



# Progressive Acquisition of SVBRDF and Shape in Motion

Hyunho Ha,  Seung-Hwan Baek,  Giljoo Nam  and Min H. Kim 

School of Computing, KAIST, Daejeon, Korea  
{hhha, shwbaek, gjnam}@vclab.kaist.ac.kr, minhkim@kaist.ac.kr

## Abstract

To estimate appearance parameters, traditional SVBRDF acquisition methods require multiple input images to be captured with various angles of light and camera, followed by a post-processing step. For this reason, subjects have been limited to static scenes, or a multiview system is required to capture dynamic objects. In this paper, we propose a simultaneous acquisition method of SVBRDF and shape allowing us to capture the material appearance of deformable objects in motion using a single RGBD camera. To do so, we progressively integrate photometric samples of surfaces in motion in a volumetric data structure with a deformation graph. Then, building upon recent advances of fusion-based methods, we estimate SVBRDF parameters in motion. We make use of a conventional RGBD camera that consists of the colour and infrared cameras with active infrared illumination. The colour camera is used for capturing diffuse properties, and the infrared camera-illumination module is employed for estimating specular properties by means of active illumination. Our joint optimization yields complete material appearance parameters. We demonstrate the effectiveness of our method with extensive evaluation on both synthetic and real data that include various deformable objects of specular and diffuse appearance.

**Keywords:** Inverse rendering, SVBRDF, 3D reconstruction

**ACM CCS:** • Computing methodologies → Reflectance modelling

## 1. Introduction

Capturing material appearance has been a long-lasting research problem in computer graphics. Many specialized hardware systems and software have been proposed to capture appearance parameters, which can be used for photorealistic rendering of real-world objects [DRS10]. Besides the cost of building a specialized hardware setup, a long process of material acquisition is required. It begins with capturing various photometric observation samples with diverse angles of light and camera, resulting in hundreds of images [GHP\*08, HLZ10, NLW\*16, TFG\*13, SSWK13, FHW\*11, TAL\*07, LWS\*13]. The process is then followed by heavy computational processes that include calibration, registration, inverse rendering and so on, often resulting in computation taking numerous hours.

In addition to the system-building cost and the long hours of processing, the entire input images with different angles of light and camera should be captured in advance to be processed all together for the optimization process of inverse rendering. This setup forces the captured object to be static. If the object moves with motion or

is deformed into a different shape, the registration and geometric relationship of input images are broken so that the entire input images should be recaptured from scratch to estimate appearance parameters. The state-of-the-art material appearance acquisition techniques assume that a target object is both *rigid* and *static*. Neither deformation nor motion has been allowed in traditional acquisition methods. To address the aforementioned drawbacks of the existing solutions, we were motivated to capture the material appearance of a dynamic object in motion like a person or any deformable object such as cloth.

Acquiring the appearance of dynamic objects has been achieved by developing a specialized multiview video system [TAL\*07, FHW\*11, LWS\*13]. However, these systems are limited to capturing subjects placed within the multiple light-camera stage. Also, these systems are significantly more expensive than practical acquisition methods. In contrast, we were motivated to devise a practical acquisition solution without requiring any specialized hardware setup, such as a mechanical gantry with two robotic arms or a multiview camera-light stage. To this end, we decided to make use of a conventional RGBD camera for our acquisition setup, following



**Figure 1:** (a) We provide the first-ever method to simultaneously estimate the SVBRDF, shape and motion of dynamic objects using a single RGBD camera. (b)–(e) We obtain both diffuse and specular appearance with our novel joint optimization scheme, based on our hierarchical data structure, which allows us to render captured scenes under novel view and light conditions. Refer to the supplemental video for more results.

the trend of state-of-the-art practical techniques [AWL15, RPG16, HSL\*17, RRF17, WZ15, WWZ16, PNS18, NLGK18].

The conventional RGBD camera that we used in this work is a Kinect 2 sensor that consists of two camera modules: a colour imaging module is an ordinary colour camera to capture red, green and blue colours of objects, and a time-of-flight (TOF) imaging module is an infrared camera to capture the depth information with active illumination of an infrared light module. We utilize the colour camera for capturing diffuse properties and use the pair of the infrared camera and the infrared illumination module to estimate specular properties.

In this work, we propose a progressive estimation of the spatially varying bidirectional reflectance distribution function (SVBRDF) and the shape of a deformable object in motion using a single RGBD camera. As we are using a depth camera, we can estimate the shape and motion vectors of the target object simultaneously while estimating appearance. We introduce a novel architecture to progressively integrate photometric observation samples in motion in a volumetric structure through a deformation graph. Existing works using a single camera can capture SVBRDFs of static objects based on a hierarchical data structure that consists of multiple clusters of similar appearance. To the best of our knowledge, none of these methods can acquire SVBRDF and surface geometry with motion simultaneously. Our method estimates not only geometry with motion but also SVBRDFs.

In addition, the traditional material acquisition methods [GHP\*08, HLZ10, NLW\*16, TFG\*13, SSWK13, FHW\*11, TAL\*07, LWS\*13] require several hours to capture input images of rigid objects. Our novel inverse rendering framework allows us to estimate SVBRDF parameters and shape information *progressively* in interactive time as we build our framework by combining the recent advances of fusion-based methods [NFS15, IZN\*16, GXY\*17] and the practical inverse rendering technique that captures SVBRDF with active illumination [NLGK18, WZ15]. Our progressive acquisition approach does not need to wait for several hours to capture input images. From an application perspective, it does not force the target object to be static until all of the input images are captured. Our method can progressively update both appearance and shape parameters simultaneously. Processing each frame takes less than a half second with a single GPU to estimate

every parameter from photometric samples accumulated through motion vectors.

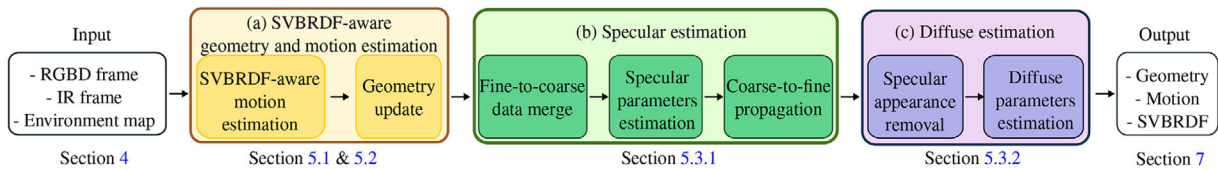
In summary, our method is the first to bridge the gap between SVBRDF acquisition of rigid objects and fusion-based dynamic scanning of diffuse colours, allowing for simultaneous acquisition of SVBRDF and shape in motion. Our main contributions are summarized as follows:

- an architecture to accumulate photometric samples of a dynamic object in a volumetric structure through a deformation graph of motion,
- a joint optimization framework that can estimate SVBRDF, shape and motion simultaneously, and
- a progressive appearance computation framework for inverse rendering.

## 2. Related Work

**Appearance Acquisition of Static Objects.** Traditionally, material appearance of static objects has been effectively acquired with specialized hardware systems that consist of multiple lights or cameras [GHP\*08, HLZ10, TFG\*13, SSWK13, GCHS10, NLW\*16, RRF17, BJTK18]. However, the building cost of such systems is too high to make the acquisition process not available for casual users to have access to this acquisition process. To resolve this issue, practical methods using a single camera have been introduced [AWL15, RPG16, HSL\*17, RRF17, WZ15, WWZ16, SWK19, PNS18, NLGK18]. These methods can capture material appearance by inferring diffuse and specular appearance parameters from multiple observations with different view/light angles. Although being effective for appearance and shape estimation, these methods are limited to capturing *static* objects, meaning objects without any motion. In contrast, we extend the target objects of appearance acquisition to *dynamic* objects through a joint estimation of appearance, shape and motion.

**Multi-Camera Acquisition of Dynamic Objects.** To capture the geometry and appearance of dynamic objects, various specialized multi-camera systems have been proposed. Most previous systems only target diffuse appearance, neglecting specular appearance [WVT12, DKD\*16, DDF\*17, XSH\*19]. There have been few attempts to estimate the complete appearance of diffuse and specular



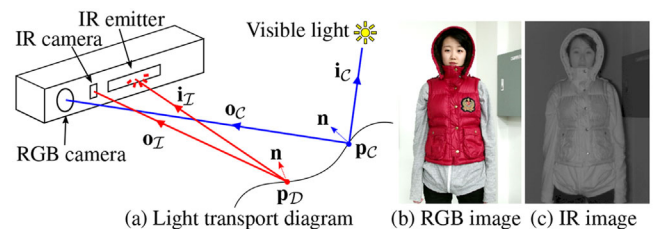
**Figure 2:** For each frame, our method takes inputs of RGB, IR and depth images from a conventional RGBD sensor (Kinect 2), in addition to the static environment map (only captured once at the calibration stage). (a) We first estimate motion fields and scene geometry with consideration of SVBRDF. (b) Specular parameters are then estimated by exploiting the hierarchical data structure. (c) Given the specular estimates, residual observation is fitted to the diffuse component, resulting in the diffuse albedo estimates. This framework runs in an online manner, producing geometry, motion and SVBRDF per frame as output.

components simultaneously [FHW\*11, TAL\*07, LWS\*13]. However, these multiview methods require very expensive acquisition systems with multiple cameras and lights. They are also offline methods with high computational costs. In contrast, our method progressively estimates diffuse/specular parameters, geometry and motion in an *online* manner using a *single* RGBD camera, which makes them more practical.

**Single-Camera Acquisition of Dynamic Objects.** Although estimating the shape and motion of dynamic objects from a single camera has been extensively studied, estimating appearance simultaneously is known to be challenging [NFS15, SBC117, LZG18, YGX\*17, YZG\*18, ZYL\*18, YGX\*17, YZG\*18, ZYL\*18, YZZ\*19]. Only a few studies have been attempted to capture diffuse components either in diffuse albedo [GXY\*17] or shaded diffuse colours [DDF\*17, IZN\*16, SBI18]. In addition, there are practical acquisition methods that allows users to capture appearance in a simple setup. Lin *et al.* [LPG19] estimate appearance parameters by simply capturing HDR images of an object and a light probe. Dong *et al.* [DCP\*14] capture SVBRDFs with known geometry from an input video with motion. The main technical challenge of simultaneously estimating *specular appearance* is that the number of light/view samples in each frame is not sufficient for appearance estimation. In order to overcome this, we utilize the active infrared illumination in the TOF camera for estimating specular parameters and integrate photometric samples into a hierarchical data structure. This enables us to reconstruct all appearance parameters, geometry and motion of dynamic objects simultaneously.

### 3. Overview

Our method progressively estimates SVBRDF, geometry and motion of a deformable object in a frame-by-frame manner. Based on the traditional fusion framework [NFS15, GXY\*17], using an RGBD camera, our method accumulates photometric samples of the target object in our hierarchical data structure. The data structure allows us to estimate the full appearance of the object using a small number of frames. Using the estimated appearance parameters of the object, our algorithm progressively updates the appearance information more accurately over time through the object's motion vectors. Figure 2 describes the overview of our method.



**Figure 3:** (a) Environmental scene illumination in visible RGB channels is reflected at object surfaces and captured by the colour camera on the RGBD sensor. Estimating the specular component from the RGB image (b) is challenging due to lack of the view/light direction information. In contrast, the point IR illumination of the depth camera enables effective reconstruction of specular appearance from the IR image (c).

