fluctuation of probability could be observed at error-prone regions as an indication of uncertainty propagation.

METHOD	MEAN (m)	STD (m)
<i>PSDF</i>	0.011692	0.015702
<i>TSDF</i>	0.022556	0.076120

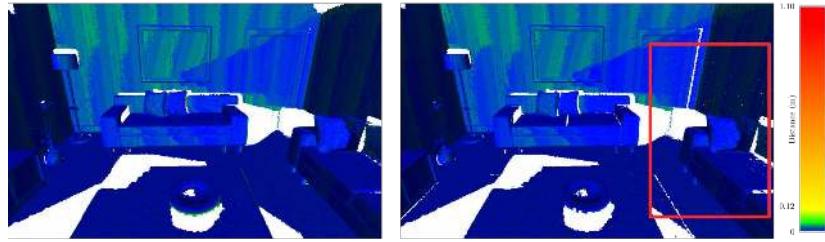
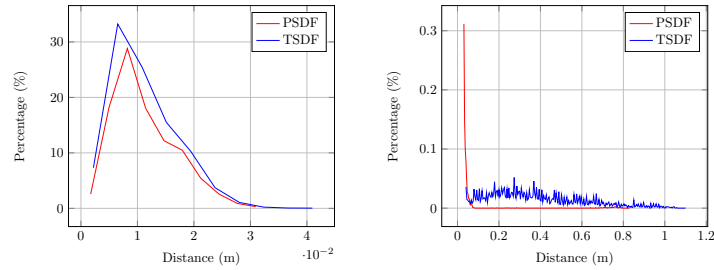
Table 1. Statistics of the point-to-point distances from the reconstructed model vertices to the ground truth point cloud. *PSDF* yields better reconstruction accuracy.

while the increasing rate is strictly constrained in the *lounge* sequence where there is no loop closure.

Time. To make a running time analysis, we take real world *lounge* as a typical sequence where noise is common while the camera trajectory fits scanning behavior of humans. As we have discussed, evaluation of inlier ratio was performed, increasing total time of the fusion stage. However we find on a GPU, even with a relatively high resolution, the average increased time is at the scale of ms (see Fig.8(c) and Table.2) and can be accepted.

When we come to meshing, we find that by taking the advantage of PSDF fusion and inlier ratio evaluation, unnecessary computations can be avoided and *PSDF* method runs faster than *TSDF*, as plotted in Fig.8(d). The meshing stage is the runtime bottleneck of the approach, in general the saved time compensate for the cost in fusion stage, see Fig.8(e) and Table.2.

We also compare the time of meshing to the widely used ray-casting that renders surfaces in real-time. According to Table 2, in some scenes where sensor is close to the surfaces performing scanning, less blocks are allocated in viewing frustum and the meshing speed could be comparative to ray-casting, as illustrated in Fig.8(f). As for other scenes requiring a large scanning range, especially *frei3-long_office* where more blocks in frustum have to be processed, ray-casting shows its advantage. We argue that in applications that only require visualization, ray-casting can be adopted; otherwise meshing offers more information and is still preferable.

(a) Point-to-point distance heatmap. Left, *PSDF*. Right, *TSDF*.

(b) Histogram head: 0 to 4cm

(c) Histogram tail: 4cm to ∞

Fig. 7. Comparison of output mesh quality of *PSDF* and *TSDF*. First row: heatmap of the point-to-point distance from the reconstructed model to the ground truth. Notice the outliers in the red box. Second row, the distance histogram divided into two parts to emphasize the existence of outliers. As shown in the histogram tail, *TSDF* generates many outliers, while *PSDF* avoids such problems. Best viewed in color.

6 Conclusions

We propose PSDF, a joint probabilistic distribution to model the 3D geometries and spatial uncertainties. With the help of Bayesian updating, parameters of the distribution could be incrementally estimated. Built upon a hybrid data structure, our framework can iteratively generate surfaces from the volumetric PSDF field and update PSDF values through reliable probabilistic data fusion supported by reconstructed surfaces. As an output, high-quality mesh can be generated in real-time with duplicates removed and noise cleared.

In the future, we seek to improve our framework by employing more priors to enrich the PSDF distribution. Localization modules will also be integrated in the probabilistic framework for a complete SLAM system.

Acknowledgments

This work is supported by the National Natural Science Foundation of China (61632003, 61771026), and National Key Research and Development Program of China (2017YFB1002601).