### 4. Acquisition Setup

In order to make our acquisition system practical, we make use of two off-the-shelf imaging devices: a conventional RGBD camera (Kinect 2) and a 360° camera (Ricoh Theta). The RGBD camera is the main device to capture SVBRDF and shape in motion, and the spherical camera is used to capture the environment illumination of the scene.

We chose the RGBD camera because the camera consists of the colour and infrared cameras with active infrared illumination. First, the colour camera can be used for capturing diffuse colour properties under the scene ambient illumination. Second, unlike the previous generation of RGBD cameras (Kinect 1 or PrimeSense), the second generation of the Kinect sensor includes the TOF camera module to estimate depth. The camera API allows us to access to raw infrared image data, time-modulated phase images under active infrared (IR) illumination, without having spatial modulation artefacts shown in the previous generation. The clear infrared image data under the known active illumination can be utilized to estimate view-light-dependent reflectance property, that is, specular albedo and surface roughness. As the angle between active infrared illumination and infrared TOF sensor in the RGB-D camera is approximately 5° at a distance of around 1 m, this could be sufficiently wide to capture most SVBRDF except the Fresnel effect, as discussed in [NJR15, NLGK18].

**Table 1:** Symbols and notations used in the paper.

	Symbol	Description	
Image	$t$	Frame number	
	$u$	Image pixel	
	$\mathcal{P}$	Pixel domain	
	$\mathcal{P}_D^t, \mathcal{P}_C^t$	Set of visible pixels at the depth camera space and the colour camera space at $t$	
	$\tilde{u}_D, \tilde{u}_C$	Corresponding pixel of a rendered image pixel $u$ in the depth and the colour camera space	
	$\tilde{u}_{x_D}$	Corresponding depth pixel of a voxel $x$ in the depth camera space	
	$D^t, C^t, I^t$	Depth, colour and IR image at the frame $t$	
	$\Gamma, Y$	Chromaticity and luminance of the colour image	
	$\tilde{V}_D^t, \tilde{N}_D^t$	Vertex map and normal map of the warped mesh at the depth camera space at $t$	
	$\tilde{V}_C^t, \tilde{N}_C^t$	Vertex map and normal map of the warped mesh at the colour camera space at $t$	
	$V_D^t, N_D^t$	Vertex map and normal map of the depth image at $t$	
	$O_C^t$	View direction of $\tilde{V}_C^t$ to the colour camera at $t$	
	Transformation	$\mathcal{K}, \mathcal{D}, \mathcal{C}$	Canonical, depth (IR) and colour camera space
		$P$	Perspective projection
$T_i$		Deformation graph transformation matrix at the node $i$	
$K_D, K_C$		Depth (IR) camera, Colour camera intrinsic matrix	
$T_{\mathcal{K} \rightarrow \mathcal{D}}^t$		Canonical space to depth camera space transformation matrix at $t$	
$T_{\mathcal{D} \rightarrow \mathcal{K}}^t$		Depth camera space to canonical space transformation matrix at $t$	
$T_{\mathcal{K} \rightarrow \mathcal{C}}^t$		Canonical space to colour camera space transformation matrix at $t$	
$T_{\mathcal{D} \rightarrow \mathcal{C}}^t$		Depth camera space to colour camera space transformation matrix at $t$	
Geometry	$p_{\mathcal{K}}^t, \tilde{p}_{\mathcal{K}}^t, n_{\mathcal{K}}^t, \tilde{n}_{\mathcal{K}}^t$	Point in the canonical space and its warped point and its normal at $t$	
	$x, x_{\mathcal{K}}$	Voxel and its canonical position	
	$\tilde{x}_D^t, \tilde{x}_C^t, \tilde{n}_D^t, \tilde{n}_C^t$	Position of a voxel $x$ at the depth camera and the colour camera space and its normal at $t$	
	$\mathcal{T}$	TSDF structure	
	$d_{\mathcal{T}}, \omega_{\mathcal{T}}$	Signed distance value and its weight	
	$\mathcal{T}_D^t, \mathcal{T}_C^t$	Set of visible voxel at the depth space and the colour space at $t$	
	$\mathcal{V}_{\mathcal{K}}, \tilde{\mathcal{V}}_{\mathcal{K}}$	Canonical frame and warped mesh	
Motion	$\mathcal{G}$	Deformation graph	
	$\mathcal{W}^t$	Motion field at frame $t$	
	$q_i, \sigma_i, w_i$	Position, radius and the weight of the $i$ th deformation graph node	
	$\Phi$	Kernel function	
Parameters	$\lambda$	Regularize parameter	
	$\kappa$	IR emitter illumination	
	$\gamma_C, \gamma_I$	Gamma value of the colour and the IR camera	
	$v_x^t, \omega_x^t$	Half angle buffer value and its weight of a voxel $x$ at $t$	
	$v_m^t, \omega_m^t$	Half angle buffer value and its weight of a cluster $m$ at $t$	
Reflectance	$\mathcal{A}$	Appearance of the canonical space model	
	$\rho_d, \rho_s, \alpha$	Diffuse albedo, specular albedo and specular roughness	
	$\tilde{o}_D^t, \tilde{o}_C^t$	View direction of $\tilde{x}_D^t$ to the depth camera at $t$ and $\tilde{x}_C^t$ to the colour camera at $t$	
	$\tilde{l}_{L,D}^t$	Light direction of $\tilde{x}_D^t$ to the IR emitter at $t$	
	$\theta_l, \theta_o, \theta_h$	Zenith angle between the normal plane and the light, view and half vector direction	
	$f_r, f_d, f_s$	Reflectance, diffuse reflectance, and specular reflectance function	
	$H_k, l_k$	Spherical harmonics basis function and coefficient	
	$B, S$	Diffuse and specular reflection	
	$\mathcal{M}$	Cluster	

Figure 1(a) presents our acquisition setup. Figure 3 depicts light transport in our imaging setup. First, to estimate the incident illumination of the scene, we capture an environment map using a 360° camera. Second, for each frame, an RGB colour frame is captured by the colour camera module in Kinect 2 under the ambient scene illumination. An infrared frame is captured by the TOF camera module under active illumination of the TOF camera module in the RGBD camera. We use both colour and infrared frames in the video stream to estimate the diffuse and specular parameters of SVBRDF.

## 5. Capturing SVBRDF and Shape in Motion

To capture SVBRDF on a non-rigid object using a single RGBD camera, we introduce a two-step framework of dynamic inverse rendering. First, we calculate motion fields by comparing both the appearance and geometry of the current frame with those properties of the static model continuously accumulated from previous frames. Second, using the estimated motion fields, we update three properties sequentially in each frame: geometry, specular reflectance

parameters and diffuse albedo, sequentially in each frame. Refined geometry is used for estimating the parameters of specular albedo and surface roughness from photometric samples under active infrared illumination of the TOF camera. We cluster these parameters in the deformation graph structure to estimate SVBRDFs per cluster. Our SVBRDF acquisition progressively captures diffuse albedo, specular coefficient, specular roughness, geometry and motion frame by frame under visible environment illumination.

## 5.1. Background

### 5.1.1. Voxel Grid and SVBRDF Model

**Voxel Grid.** We make use of a truncated signed distance function (TSDF) volume  $\mathcal{T}$  [CL96] to store the estimated appearance and shape information. We reconstruct actual geometry and appearance properties in the 3D voxel grid of TSDF, which is a set of voxel  $\mathbf{x} \in \mathbb{N}^3$  that consists of two sets of properties:  $\mathcal{T} = \{\mathcal{V}, \mathcal{A}\}$ . First, geometry  $\mathcal{V} = \{[d_{\mathcal{T}}, \omega_{\mathcal{T}}]\}^t$  at frame  $t$  is defined as a signed distance value  $d_{\mathcal{T}}$  and its weight  $\omega_{\mathcal{T}}$ . Second, appearance  $\mathcal{A} = \{[\rho_d, \rho_s, \alpha]\}^t$  is a set of diffuse albedo  $\rho_d$ , specular albedo  $\rho_s$  and roughness parameter  $\alpha$ . As appearance parameters are surface properties, we only accumulate the appearance parameters to the voxels near to surface:  $|d_{\mathcal{T}}| < 0.5\tau$ , where  $\tau$  is the truncate value.

**Reflectance.** We formulate appearance as SVBRDF, where the reflectance function  $f_r$  represents the isotropic Ward model [War92] at vertex point  $\mathbf{p}$  in the voxel grid. The diffuse term  $f_d$  represents individual diffuse albedo  $\rho_d$  per point, and the specular term  $f_s$  shares specular albedo (a.k.a. specular coefficient)  $\rho_s$  and roughness parameter  $\alpha$  per cluster  $\mathcal{M}$  as follows:

$$f_r(\mathbf{i}, \mathbf{o}; \rho_d, \rho_s, \alpha, \mathbf{n}, \mathbf{p}) = f_d(\rho_d, \mathbf{p}) + f_s(\mathbf{i}, \mathbf{o}; \rho_s, \alpha, \mathbf{n}, \mathbf{p}),$$

$$= \frac{\rho_d}{\pi} + \frac{\rho_s}{4\pi\alpha^2\sqrt{\cos\theta_i\cos\theta_o}} e^{-(\tan^2(\theta_h)/\alpha^2)}, \quad (1)$$

where  $\mathbf{i}$  and  $\mathbf{o}$  are the incident light vector and the view vector,  $\mathbf{h} = (\mathbf{i} + \mathbf{o})/|\mathbf{i} + \mathbf{o}|$  is the half-angle vector,  $\theta_i$ ,  $\theta_o$  and  $\theta_h$  are an angle between the normal  $\mathbf{n}$  and each vector  $\mathbf{i}$ ,  $\mathbf{o}$  and  $\mathbf{h}$ , respectively, at point  $\mathbf{p}$ .

**Rendering.** With an objective of per-frame inverse rendering, we capture an HDR environment map as scene illumination over solid angle  $\Omega$  as input. Suppose we have incident light  $L(-\mathbf{i}; \mathbf{p})$  over angle  $\Omega$ . Using the rendering equation [Kaj86], we calculate reflected light  $L(\mathbf{o}; \mathbf{p})$  as

$$L(\mathbf{o}; \mathbf{n}, \mathbf{p}) = \int_{\Omega} L(-\mathbf{i}; \mathbf{p}) f_r(\mathbf{i}, \mathbf{o}; \rho_d, \rho_s, \alpha, \mathbf{n}, \mathbf{p}) (\mathbf{n} \cdot \mathbf{i}) d\mathbf{i}$$

$$\approx B(\rho_d, \mathbf{n}, \mathbf{p}) + S(\mathbf{o}; \rho_s, \alpha, \mathbf{n}, \mathbf{p}). \quad (2)$$

First, for computational efficiency, we approximate diffuse reflection as spherical harmonics of radiosity [WZN\*14, RH01] from given normals, assuming fixed environment illumination:  $B(\rho_d, \mathbf{n}, \mathbf{p}) = \rho_d \sum_{k=0}^8 l_k H_k(\mathbf{n})$ , where  $l_k$  are the nine spherical harmonics coefficients of incident environment illumination (up to the second order) over  $\Omega$ , and the spherical harmonics basis functions  $H_k(\mathbf{n})$  take normals  $\mathbf{n}$  as input to calculate diffuse shading in the global space. Second, we calculate specular reflec-

tion  $S(\mathbf{o}; \rho_s, \alpha, \mathbf{n}, \mathbf{p}) = \int_{\Omega} L(-\mathbf{i}, \mathbf{p}) f_s(\mathbf{i}, \mathbf{o}; \rho_s, \alpha, \mathbf{n}, \mathbf{p}) (\mathbf{n} \cdot \mathbf{i}) d\mathbf{i}$  by integrating the spherical illumination map using uniform sampling of the upper hemisphere in the normal space.

### 5.1.2. Capturing Shape

Simultaneously estimating SVBRDF, geometry and motion is a chicken-and-egg problem because they are strongly correlated. Once the first-frame observation of the RGBD camera is stored in the canonical space, we begin with estimating the per-frame motion field by formulating the following optimization problems. Before explaining SVBRDF estimation in motion, we briefly explain how to estimate the motion field to accumulate dynamic photometric samples in our hierarchical data structure. This motion part is inherited from the traditional fusion-based framework [NFS15]. Refer to Table 1 for symbols and notations used in this paper.

**Global Registration.** To improve robustness, we first estimate global transformation that registers the input frame of a depth camera to the voxel grid in each frame, which is formulated as a 6-DOF rigid body transformation (RBT) matrix  $\mathbf{T}_{\mathcal{D} \rightarrow \mathcal{K}}^t \in \mathbf{SE}(3)$  such that point  $\mathbf{p}_{\mathcal{D}}^t$  in the depth camera space  $\mathcal{D}$  at frame  $t$  is transferred into the canonical space of voxel grid  $\mathcal{K}$  via  $\mathbf{p}_{\mathcal{K}}^t = \mathbf{T}_{\mathcal{D} \rightarrow \mathcal{K}}^t \mathbf{p}_{\mathcal{D}}^t$ . The matrix can be optimized by solving the iterative closest point (ICP) method [RL01].

**Capturing Shape via Motion.** Following the previous work of DynamicFusion [NFS15], we first estimate the local non-rigid motion fields per frame and update the shape of the deformable objects. We define a motion field  $\mathcal{W}$  from the canonical space  $\mathcal{K}$  to the current warped frame  $t$  as  $\mathcal{W} = \{[\mathbf{q}_i, \sigma_i, \mathbf{T}_i]\}^t$ , where  $\mathbf{q}_i$  is a position of  $i$ th node from the total  $N$  number of nodes ( $i \in \{1, \dots, n\}$ ) in the deformation graph  $\mathcal{G}$ ,  $\sigma_i \in \mathbb{R}^+$  is a radius parameter for the distance weight  $w_i$  between node  $\mathbf{q}_i$  and point  $\mathbf{p}_{\mathcal{K}}$  in the canonical space:  $w_i(\mathbf{p}_{\mathcal{K}}, \sigma_i) = \exp(-\|\mathbf{p}_{\mathcal{K}} - \mathbf{q}_i\|^2 / (2\sigma_i^2))$  and  $\mathbf{T}_i \in \mathbf{SE}(3)$  is a 6-DOF RBT of the  $i$ th node. The motion field  $\mathcal{W}^t$  at a point  $\mathbf{p}_{\mathcal{K}}$  is defined by dual-quaternion blending [KCvO07] using the  $k$ -nearest neighbour nodes with its convex weights. The motion field  $\mathcal{W}^t$  warps a point  $\mathbf{p}_{\mathcal{K}}$  and a normal  $\mathbf{n}(\mathbf{p}_{\mathcal{K}})$  in the canonical space by  $[\tilde{\mathbf{p}}_{\mathcal{K}}^t, 1]^T = \mathcal{W}^t(\mathbf{p}_{\mathcal{K}})[\mathbf{p}_{\mathcal{K}}^t, 1]^T$  and  $[\tilde{\mathbf{n}}(\mathbf{p}_{\mathcal{K}})^T, 0]^T = \mathcal{W}^t(\mathbf{p}_{\mathcal{K}})[\mathbf{n}(\mathbf{p}_{\mathcal{K}})^T, 0]^T$ . Given depth image  $\mathbf{D}^t$  and the estimated warp motion field, we obtain a weighted average of the projective TSDF values for every voxel  $\mathbf{x}$  to reconstruct the shape. Finally, we conduct the marching cube algorithm on the TSDF volume to create a polygonal mesh model per frame and update deformation graph. For more detail, refer to the previous work [NFS15] and the supplemental material.

## 5.2. Estimating Motion with SVBRDF

State-of-the-art fusion methods [ZNI\*14, DNZ\*17, NFS15, GXY\*17] evaluate only diffuse colour and geometry differences to estimate motion field. In contrast, we can estimate the current motion field  $\mathcal{W}^t$  by minimizing the following energy function making use of given geometry and SVBRDF:

$$E_{\text{motion}}(\mathcal{W}^t) = E_{\text{depth}} + \lambda_{\text{dreg}} E_{\text{dreg}} + \lambda_{\text{pcolour}} E_{\text{pcolour}}, \quad (3)$$

where  $E_{\text{depth}}$  and  $E_{\text{dreg}}$  are the data term and its regularizer for geometry,  $E_{\text{pcolour}}$  is our novel data term for SVBRDF.  $\lambda_{\text{dreg}}$  and  $\lambda_{\text{pcolour}}$  are the corresponding weights.

**Geometric Energy.** Our geometric energy terms  $E_{\text{depth}}$  and  $E_{\text{dreg}}$  are similar to those terms used in [NFS15].  $E_{\text{depth}}$  optimizes the motion parameter by minimizing the plane-normal distance between the warped mesh from the previous frame and its correspondence point in the current depth image. To enforce the local smoothness of motion and prevent overfitting,  $E_{\text{dreg}}$  minimizes the distance when the node is warped by its own motion parameter and when it is warped by the motion of the neighbouring nodes. Refer to [NFS15] or the supplemental document for more details.

**Colour Energy.** Assuming that SVBRDF of the captured object does not change over time, our novel motion estimation term  $E_{\text{pcolour}}$  considers object appearance to enforce the photometric consistency of object surfaces at the  $i$ th node in the camera space  $\mathcal{C}$  as follows:

$$E_{\text{pcolour}}(\mathcal{W}^i) = \sum_{u \in \mathcal{P}_C^i} \|\mathbf{C}^i(\tilde{u}_C) - L^i(\tilde{\mathbf{O}}_C^i(u); \tilde{\mathbf{N}}_C^i(u), \tilde{\mathbf{V}}_C^i(u))\|_2^2, \quad (4)$$

where  $\mathcal{P}_C^i$  is a set of visible pixels  $u$  obtained by rendering the warped static model to the current colour camera space  $\mathcal{C}^i$ ,  $\tilde{\mathbf{V}}_C^i: \mathbb{N}^2 \rightarrow \mathbb{R}^3$  is the vertex map of the warped mesh  $\tilde{\mathcal{V}}_C^i$  transformed by  $\mathbf{T}_{\mathcal{K} \rightarrow \mathcal{C}}^i$  from the canonical space to the current colour camera space,  $\tilde{\mathbf{O}}_C^i$  is the view direction of  $\tilde{\mathbf{V}}_C^i$  to the colour camera,  $\tilde{\mathbf{N}}_C^i: \mathbb{N}^2 \rightarrow \mathbb{R}^3$  is the normal map of  $\tilde{\mathcal{V}}_C^i$  transformed by  $\mathbf{T}_{\mathcal{K} \rightarrow \mathcal{C}}^i$ ,  $\tilde{u}_C = P(\mathbf{K}_C \tilde{\mathbf{V}}_C^i(u))$  is the pixel in the colour image  $\mathcal{C}^i$  that corresponds to  $u$ ,  $\mathbf{K}_C$  is the intrinsic matrix of the colour camera, and the reflected light  $L^i = B^i + S^i$  is rendered by Equation (2). As unestimated specular components degrade the estimate quality of the estimating motion, this term helps to correctly estimate the photometric difference between a colour image and our reconstructed objects. Refer to Figure 9 to see how geometric accuracy has been improved by accounting for SVBRDF in estimating motion.

**Motion Optimization.** In order to solve Equation (3), we reformulate it as the sum of squared residuals  $f$  so that we can define a new vector field  $\mathbf{F}$  to find out the vector of motion parameters  $\mathcal{X}$ , satisfying:  $E(\mathcal{X}) = \sum f(\mathcal{X})^2 = \|\mathbf{F}(\mathcal{X})\|^2$ . Then, the optimization formulation can be solved by the Gauss–Newton method. The reformulated optimization needs the linearization of three terms of motion, diffuse reflectance and specular reflectance.

For the first two approximation steps of motion and diffuse colour, we follow an existing method of using twist representation [MSZ94] that represents each node’s motion parameters  $\mathcal{X}$  (3D for rotation and 3D for translation), and converting it to  $\mathbf{SE}(3)$  using an exponential map. We also linearize the diffuse colour image using the first-order Taylor approximation [WVT12, NFS15, GXY\*17].

However, linearizing our novel SVBRDF term is not trivial. Different from view-invariant diffuse reflection  $B^i$  at frame  $t$ , specular reflection  $S^i$  at vertex  $\tilde{\mathbf{V}}_C^i$  depends on the outgoing angle variable  $\tilde{\mathbf{O}}_C^i$  with appearance parameters  $(\rho_s, \alpha, \mathbf{n})$  and also is formulated by the integration of the incident light (Equations (1) and (2)). Therefore, the computational cost for the direct minimization of Equation (4) with the SVBRDF term is highly expensive. Instead, we first ren-



**Figure 4:** (a) and (b) Input photographs of 720th and 740th frames. (c) Our estimated motion fields showing the deformation of the cloth at the 740th frame.

der specular reflection  $S^i$  with given environment illumination in the current colour camera space  $\mathcal{C}^i$  and then substitute  $S^i$  from captured colour image  $\mathbf{C}^i$  for comparison with pure radiosity  $B^i$ , based on Equation (2). This solution increases colour optimization very efficiently and enabling us to consider SVBRDF when estimating motion fields.

Finally, in each Gauss–Newton iteration, we find parameters of  $\Delta\mathcal{X}$  by solving a linear least-squares problem [DNZ\*17]:

$$\Delta\hat{\mathcal{X}} = \arg \min_{\Delta\mathcal{X}} \|\mathbf{F}(\mathcal{X}^{(j-1)}) + \mathbf{J}_F(\mathcal{X}^{(j-1)}) \cdot \Delta\mathcal{X}\|. \quad (5)$$

To obtain  $\Delta\hat{\mathcal{X}}$ , we set the partial derivatives of the above equation with respect to  $\Delta\mathcal{X}$  as zero to solve the following equation:  $\mathbf{J}_F^T(\mathcal{X}^{(j-1)})\mathbf{J}_F(\mathcal{X}^{(j-1)}) \cdot \Delta\hat{\mathcal{X}} = -\mathbf{J}_F^T(\mathcal{X}^{(j-1)})\mathbf{F}(\mathcal{X}^{(j-1)})$ . We solve this problem with pre-conditioned conjugate gradient method (Section 6). Finally, we update motion field as follows:  $\mathbf{T}_i^j = e^{\Delta\hat{\mathcal{X}}} \cdot \mathbf{T}_i^{j-1}$ . Figure 4 shows an example of the estimated motion field using our SVBRDF-aware motion optimization.

### 5.3. Capturing SVBRDF in Motion

The state-of-the-art methods for estimating material appearance have focused on SVBRDF of static objects [PNS18, NLGK18] or only diffuse albedo of dynamic objects [GXY\*17]. As we estimate per-vertex motion and shape, we then estimate complete SVBRDF parameters per vertex in a progressive way through our novel optimization method.