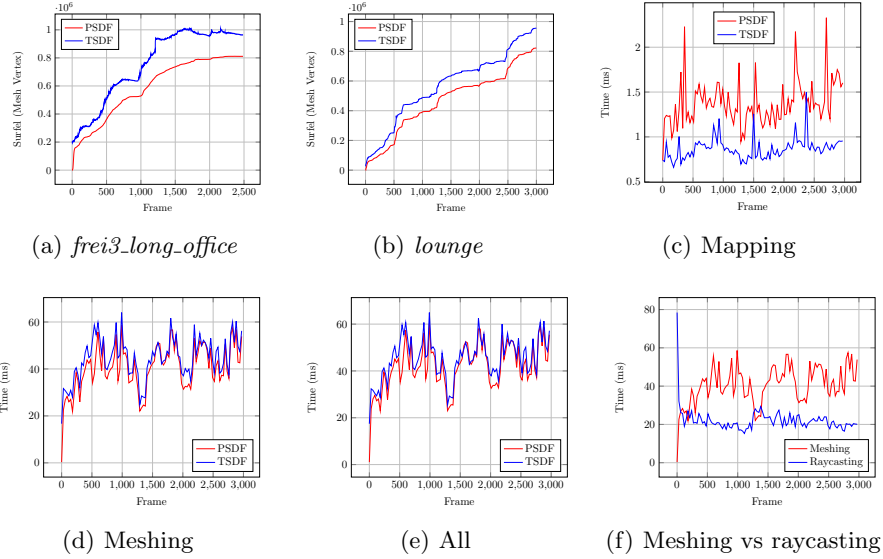


Fig. 8. (a)-(b), increasing trend of mesh vertices (surfels). (a), *frei3_long_office* scene in which a loop closure occurs at around 1500 frames. (b), *lounge* scene without loop closures. Mesh generated by *PSDF* consumes approximately 80% the memory of *TSDF*. (c)-(f), running time analysis of our approaches and compared methods on *lounge* scene. *PSDF* is 1ms slower on GPU due to additional inlier prediction stage, but will save more time for surface extraction, causing a faster speed in general. The speed of online meshing is slower than raycasting but comparable.

Dataset	Frames	MEMORY (vertex count)		TIME (ms)				Ray-casting
		PSDF	TSDF	PSDF		TSDF		
				Mapping	Meshing	Mapping	Meshing	
<i>burghers</i>	11230	1362303	1438818	1.053	39.691	0.771	38.907	27.761
<i>copyroom</i>	5490	1222196	1328397	1.329	39.090	0.909	41.595	21.292
<i>garden</i>	6152	922534	978206	1.473	68.224	0.907	66.680	25.479
<i>lounge</i>	3000	821360	955218	1.331	41.270	0.881	45.248	22.097
<i>livingroom1</i>	965	529305	518885	1.255	33.090	0.743	32.865	27.248
<i>livingroom2</i>	880	609421	683008	1.407	33.446	0.759	40.230	29.069
<i>office1</i>	965	667614	674034	1.264	24.933	0.799	27.654	26.957
<i>office2</i>	880	712685	883138	1.322	35.029	0.767	45.923	29.981
<i>frei1_xyz</i>	790	212840	352444	1.193	45.163	1.149	71.670	19.796
<i>frei3_long_office</i>	2486	811092	963875	2.424	159.417	1.375	161.485	26.545

Table 2. Evaluation of memory and time cost on various datasets. *PSDF* reduces model’s redundancy by rejecting false surfaces and noise. The mapping stage of *TSDF* is faster, but in general *PSDF* spends less time considering both mapping and meshing stages.