#### 5.3.1. Specular Parameters

There are two main technical challenges for estimating specular parameters: First, specular reflectance depends on both light and view directions, whereas diffuse reflectance is a constant. In particular, specular parameter estimation requires a set of multiple photometric samples with known light and view directions before optimization. Second, per-frame progressive optimization of specular parameters is therefore supposed to suffer from a lack of samples more than the traditional offline methods. The appearance parameters of the same materials need to be shared with spatial and temporal coherence for efficient sampling. We handle these challenges as follows.

**Point-Light Illumination for Specular Reflection.** As mentioned earlier, in an RGBD camera, there is a TOF camera module that consists of an infrared light and an infrared camera to measure depth

(Figure 3). We utilize the pair of the infrared illumination and the infrared camera module to capture photometric samples to estimate specular parameters.

First, we have geometrically calibrated these two devices beforehand to obtain the light and view vectors ( $\mathbf{i}_I$ ,  $\mathbf{o}_I$ ). The relative position and orientation of both  $\mathbf{i}_I$  and  $\mathbf{o}_I$  with respect to the surface geometry are obtained using the estimated motion field. Given the known light and view vectors in the normal space, we can remove the integral over hemisphere  $\Omega$  in Equation (2) using the point light assumption:

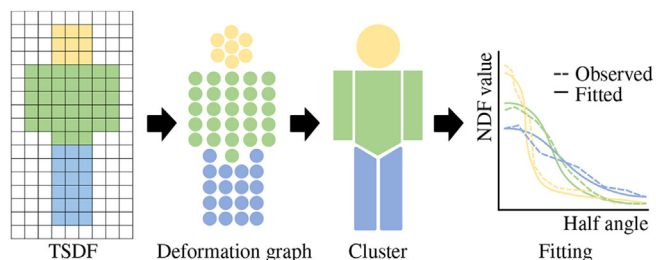
$$S(\mathbf{i}_I, \mathbf{o}_I; \rho_s, \alpha, \mathbf{n}, \mathbf{p}) = L_{i_I}(\mathbf{p})f_s(\mathbf{i}_I, \mathbf{o}_I; \rho_s, \alpha, \mathbf{n}, \mathbf{p})(\mathbf{n} \cdot \mathbf{i}_I). \quad (6)$$

It allows us to solve the inverse problem efficiently per frame, assuming that the surface roughness of microfacets is consistent in both visible and infrared illumination, following [WZ15, PNS18].

**Hierarchical Data Structure.** Different from capturing diffuse albedo, estimating specular parameters requires dense observation samples, and thus existing SVBRDF methods [TAL\*07, LWS\*13, WZ15, PNS18] have used a hierarchical data structure to accumulate sparse samples of specular appearance per each cluster to be used for inferring the specular parameters. In addition, existing dynamic fusion methods [NFS15, GXY\*17] make use of a hierarchical data structure to regularize motion vectors of moving objects. In this work, to estimate the SVBRDF of objects in motion, we combine these two data structures into a novel hierarchical data structure that allows us to estimate motion vectors and appearance parameters together. The structure consists of three main components: surface clusters, deformation graphs and a TSDF volume, where diffuse albedo is estimated per voxel, and specular parameters are estimated per cluster (a set of deformation graph nodes that are associated with motion fields) by assuming that surfaces are dichromatic, and that roughness is locally smooth [WZ15, PNS18].

Once these attributes are optimized per frame, they are interpolated to each vertex in the static model. Our hierarchical structure is beneficial in two aspects: First, we can efficiently estimate both appearance and geometry *in motion* per frame, which requires *expensive* optimization, by working on the small number of clusters compared to the number of voxel grids. Second, we can achieve observations with various angles of  $\theta_h$  to optimize SVBRDF parameters per frame by working on a large range of surfaces with potentially different angles of  $\theta_h$ . Figure 5 visualizes our hierarchical data structure that accumulate photometric samples.

**Fine-to-Coarse Sample Accumulation.** We accumulate these photometric samples in the hierarchical data structure of the *half-angle buffer* based on spatiotemporally coherent clustering using the motion fields. We first store the reflection observations of the infrared point light in the fine-grained TSDF voxel grid. To this end, we first warp the positions of the canonical voxels  $\mathbf{x}_K$  into the current depth camera frame via  $\tilde{\mathbf{x}}_D = \mathbf{T}_{K \rightarrow D}^t \mathcal{W}^t(\mathbf{x}_K) \mathbf{x}_K$ . We then calculate the perspective projection of  $\tilde{\mathbf{x}}_D$  to check visibility and correspondence of  $\tilde{u}_{x_D} = P(\mathbf{K}_D \tilde{\mathbf{x}}_D)$  with respect to camera pixels  $\mathbf{I}^t$ . Once we find out the corresponding camera intensity  $\mathbf{I}^t(\tilde{u}_{x_D})$ , we calculate the specular intensity  $v$  with respect to the half-angle vector angle  $\theta_h$  (a.k.a. the discrete normal distribution function (NDF)) by normalizing the



**Figure 5:** We accumulate shape and SVBRDF parameters in a hierarchical data structure. First, we store every observation from the RGBD camera into the high-resolution TSDF structure. We then transfer the observation into the deformation graph structure for efficient appearance estimation. Nodes are associated with motion fields to yield the spatiotemporal coherence of appearance estimates. Finally, the deformation nodes are clustered, providing enough samples for fitting BRDF parameters for each cluster.

gamma-corrected intensity with shading  $1/(\mathbf{n} \cdot \mathbf{i}_I)$  and distance  $d^2$  at point  $\tilde{\mathbf{x}}_D$  as follows:

$$v = \frac{d^2(\tilde{\mathbf{x}}_D)}{\kappa} \cdot \frac{(\mathbf{I}^t(\tilde{u}_{x_D}))^{\gamma_I}}{\tilde{\mathbf{n}}_D^t \cdot \tilde{\mathbf{i}}_{I,D}^t}, \quad (7)$$

where  $\tilde{\mathbf{n}}_D$  is a normal at  $\tilde{\mathbf{x}}_D$ ,  $\tilde{\mathbf{i}}_{I,D}^t$  is incident IR illumination vector at  $\tilde{\mathbf{x}}_D$ ,  $\gamma_I$  is the infrared camera gamma and  $\kappa$  is a normalization constant. Both  $\gamma_I$  and  $\kappa$  are calibrated, following [PNS18]. We assume that the infrared emitter and receiver are close enough that both  $\mathbf{i}$  and  $\mathbf{o}$  are the same as  $\mathbf{h}$  to simplify Equation (7) similar to [WZ15]. Per-voxel specular reflectance,  $v_x^t$ , of point  $\mathbf{x}$  at current frame  $t$  is updated in the half-angle buffer through weighted average in the static model:

$$v_x^t(\theta_h) = \frac{v \cdot \omega + v_x^{t-1}(\theta_h) \cdot \omega_x^{t-1}(\theta_h)}{\omega + \omega_x^{t-1}(\theta_h)}, \quad (8)$$

where  $\omega = \Phi_{\text{bell}}(u) \cdot \tilde{\mathbf{n}}_D^t \cdot \tilde{\mathbf{o}}_D^t$ ,  $\tilde{\mathbf{o}}_D^t$  is camera view direction at  $\tilde{\mathbf{x}}_D$ ,  $\Phi_{\text{bell}}$  is the bell-shaped filter kernel to suppress extreme noise at the edge of the image. We also update the corresponding weight as follows:  $\omega_x^t(\theta_h) = \omega + \omega_x^{t-1}(\theta_h)$ . As we estimate specular parameters in the hierarchical data structure, we lift the discrete NDF values stored in the high-resolution TSDF structure to the deformation graph's nodes. Specifically, we assign the target deformation node to a TSDF voxel based on the diffuse albedo values of the node and the voxel. We then cluster deformation graph nodes  $\mathbf{q}_i$  with normalized diffuse albedo using the  $k$ -mean clustering algorithm ( $k$  varies up to eight).

**Specular Parameters Optimization.** For each cluster  $m \in \mathcal{M}^t$ , we estimate infrared diffuse albedo  $\rho_{d_I}$  by finding out the minimum value of  $v_m^t(\theta_h)$  such that  $\frac{\Phi_{\text{box}}(v_m^t(\theta_h+1))}{\Phi_{\text{box}}(v_m^t(\theta_h))} \geq 1 + \epsilon$ , where  $\Phi_{\text{box}}$  is the box filter kernel, and  $\epsilon$  is a user-defined value (0–0.01). We then estimate  $\hat{\alpha}(m)$  and initial  $\hat{\rho}_s(m)$  of each cluster  $m$  by minimizing the objective function:

$$\underset{\alpha, \rho_s}{\text{minimize}} \sum_{\theta_h=0}^{\pi/3} \left( \omega' |v_m^t(\theta_h) - \rho_{d_I} - f_s(\theta_h, \alpha, \rho_s)|^2 \right), \quad (9)$$

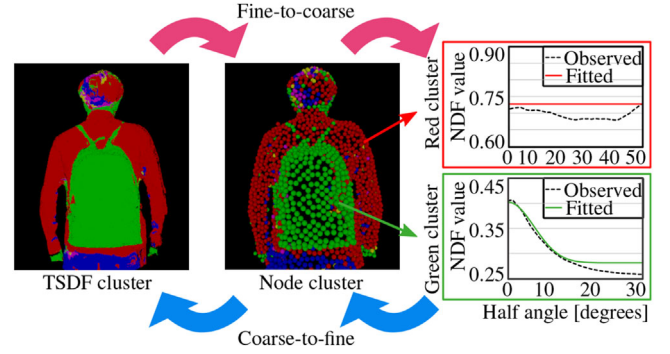
where  $\omega' = \cos^2(\theta_h) \max(\omega_{\max}, \omega_z^t(\theta_h))$  when  $\omega_{\max}$  is a pre-defined clamp parameter and  $\omega_z^t$  is the sums of  $\omega_x^t$  that belong to the cluster  $z$ . Here,  $\omega$  indicates the level of confidence of observation. We set  $\omega_{\max}$  to a certain level empirically to exclude the diffuse-like observation from the regression of the specular parameter. We solve the optimization problem through the brute-force search with a tabulated function  $f_s$  (70 levels:  $0 - 70^\circ$ ) in Equation (1), pre-computed with discrete  $\alpha$  (66 levels:  $0.05 - 0.7$ ) and  $\rho_s$  (100 levels:  $0.01 - 1.00$ ). Note that the deformation graph's nodes are associated with the motion field, allowing for spatiotemporal coherence. Figure 6 shows the estimated clusters, accumulated observations for each cluster and our fitting results.

**Specular Albedo Adjustment.** We utilize the active infrared illumination and the infrared camera to estimate the specular parameters: specular albedo  $\hat{\rho}_s$  and roughness  $\hat{\alpha}$ . However, the albedo of the infrared wavelength is independent of that of the visible wavelength. To estimate specular reflection captured by the RGB colour camera, we estimate the albedo scalar  $\lambda_{I \rightarrow C}$  that adjusts the infrared albedo to the visible specular albedo, that is,  $\lambda_{I \rightarrow C}$  is then multiplied to the infrared specular albedo  $\rho_s$ , yielding visible albedo  $\hat{\rho}_s$ . Note that the infrared roughness parameter  $\alpha$  is independent of albedo so that the same value is copied to the visible roughness  $\hat{\alpha}$ . The albedo scalar  $\lambda_{I \rightarrow C}(m)$  of each cluster  $m$  is estimated as

$$\lambda_{I \rightarrow C}(m) = \frac{\sum_{\mathbf{x} \in \mathcal{T}_C^t \cap \mathcal{M}_m} (\max(Y(\mathbf{C}^t(\tilde{\mathbf{u}}_{\mathbf{x}_c})) - Y(\mathbf{B}^t(\tilde{\mathbf{x}}_c^t)), 0))}{\sum_{\mathbf{x} \in \mathcal{T}_C^t \cap \mathcal{M}_m} S_C^t(\hat{\rho}_{s,I}(m), \hat{\alpha}_I(m), \tilde{\mathbf{x}}_c^t)}, \quad (10)$$

where  $\mathcal{T}_C^t$  is a set of visible surface voxels  $\mathbf{x}$  warped to the current  $\mathbf{C}^t$ ,  $\mathcal{M}_m$  is a set of voxels which cluster to  $m$ ,  $\tilde{\mathbf{x}}_c^t = \mathbf{T}_{\mathcal{D} \rightarrow \mathcal{C}} \tilde{\mathbf{x}}_D^t$  is a voxel transformed from the current depth camera space  $\mathcal{D}^t$  to the colour camera space  $\mathcal{C}^t$ ,  $\mathbf{B}^t(\tilde{\mathbf{x}}_c^t) = \mathbf{B}^t(\rho_d(\tilde{\mathbf{x}}_c^t), \mathbf{n}, \tilde{\mathbf{x}}_c^t)$  is diffuse shading rendering using the diffuse albedo of voxel at  $\mathbf{C}^t$ ,  $Y(\cdot)$  is a luminance function that converts a colour to the luminance intensity,  $Y(\mathbf{C}^t) - Y(\mathbf{B}^t)$  is the difference between the captured colour and rough diffuse albedo of voxels subject to  $Y(\mathbf{C}^t) > Y(\mathbf{B}^t)$ , yielding initial specular shading in the colour camera and  $S_C^t(\tilde{\mathbf{x}}_c^t) = S_C^t(\mathbf{0}; \hat{\rho}_{s,I}, \hat{\alpha}_I, \mathbf{n}, \tilde{\mathbf{x}}_c^t)$  is specular shading rendered at  $\mathbf{C}^t$  with the IR specular parameters using Equation (2). In order to calculate the diffuse shading image, we use the  $(t - 1)$  frame estimated diffuse albedo. Our algorithm refines the diffuse albedo and the specular albedo progressively over time.

**Coarse-to-Fine Propagation of Parameters.** Before we render the specular shading of each voxel  $S_C^t(\tilde{\mathbf{x}}_c^t)$ , we propagate the visible specular parameters from the deformation graphs to the resolution of TSDF. Each deformation-graph node takes the appearance values from its associated cluster. Each TSDF voxel obtains the parameters from the deformation nodes based on the  $k$ -nearest neighbours classified by diffuse albedo. As every voxel  $\mathbf{x}$  is associated with four  $k$ -nearest neighbour nodes, we propagate per-cluster  $\hat{\alpha}(\mathbf{x})$  and  $\hat{\rho}_s(\mathbf{x})$  to every voxel  $\mathbf{x}$  by the minimum difference of albedos in each voxel and the node within the  $k$ -nearest neighbour.



**Figure 6:** We accumulate all the photometric samples from the fine to the coarse levels: TSDF, deformation graph and cluster. After we estimate specular appearance per cluster, we propagate the estimated appearance from the coarse to the fine levels.

### 5.3.2. Diffuse Albedo Estimation

Existing fusion-based methods that estimate appearance account for diffuse reflection, assuming that surfaces have pure diffuse albedo only. The traditional fusion-based methods can integrate averaged photometric observations as diffuse albedos per voxel without separating specular reflection from them [NFS15]. The state-of-the-art method [GXY\*17] accounts for shading when calculating diffuse albedos by capturing the environment illumination additionally. However, these methods still cannot account for specular reflection from diffuse albedo computation. In contrast, our method separates specular reflection from the entire reflection, yielding pure diffuse reflection.

### 5.3.3. SVBRDF Optimization

Given the motion field  $\mathcal{W}^t$ , we estimate the surface properties of SVBRDF  $\mathcal{A}^t = \{[\rho_d, \rho_s, \alpha]^t\}$ : diffuse albedo, specular albedo and surface roughness per voxel  $\mathbf{x}$  in the TSDF volume  $\mathcal{T}$  by formulating the following energy function:

$$E_{\text{SVBRDF}}(\mathcal{A}^t) = E_{\text{vcolour}} + \lambda_{\text{treg}} E_{\text{treg}} + \lambda_{\text{sreg}} E_{\text{sreg}}, \quad (11)$$

where  $E_{\text{vcolour}}$  is the per-voxel colour data term,  $E_{\text{treg}}$  is the temporal regularizer and  $E_{\text{sreg}}$  is the spatial regularizer for the diffuse SVBRDF parameters.

The colour data term  $E_{\text{vcolour}}$  enforces photometric consistency of the SVBRDF parameters (on each voxel warped to the camera  $\mathbf{x}_c^t$ ) to make rendering with them satisfy given camera observation  $\mathbf{C}^t$ :

$$E_{\text{vcolour}} = \sum_{\mathbf{x} \in \mathcal{T}_C^t} \Phi(\|\tilde{\mathbf{n}}_c^t - \tilde{\mathbf{o}}_c^t\|) \|\mathbf{C}^t(\tilde{\mathbf{u}}_{\mathbf{x}_c}) - L^t(\tilde{\mathbf{x}}_c^t)\|_2^2, \quad (12)$$

where  $\tilde{\mathbf{u}}_{\mathbf{x}_c} = P(\mathbf{K}_C \tilde{\mathbf{x}}_c^t)$  is a corresponding pixel of  $\tilde{\mathbf{x}}_c^t$  at the current colour image  $\mathbf{C}^t$ ,  $\tilde{\mathbf{n}}_c^t$  and  $\tilde{\mathbf{o}}_c^t$  are normals and camera vectors at  $\tilde{\mathbf{x}}_c^t$ , respectively and  $\Phi$  is a robust kernel where  $\Phi(x) = 1/(1 + 5x)^3$ , following [ZDI\*15]. Here,  $L^t(\tilde{\mathbf{x}}_c^t) = L^t(\tilde{\mathbf{o}}_c^t; \tilde{\mathbf{n}}_c^t, \tilde{\mathbf{x}}_c^t)$  is the outgoing radiance under visible environment illumination, which is the sum of diffuse radiosity  $B^t$  and specular reflection  $S^t$  of the voxel in the colour camera space (Equation (2)).



Regularizer  $E_{\text{ureg}}$  in Equation (11) suppresses the temporal overfit of diffuse albedo  $\rho_d$  towards specular reflection:

$$E_{\text{ureg}} = \sum_{\mathbf{x} \in \mathcal{T}_C^t \cap \mathcal{T}_C^{t-1}} \|\rho_d^t(\mathbf{x}) - \rho_d^{t-1}(\mathbf{x})\|_2^2, \quad (13)$$

where  $\mathcal{T}_C^{t-1}$  is a set of visible surface voxels  $\mathbf{x}$  at the previous colour camera frame  $C^{t-1}$ .

In addition to the colour data term, we enforce local smoothness of diffuse albedo by formulating  $E_{\text{sreg}}$ :

$$E_{\text{sreg}} = \sum_{\mathbf{x} \in \mathcal{T}_C^t} \sum_{\mathbf{y} \in N(\mathbf{x}) \cap \mathcal{T}_C^t} \Phi(\|\Gamma(\tilde{\mathbf{u}}_{\mathbf{x}_C}) - \Gamma(\tilde{\mathbf{u}}_{\mathbf{y}_C})\|) \|\rho_d^t(\mathbf{x}) - \rho_d^t(\mathbf{y})\|_2^2, \quad (14)$$

where  $N(\mathbf{x})$  is a set of the neighbouring voxels  $\mathbf{x}$ ,  $\tilde{\mathbf{u}}_{\mathbf{x}_C}$  and  $\tilde{\mathbf{u}}_{\mathbf{y}_C}$  are pixels obtained by transforming voxels  $\mathbf{x}$  and  $\mathbf{y}$  to the current colour camera space  $C^t$ , respectively,  $\Gamma = C^t/Y(C^t)$  is the ratio of chromaticity to luminance  $Y$  of each pixel.

To implement this optimization progressively, we render visible specular reflection  $S^t$  with the specular parameters  $\hat{\rho}_s$  and  $\hat{\alpha}$  at voxel  $\mathbf{x}^t$  that we have estimated in Section 5.3.1, using  $\mathbf{i}_C$  and  $\mathbf{o}_C$  under visible environment illumination (captured by a 360 camera). We then subtract the estimated specular components from the captured image so that Equation (11) can be optimized only with respect to the pure diffuse albedo. This can be solved with the pre-conditioned conjugate gradient optimization as it becomes a least-square problem.

## 6. Implementation Details

**Radiometric Calibration.** We have conducted radiometric calibration for the RGB camera module and the infrared TOF camera module inside an RGBD device, Kinect 2 (Figure 3) in order to quantify the sensor responses in the red, green, blue and infrared channels. First, we estimate the RGB irradiance of the illumination ( $r_n, g_n, b_n$ ) by capturing the standard reflectance tile, Spectralon (Labsphere SRM99) for white balancing with the gamma value of  $\gamma_C = 2.2$ . Then, we calibrate the infrared camera parameters by solving the following optimization [PNS18]:

$$\min_{\kappa, \gamma} \sum_{u \in \mathcal{P}_S} \left( \mathbf{I}(u) - \left( \kappa \cdot \psi \frac{\mathbf{n}(u) \cdot \mathbf{i}(u)}{\pi \cdot d^2(u)} \right)^{\gamma_I} \right)^2, \quad (15)$$

where  $\mathcal{P}_S$  is a set of pixels  $u$  in the region where the spectralon is captured,  $\kappa$  is the illumination intensity of the infrared emitter in the Kinect 2 sensor,  $\gamma_I$  is the gamma exponent of the infrared camera,  $\mathbf{I}(u)$  is the infrared value at the pixel  $u$ ,  $\mathbf{n}(u)$  is the normal of the pixel  $u$ ,  $\mathbf{i}(u)$  is the incident light direction of the pixel  $u$  and  $d(u)$  is the distance between the IR emitter and the pixel  $u$ . We have estimated the values of  $\kappa$  and  $\gamma_I$  as 0.46 and 0.92 through nonlinear optimization [BGN00]. Given the radiometric parameters  $r_n, g_n, b_n, \kappa, \gamma_C$  and  $\gamma_I$  in the pre-processing of calibration, we linearize each RGB and infrared images and normalize them with irradiance.

**Pre-conditioned Conjugate Gradient for GPU.** We have implemented a GPU-based data-parallel pre-conditioned conjugate gradient (PCG) solver [WBS\*13]. The main computational bottleneck is the part of calculating matrix–vector multiplication. Fol-

**Table 2:** Per-frame processing time of our method. Our method takes 456 ms in total to process each frame inputs.

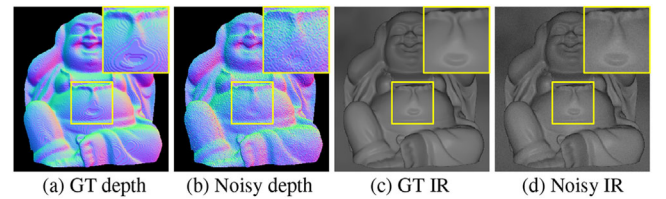
Algorithm	Processing time (ms)
Global registration	8
Motion estimation	224
TSDF integration	26
Specular estimation	89
Diffuse estimation	43
Marching cube	63
Etc.	2
Total	456

lowing [ZNI\*14], we have made use of two sparse matrix–vector multiplication kernels.