References

1. Botsch, M., Kobbelt, L.: High-quality Point-based Rendering on Modern GPUs. In: Proceedings of Pacific Conference on Computer Graphics and Applications. pp. 335–343 (2003) [2](#)
2. Choi, S., Zhou, Q.Y., Koltun, V.: Robust Reconstruction of Indoor Scenes. In: Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (2015) [3](#)
3. CloudCompare-project: CloudCompare. <http://www.cloudcompare.org/> [11](#)
4. Curless, B., Levoy, M.: A Volumetric Method for Building Complex Models from Range Images. In: Proceedings of ACM SIGGRAPH. pp. 303–312 (1996) [3](#), [4](#)
5. Dai, A., Nießner, M., Zollhöfer, M., Izadi, S., Theobalt, C.: BundleFusion: Real-time Globally Consistent 3D Reconstruction using On-the-fly Surface Reintegration. *ACM Transactions on Graphics* **36**(3), 24 (2017) [1](#), [3](#), [4](#)
6. Dong, W., Shi, J., Tang, W., Wang, X., Zha, H.: An Efficient Volumetric Mesh Representation for Real-time Scene Reconstruction using Spatial Hashing. In: Proceedings of IEEE International Conference on Robotics and Automation (2018) [3](#), [4](#), [10](#)
7. Dzitsiuk, M., Sturm, J., Maier, R., Ma, L., Cremers, D.: De-noising, Stabilizing and Completing 3D Reconstructions On-the-go using Plane Priors. In: Proceedings of IEEE International Conference on Robotics and Automation. pp. 3976–3983 (2017) [4](#)
8. Forster, C., Pizzoli, M., Scaramuzza, D.: SVO: Fast Semi-direct Monocular Visual Odometry. In: Proceedings of IEEE International Conference on Robotics and Automation. pp. 15–22 (2014) [4](#)
9. Handa, A., Whelan, T., McDonald, J., Davison, A.J.: A Benchmark for RGB-D Visual Odometry, 3D Reconstruction and SLAM. In: Proceedings of IEEE International Conference on Robotics and Automation. pp. 1524–1531 (2014) [9](#), [10](#), [11](#)
10. Hornung, A., Wurm, K.M., Bennewitz, M., Stachniss, C., Burgard, W.: OctoMap: an Efficient Probabilistic 3D Mapping Framework based on Octrees. *Autonomous Robots* **34**(3), 189–206 (2013) [3](#)
11. Kähler, O., Prisacariu, V., Valentin, J., Murray, D.: Hierarchical Voxel Block Hashing for Efficient Integration of Depth Images. *IEEE Robotics and Automation Letters* **1**(1), 192–197 (2016) [3](#)
12. Kähler, O., Prisacariu, V.A., Ren, C.Y., Sun, X., Torr, P., Murray, D.: Very High Frame Rate Volumetric Integration of Depth Images on Mobile Devices. *IEEE Transactions on Visualization & Computer Graphics* **21**(11), 1241–1250 (2015) [1](#), [3](#)
13. Keller, M., Lefloch, D., Lambers, M., Weyrich, T., Kolb, A.: Real-time 3D Reconstruction in Dynamic Scenes using Point-based Fusion. In: Proceedings of International Conference on 3DTV. pp. 1–8 (2013) [1](#), [3](#), [4](#)
14. Klingensmith, M., Dryanovski, I., Srinivasa, S.S., Xiao, J.: CHISEL: Real Time Large Scale 3D Reconstruction Onboard a Mobile Device using Spatially-Hashed Signed Distance Fields. In: Proceedings of Robotics: Science and Systems. pp. 1–8 (2015) [2](#), [3](#), [4](#)
15. Kolev, K., Tanskanen, P., Speciale, P., Pollefeys, M.: Turning Mobile Phones into 3D Scanners. In: Proceedings of IEEE Conference on Computer Vision and Pattern Recognition. pp. 3946–3953 (2014) [1](#), [4](#), [8](#)