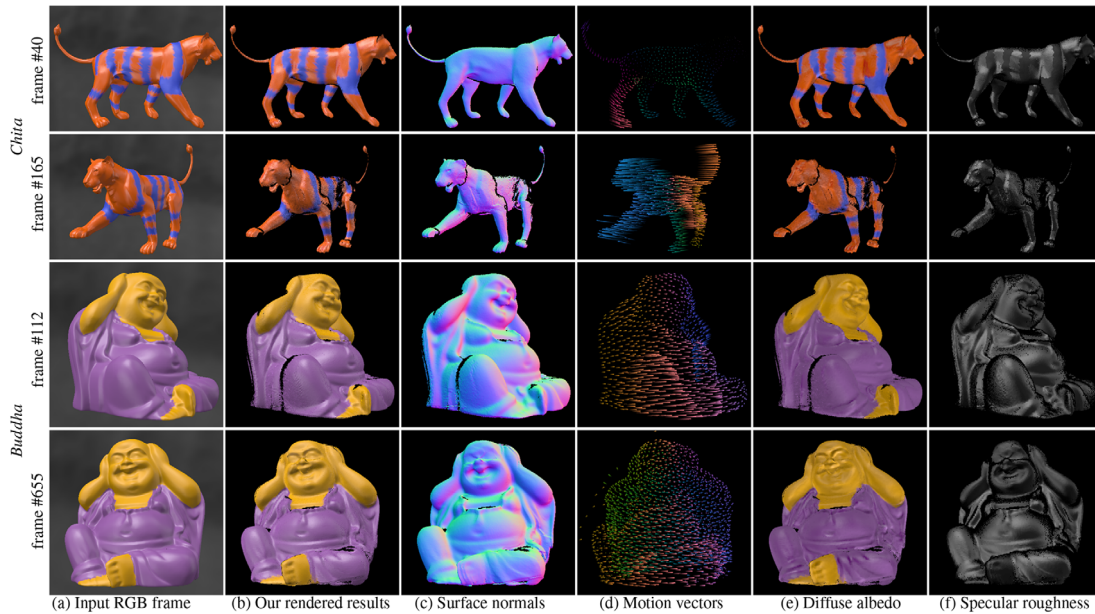
**Environment Map Capture.** To estimate incident illumination of scenes, we have captured scene environment maps as high-dynamic-range (HDR) radiance maps using a 360 camera (Ricoh Theta) with multiple exposures. In this paper, we have used monochromatic illumination maps by converting RGB radiance maps to luminance maps for computational efficiency. We then represent the environment maps with spherical harmonics coefficients for efficiently computing shading.

## 7. Results

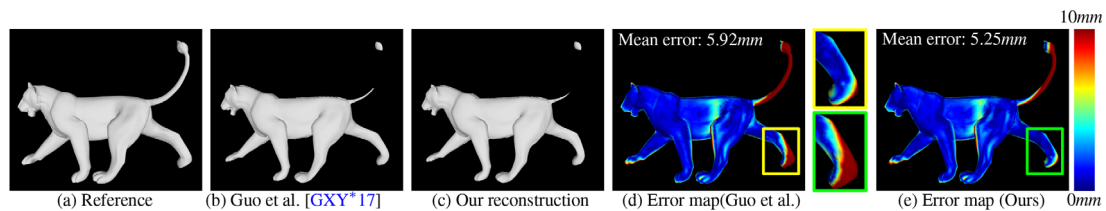
We built our capture setup using a Kinect 2 RGBD camera that consists of both an RGB and an infrared camera with an infrared illuminator in the TOF camera module (see Figures 1 and 3). Our method is implemented in C++, where CUDA-based GPU acceleration is extensively used for parallel processing, along with the OpenGL Shading Language for rendering. We set the resolution of the TSDF volume as  $512 \times 512 \times 512$ , and each TSDF voxel is defined as a cube with a width of 2 mm. Each node in the deformation graph has a radius of 20 mm. For the ground-truth data, we use 1.5 mm voxel size and 15 mm deformation graph radius. The truncation range for TSDF is five times wider than the voxel size. We pre-compute a discrete table of the BRDF function for pre-defined samples of parameters: The half-angle is sampled from 0 to  $60^\circ$  with a step size of  $1^\circ$ . Then, the Ward BRDF model is pre-computed with the values of  $\alpha$  and  $\rho_s$  from 0.05 to 0.70 and 0.01 to 1 both with 0.01 intervals, respectively. We tested our algorithm on a desktop computer with an Intel Core i7-7700K 4.20 GHz and a graphics card of



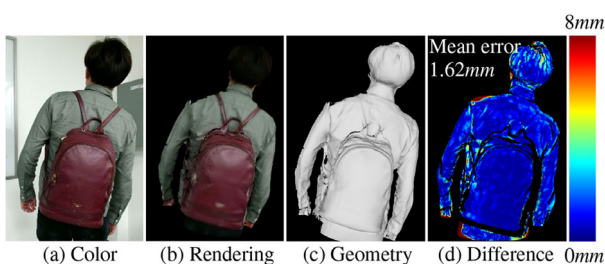
**Figure 7:** Synthetic input example. (a) Ground-truth depth image with normals. (b) Synthetic depth image with Gaussian noise. (c) GT IR image. (d) Synthetic IR image with Gaussian noise.



**Figure 8:** We have evaluated the reconstruction accuracy of geometry and appearance parameters, compared with the ground truth dataset synthetically created with Gaussian noise. (a) Input RGB images, (b) our results of reconstructed 3D models at each frame, followed by (c) our results of surface normals, (d) motion vectors, (e) diffuse albedo and (f) specular roughness.



**Figure 9:** We compared the geometric accuracy of our method with a state-of-the-art method, Guo et al. [GXY\*17], implemented by ourselves. We used the synthetic dataset that we created with Gaussian noise. Our method accurately reconstructs motion and geometry resulting in a low geometric error of average 5.25 mm. (a) Reference ground truth geometry. (b) Result by Guo et al. (c) Our reconstructed geometry. (d) Error map of Guo et al. compared with the GT geometry. (e) Error map of our results compared with the GT. Close-up boxes compares our method with that of Guo et al.

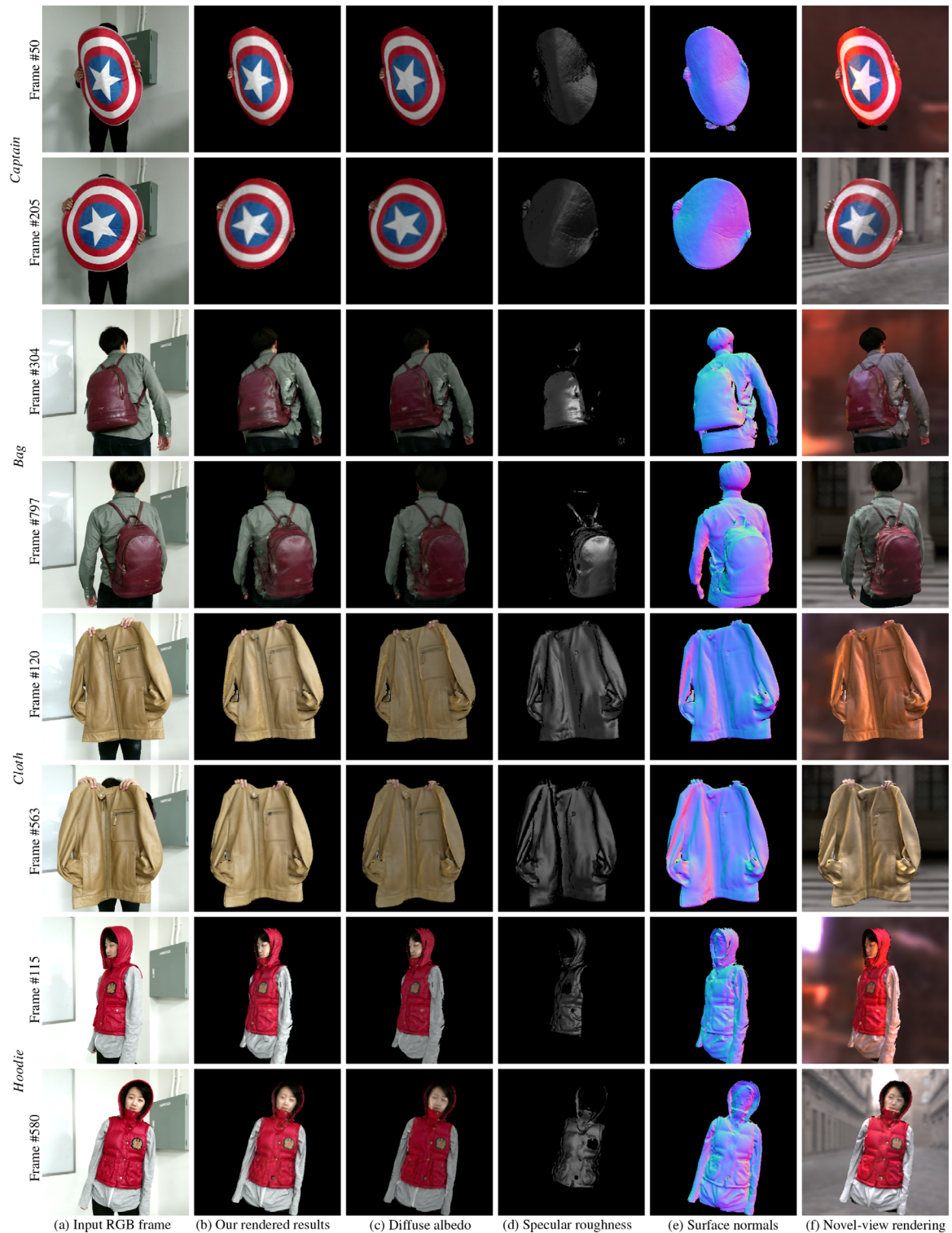


**Figure 10:** We evaluate the accuracy of the reconstructed geometry, comparing the warped geometry with the depth map of the current frame. (a) Input RGB image, (b) rendered result, (c) warped geometry, and (d) a difference map between (c) and the current depth map. The mean error of the depth values is just 1.62 mm.

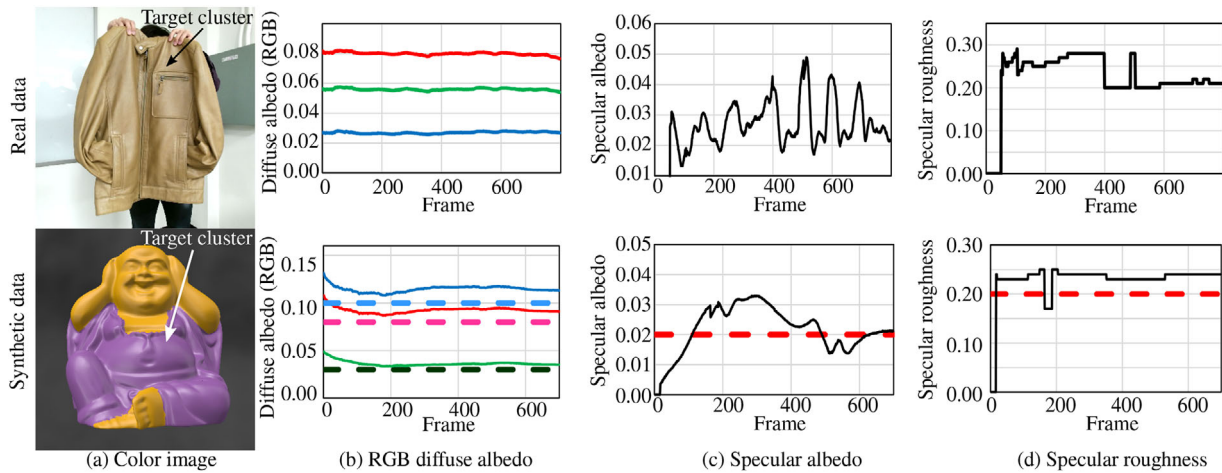
NVIDIA Titan V (12 GB). Our entire algorithm took 456 ms to process and render each frame. Our method is designed to be online, progressively processing input frames. Table 2 shows the detailed timestamps of our method taking 456 ms to process each frame. We provide our experimental results and comparison as follows.

**Quantitative Evaluation.** We created a synthetic dataset with known shape, SVBRDF and motion of different objects using OpenGL rendering. To make our synthetic dataset closer to the real sensor input, we also added Gaussian noise to the ground-truth (GT) depth images with  $\mathcal{N}(0, 0.002^2)$  and the GT IR images with  $\mathcal{N}(0, 1000^2)$ , as shown in Figure 7.

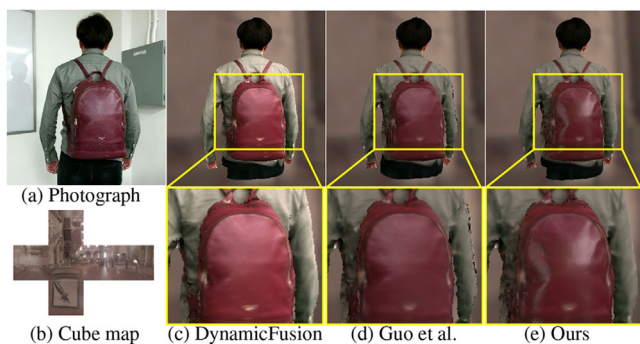
Using the synthetic dataset created with Gaussian noise, we have evaluated the accuracy of our reconstruction algorithm compared with the ground truth. Figure 8 presents our reconstruction results compared with the ground-truth SVBRDF, shape and motion. In order to quantitatively evaluate the reconstruction accuracy of shape



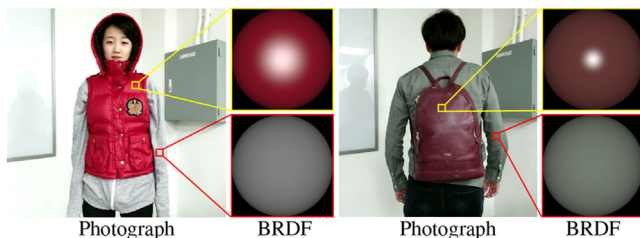
**Figure 11:** Our method faithfully reconstructs SVBRDF, geometry and motion of real-world dynamic objects. (a) Input RGB frame. (b) our results rendered with (c) reconstructed diffuse albedo, (d) specular roughness and (e) surface normals. (f) Results rendered under novel environment illumination. Note that we have multiplied the specular intensities by a factor of two for visualization purpose only. Refer to the supplemental video for more results.



**Figure 12:** Our method progressively estimates SVBRDF by accumulating observations. We evaluate the stability of our SVBRDF reconstruction for the real and synthetic dataset. (a) Real (top) and synthetic (bottom) colour photographs and target clusters noted by arrows. (b)–(d) Estimated SVBRDF parameter values for the target clusters over frames. We plot our estimates with solid lines and the ground truth values with dotted lines (synthetic only).



**Figure 13:** We compare our method with two online scanning methods. (a) Reference photograph. (b) An environment map. (c) DynamicFusion [NFS15] presents fixed specular reflection while (d) Guo et al. [GXY\*17] can reconstruct only diffuse shading. (e) Our method can acquire both specular and diffuse appearance. Note that specular reflection changes realistically in our results when the environment illumination spins. Refer to the supplemental video for more results.



**Figure 14:** We present captured BRDFs of dynamic objects, reconstructed by our method. Even though objects are dynamic, we reconstruct diffuse colour and specular lobe successfully.

and motion together, we warp the estimated geometry to the current camera frame with motion estimates for each frame. The averaged error between the ground-truth shape and the reconstructed shape with motion is very low at just 5.25 mm in the Hausdorff distance, as shown in Figure 9.

Moreover, we evaluate the accuracy of the estimated flow by our method in the real scene. We compare the differences between the actual depth map at #400 frame (of the *Bag* scene) and the warped geometry of our reconstructed model in Figure 10. The average distance error of the entire human body is just 1.62 mm. It is not surprising that there are some large errors around challenging geometry, such as hair and the silhouette of the body. Overall, our method successfully reconstructs SVBRDF, shape and motion for not only synthetic data but also real data. Refer to the supplementary video for every reconstruction results.

**Qualitative Evaluation.** Figure 11 presents the results of real-world dynamic objects. We present (a) input colour frame, (b) our reconstructed 3D objects rendered with a point light, (c) diffuse albedo, (d) specular roughness, (e) surface normals and (f) novel light-and-view rendering with an environment illumination map to qualitatively evaluate the overfitting problem of inverse rendering. The results of diffuse albedo and specular roughness demonstrate the effectiveness of our decomposition of material properties. Our results rendered under novel environmental lighting and view conditions present no typical blinking artefacts of overfitting. Refer to the supplemental video for every reconstruction result.

**Progressive Reconstruction.** Our method reconstructs SVBRDF and geometry progressively per frame. Figure 12 shows our estimates of specular albedo, specular roughness and diffuse albedo for each frame for a real scene and a synthetic scene. For the real-world case, the optimization for diffuse albedo converges fast thanks to our robust clustering. In contrast, optimization for specular

parameters requires a long iteration to be stabilized, showing fluctuation at an early stage. Our method provides fast and accurate convergence when optimizing both diffuse and specular appearance parameters: RGB diffuse albedo, specular albedo and specular roughness.

**Comparison.** Figure 13 compares our method with other two fusion-based methods of capturing objects using a single RGBD camera: DynamicFusion [NFS15] and Guo *et al.* [GXY\*17]. As there is no available public source, we implemented both methods. DynamicFusion [NFS15] does not separate diffuse and specular reflection, that is, it stores the sum of diffuse and specular colours as a single colour while geometry and motion can be faithfully recovered. Guo *et al.* [GXY\*17] extend the estimation of geometry and motion to capture the diffuse appearance of objects. Their method can cover only diffuse shading rendering, missing specular reflection. Our method is capable of estimating the full SVBRDF appearance of diffuse and specular components, in addition to the motion and geometry of dynamic objects. Figure 14 presents our reconstructed BRDFs of the dynamic objects that present very different characteristics.

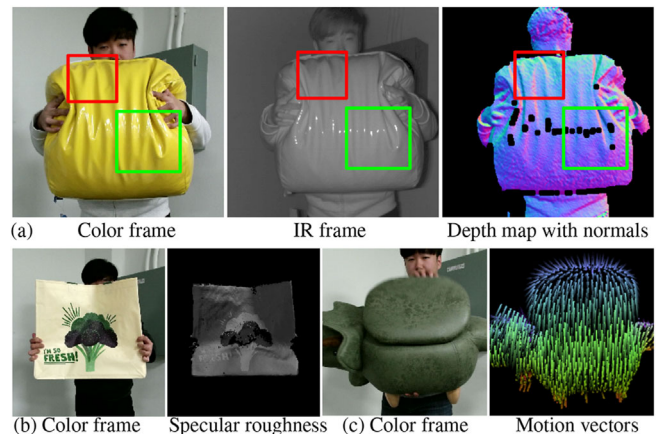
## 8. Discussion

Although our method can handle various scenes, it is not free from limitations.

**Infrared Illumination & Camera.** In order to estimate appearance parameters, we have made use of the infrared illumination and camera in the TOF module of Kinect 2. As these light and sensor modules are originally designed for measuring depth information, the RGBD camera API does not provide any custom control of exposure level of the IR emitter/sensor. Therefore, we were not able to capture high-dynamic-range (HDR) images as input for estimating specular reflection. When surface normals of specular objects look at the camera directly, not only depth map but also IR images have often been saturated, resulting in suboptimal results. See Figure 15(a).

**Frequency of Texture.** For estimating specular roughness parameters per pixel, we have used the infrared camera in the TOF module. The resolution of the IR image sensor is 512-by-424, while the resolution of the RGB image sensor is 1920-by-1080. The resolution of the IR sensor is four-time lower than that of the RGB sensor. Accordingly, when a target object has high-frequency patterns, the estimated specular parameter cannot reflect the object's appearance in high frequency. Figure 15(b) shows that the high-frequency structure of diffuse surface lines over the green tree cannot be estimated properly.

**Rapid motion.** When estimating motion vectors of dynamic objects, our method inherits the traditional linearized approximation of differential motion using the twist representation [MSZ94]. We found that the performance of this approach has become suboptimal when motion occurs very dynamically, or input frames are captured with severe motion blurring. See Figure 15(c) for an example. In addition, when motion causes deformation of objects surface with



**Figure 15:** (a) Owing to the limited dynamic range of the image sensor, input signals from the IR camera are saturated when surface smoothness is high. It often results in no depth values. The performance of our method becomes suboptimal when object surfaces are very smooth. (b) The resolution of the IR camera is four-time lower than that of the RGB camera, and thus our method fails in capturing high-frequency patterns of specular roughness. (c) When motion is large or rapid, our motion estimation often suffer from suboptimal reconstruction results due to motion blur.

texture, our method cannot account for the stretch of the texture surfaces. This would be interesting future work to explore.

**Spatial Resolution.** The spatial resolution of SVBRDF and shape is determined by the spatial resolution of the TSDF volume. As we currently store this information for each vertex, the current resolution degrades the spatial resolution of the final results. Applying texture mapping to our framework would be interesting future work.

**Frame Rate.** The current frame rate is about two frames per second, which is lower than the real-time performance thus far due to challenges of heavy optimization in factorizing SVBRDF, geometry and motion simultaneously. Accelerating computation for real-time applications would be an interesting avenue to explore.

## 9. Conclusions

We have presented a novel material acquisition method that estimates SVBRDF, geometry and motion simultaneously using a single RGBD camera. We have proposed an inverse rendering framework that can efficiently estimate material appearance using the voxel grid and the deformation graph in the two different scales. We have also provided the appearance-aware motion estimation algorithm so that the specular appearance can be considered to improve the motion estimation accurately. We have experimented with real-world objects. Finally, we have carefully discussed the limitations, evaluations and comparisons with other methods to validate the performance of our method.

## Acknowledgements

Min H. Kim acknowledges Samsung Research Funding Center of Samsung Electronics (SRFC-IT2001-04), in addition to a partial support of Korea NRF grants (2019R1A2C3007229, 2013M3A6A6073718), Samsung Research, the CT Research & Development Program of KOCCA in MCST of Korea, MSIT/IITP of Korea, MSRA and Cross-Ministry Giga KOREA (GK17P0200). The authors also thank Joo Ho Lee for helping GPU implementation and helpful comments.