16. Lefloch, D., Kluge, M., Sarbolandi, H., Weyrich, T., Kolb, A.: Comprehensive Use of Curvature for Robust and Accurate Online Surface Reconstruction. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **39**(12), 2349–2365 (Dec 2017) [3](#)
17. Lorensen, W.E., Cline, H.E.: Marching Cubes : A High Resolution 3D Surface Construction Algorithm. In: *Proceedings of ACM SIGGRAPH*. vol. 6, pp. 7–9 (1987) [3](#), [9](#)
18. Marton, Z.C., Rusu, R.B., Beetz, M.: On Fast Surface Reconstruction Methods for Large and Noisy Point Clouds. In: *Proceedings of IEEE International Conference on Robotics and Automation*. pp. 3218–3223 (2009) [2](#), [3](#)
19. Newcombe, R.: Dense visual SLAM. Ph.D. thesis, Imperial College London, UK (2012) [4](#), [5](#)
20. Newcombe, R.A., Molyneaux, D., Kim, D., Davison, A.J., Shotton, J., Hodges, S., Fitzgibbon, A., Izadi, S., Hilliges, O., Molyneaux, D., Kim, D., Davison, A.J., Kohli, P., Shotton, J., Hodges, S., Fitzgibbon, A.: KinectFusion: Real-Time Dense Surface Mapping and Tracking. In: *Proceedings of IEEE and ACM International Symposium on Mixed and Augmented Reality*. pp. 127–136 (2011) [1](#), [2](#), [3](#), [4](#), [6](#), [9](#)
21. Nguyen, C.V., Izadi, S., Lovell, D.: Modeling Kinect Sensor Noise for Improved 3D Reconstruction and Tracking. In: *Proceedings of IEEE International Conference on 3D Imaging, Modeling, Processing, Visualization & Transmission*. pp. 524–530 (2012) [7](#), [10](#)
22. Nießner, M., Zollhöfer, M., Izadi, S., Stamminger, M.: Real-time 3D Reconstruction at Scale Using Voxel Hashing. *ACM Transactions on Graphics* **32**(6), 169 (2013) [3](#), [4](#), [10](#)
23. Prisacariu, V.A., Kähler, O., Golodetz, S., Sapienza, M., Cavallari, T., Torr, P.H., Murray, D.W.: InfiniTAM v3: A Framework for Large-Scale 3D Reconstruction with Loop Closure. *ArXiv e-prints* (Aug 2017) [2](#), [3](#), [4](#), [10](#)
24. Steinbrücker, F., Sturm, J., Cremers, D.: Volumetric 3D Mapping in Real-time on a CPU. In: *Proceedings of IEEE International Conference on Robotics and Automation*. pp. 2021–2028 (2014) [2](#), [3](#)
25. Sturm, J., Engelhard, N., Endres, F., Burgard, W., Cremers, D.: A Benchmark for the Evaluation of RGB-D SLAM Systems. In: *Proceedings of International Conference on Intelligent Robot Systems*. pp. 573–580 (2012) [9](#), [10](#)
26. Tanskanen, P., Kolev, K., Meier, L., Camposeco, F., Saurer, O., Pollefeys, M.: Live Metric 3D Reconstruction on Mobile Phones. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 65–72 (2013) [4](#)
27. Ulusoy, A.O., Black, M.J., Geiger, A.: Patches, Planes and Probabilities: A Non-local Prior for Volumetric 3D Reconstruction. In: *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*. pp. 3280–3289 (2016) [1](#), [2](#), [4](#), [9](#)
28. Ulusoy, A.O., Geiger, A., Black, M.J.: Towards Probabilistic Volumetric Reconstruction Using Ray Potentials. In: *Proceedings of International Conference on 3D Vision*. pp. 10–18 (2015) [4](#), [9](#)
29. Vogiatzis, G., Hernández, C.: Video-based, Real-time Multi-view Stereo. *Image and Vision Computing* **29**(7), 434–441 (2011) [1](#), [4](#), [5](#), [7](#), [8](#)
30. Whelan, T., Kaess, M., Fallon, M., Johannsson, H., Leonard, J., McDonald, J.: KinectFusion: Spatially Extended KinectFusion. *Robotics & Autonomous Systems* **69**(C), 3–14 (2012) [3](#)
31. Whelan, T., Kaess, M., Johannsson, H., Fallon, M., Leonard, J.J., McDonald, J.: Real-time Large-scale Dense RGB-D SLAM with Volumetric Fusion. *The International Journal of Robotics Research* **34**(4-5), 598–626 (2015) [3](#)

32. Woodford, O.J., Vogiatzis, G.: A Generative Model for Online Depth Fusion. In: Proceedings of European Conference on Computer Vision. pp. 144–157 (2012) [3](#), [9](#)
33. Zeng, M., Zhao, F., Zheng, J., Liu, X.: Octree-based Fusion for Realtime 3D Reconstruction. *Graphical Models* **75**(3), 126–136 (2013) [3](#)
34. Zhou, Q., Koltun, V.: Dense Scene Reconstruction with Points of Interest. *ACM Transactions on Graphics* **32**(4), 112 (2013) [3](#), [9](#), [10](#), [11](#)
35. Zienkiewicz, J., Tsiotsios, A., Davison, A., Leutenegger, S.: Monocular, Real-time Surface Reconstruction using Dynamic Level of Detail. In: Proceedings of International Conference on 3D Vision. pp. 37–46 (2016) [1](#), [3](#)