## References

- [AWL15] AITTALA M., WEYRICH T., LEHTINEN J.: Two-shot SVBRDF capture for stationary materials. *ACM Transactions on Graphics (TOG)* 34, 4 (2015), 110:1–13.
- [BGN00] BYRD R. H., GILBERT J. C., NOCEDAL J.: A trust region method based on interior point techniques for nonlinear programming. *Mathematical Programming* 89, 1 (2000), 149–185.
- [BJTK18] BAEK S.-H., JEON D. S., TONG X., KIM M. H.: Simultaneous acquisition of polarimetric SVBRDF and normals. *ACM Transactions on Graphics* 37, 6 (2018), 268:1–15.
- [CL96] CURLESS B., LEVOY M.: A volumetric method for building complex models from range images. In *SIGGRAPH '96: Proceedings of the 23rd Annual Conference on Computer Graphics and Interactive* (New York, NY, 1996), ACM Press, pp. 303–312.
- [DCP\*14] DONG Y., CHEN G., PEERS P., ZHANG J., TONG X.: Appearance-from-motion: Recovering spatially varying surface reflectance under unknown lighting. *ACM Transactions on Graphics (TOG)* 33, 6 (2014), 1–12.
- [DDF\*17] DOU M., DAVIDSON P., FANELLO S. R., KHAMIS S., KOWDLE A., RHEMANN C., TANKOVICH V., IZADI S.: Motion2fusion: Real-time volumetric performance capture. *ACM Transactions on Graphics* 36, 6 (November 2017), 246:1–246:16.
- [DKD\*16] DOU M., KHAMIS S., DEGTAREV Y., DAVIDSON P., FANELLO S. R., KOWDLE A., ESCOLANO S. O., RHEMANN C., KIM D., TAYLOR J., KOHLI P., TANKOVICH V., IZADI S.: Fusion4d: Real-time performance capture of challenging scenes. *ACM Transactions on Graphics* 35, 4 (July 2016), 114:1–114:13.
- [DNZ\*17] DAI A., NIESSNER M., ZOLLHOFER M., IZADI S., THEOBALT C.: Bundlefusion: Real-time globally consistent 3d reconstruction using on-the-fly surface reintegration. *ACM Transactions on Graphics* 36, 4 (2017), 76a.
- [DRS10] DORSEY J., RUSHMEIER H., SILLION F.: *Digital Modeling of Material Appearance*. Amsterdam, Netherlands: Elsevier, 2010.
- [FHW\*11] FYFFE G., HAWKINS T., WATTS C., MA W.-C., DEBEVEC P.: Comprehensive facial performance capture. *Computer Graphics Forum* 30, (2011), 425–434.
- [GCHS10] GOLDMAN D. B., CURLESS B., HERTZMANN A., SEITZ S. M.: Shape and spatially-varying BRDFs from photometric stereo. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 32, 6 (2010), 1060–1071.
- [GHP\*08] GHOSH A., HAWKINS T., PEERS P., FREDERIKSEN S., DEBEVEC P.: Practical modeling and acquisition of layered facial reflectance. *ACM Transactions on Graphics* 27, 5 (2008), 139.
- [GXY\*17] GUO K., XU F., YU T., LIU X., DAI Q., LIU Y.: Real-time geometry, albedo, and motion reconstruction using a single RGB-D camera. *ACM Transactions on Graphics* 36, 3 (2017), 32.
- [HLZ10] HOLROYD M., LAWRENCE J., ZICKLER T.: A coaxial optical scanner for synchronous acquisition of 3d geometry and surface reflectance. *ACM Transactions on Graphics* 29, 4 (2010), 99.
- [HSL\*17] HUI Z., SUNKAVALLI K., LEE J.-Y., HADAP S., WANG J., SANKARANARAYANAN A. C.: Reflectance capture using univariate sampling of brdfs. In *Proceedings of the IEEE International Conference on Computer Vision* (Piscataway, NJ, 2017), IEEE, pp. 5362–5370.
- [IZN\*16] INNMANN M., ZOLLHOFER M., NIESSNER M., THEOBALT C., STAMMINGER M.: Volumedeform: Real-time volumetric non-rigid reconstruction. In *European Conference on Computer Vision* (Berlin, Heidelberg, 2016), Springer, pp. 362–379.
- [Kaj86] KAJIYA J. T.: The rendering equation. In *SIGGRAPH '86: Proceedings of the 13th Annual Conference on Computer Graphics and Interactive Techniques* (New York, NY, 1986), ACM Press, pp. 143–150.
- [KCvO07] KAVAN L., COLLINS S., ŽÁRA J., O’SULLIVAN C.: Skinning with dual quaternions. In *Proceedings of the 2007 Symposium on Interactive 3D Graphics and Games* (New York, NY, USA, 2007), ACM Press, pp. 39–46.
- [LPG19] LIN Y., PEERS P., GHOSH A.: On-site example-based material appearance acquisition. *Computer Graphics Forum* 38, (2019) 15–25.
- [LWS\*13] LI G., WU C., STOLL C., LIU Y., VARANASI K., DAI Q., THEOBALT C. (2013) Capturing relightable human performances under general uncontrolled illumination. *Computer Graphics Forum*, 32 (2pt3), 275–284.
- [LZG18] LI C., ZHAO Z., GUO X.: Articulatedfusion: Real-time reconstruction of motion, geometry and segmentation using a single depth camera. In *ECCV: European Conference on Computer Vision* (Berlin, Heidelberg, 2018), Springer.
- [MSZ94] MURRAY R. M., SASTRY S. S., ZEXIANG L.: *A Mathematical Introduction to Robotic Manipulation*. CRC Press, Boca Raton, FL, 1994.
- [NFS15] NEWCOMBE R. A., FOX D., SEITZ S. M.: Dynamicfusion: Reconstruction and tracking of non-rigid scenes in real-time. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (Boston, MA, 2015), IEEE, pp. 343–352.

- [NJR15] NIELSEN J. B., JENSEN H. W., RAMAMOORTHY R.: On optimal, minimal brdf sampling for reflectance acquisition. *ACM Transactions on Graphics* 34, 6 (2015), 1–11.
- [NLGK18] NAM G., LEE J. H., GUTIERREZ D., KIM M. H.: Practical SVBRDF acquisition of 3D objects with unstructured flash photography. *ACM Transactions on Graphics* 37, 6 (2018), 267:1–12.
- [NLW\*16] NAM G., LEE J. H., WU H., GUTIERREZ D., KIM M. H.: Simultaneous acquisition of microscale reflectance and normals. *ACM Transactions on Graphics* 35, 6 (2016), 185.
- [PNS18] PARK J. J., NEWCOMBE R., SEITZ S.: Surface light field fusion. In *2018 International Conference on 3D Vision (3DV)* (Verona, Italy, 2018), IEEE, pp. 12–21.
- [RH01] RAMAMOORTHY R., HANRAHAN P.: A signal-processing framework for inverse rendering. In *SIGGRAPH '01: Proceedings of the 28th Annual Conference on Computer Graphics and Interactive Techniques* (New York, NY, 2001), ACM Press, 117–128.
- [RL01] RUSINKIEWICZ S., LEVOY M.: Efficient variants of the ICP algorithm. In *3DIM* (Washington, DC, USA, 2001), IEEE Computer Society, pp. 145–152.
- [RPG16] RIVIERE J., PEERS P., GHOSH A.: Mobile surface reflectometry. *Computer Graphics Forum* 35, 1 (February 2016), 191–202.
- [RRFG17] RIVIERE J., RESHETOUSKI I., FILIPI L., GHOSH A.: Polarization imaging reflectometry in the wild. *ACM Transactions on Graphics* 36, 6 (2017), 206.
- [SBCI17] SLAVCHEVA M., BAUST M., CREMERS D., ILIC S.: Killing-fusion: Non-rigid 3d reconstruction without correspondences. In *IEEE Conference on Computer Vision and Pattern Recognition* (Honolulu, HI, 2017), IEEE, p. 7.
- [SBI18] SLAVCHEVA M., BAUST M., ILIC S.: Sobolevfusion: 3d reconstruction of scenes undergoing free non-rigid motion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (Salt Lake City, UT, 2018), pp. 2646–2655.
- [SSWK13] SCHWARTZ C., SARLETTE R., WEINMANN M., KLEIN R.: Dome ii: A parallelized BTF acquisition system. In *Proceedings of the Eurographics 2013 Workshop on Material Appearance Modeling: Issues and Acquisition* (Goslar, Germany, 2013), Eurographics Association, pp. 25–31.
- [SWK19] STOTKO P., WEINMANN M., KLEIN R.: Albedo estimation for real-time 3d reconstruction using RGB-D and IR data. *ISPRS Journal of Photogrammetry and Remote Sensing* 150 (2019), 213–225.
- [TAL\*07] THEOBALT C., AHMED N., LENSCH H., MAGNOR M., SEIDEL H.-P.: Seeing people in different light-joint shape, motion, and reflectance capture. *IEEE Transactions on Visualization and Computer Graphics* 13, 4 (2007), 663–674.
- [TFG\*13] TUNWATTANAPONG B., FYFFE G., GRAHAM P., BUSCH J., YU X., GHOSH A., DEBEVEC P.: Acquiring reflectance and shape from continuous spherical harmonic illumination. *ACM Transactions on Graphics (TOG)* 32, 4 (2013), 109:1–12.
- [War92] WARD G. J.: Measuring and modeling anisotropic reflection. *ACM SIGGRAPH Computer Graphics* 26, 2 (1992), 265–272.
- [WBS\*13] WEBER D., BENDER J., SCHNOES M., STORK A., FELLNER D.: Efficient gpu data structures and methods to solve sparse linear systems in dynamics applications. *Computer Graphics Forum* 32, 2 (2013), 16–26.
- [WVT12] WU C., VARANASI K., THEOBALT C.: Full body performance capture under uncontrolled and varying illumination: A shading-based approach. In *Proceedings of the 12th European Conference on Computer Vision - Volume Part IV* (Berlin, Heidelberg, 2012), Springer, pp. 757–770.
- [WWZ16] WU H., WANG Z., ZHOU K.: Simultaneous localization and appearance estimation with a consumer RGB-D camera. *IEEE Transactions on Visualization and Computer Graphics* 22, 8 (2016), 2012–2023.
- [WZ15] WU H., ZHOU K.: AppFusion: Interactive appearance acquisition using a kinect sensor. *Computer Graphics Forum* 34, 6 (2015), 289–298.
- [WZN\*14] WU C., ZOLLHOFER M., NIESSNER M., STAMMINGER M., IZADI S., THEOBALT C.: Real-time shading-based refinement for consumer depth cameras. *ACM Transactions on Graphics* 33, 6 (2014), 200.
- [XSH\*19] XU L., SU Z., HAN L., YU T., LIU Y., LU F. (2019) UnstructuredFusion: Realtime 4D Geometry and Texture Reconstruction using CommercialRGBD Cameras. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 1–1. <https://doi.org/10.1109/tpami.2019.2915229>.
- [YGX\*17] YU T., GUO K., XU F., DONG Y., SU Z., ZHAO J., LI J., DAI Q., LIU Y.: Bodyfusion: Real-time capture of human motion and surface geometry using a single depth camera. In *The IEEE International Conference on Computer Vision (ICCV)* (Venice, Italy, October 2017).
- [YZG\*18] YU T., ZHENG Z., GUO K., ZHAO J., DAI Q., LI H., PONS-MOLL G., LIU Y.: Doublefusion: Real-time capture of human performances with inner body shapes from a single depth sensor. In *The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)* (Salt Lake City, UT, June 2018).
- [YZZ\*19] YU T., ZHENG Z., ZHONG Y., ZHAO J., DAI Q., PONS-MOLL G., LIU Y.: SimulCap : Single-View Human Performance Capture With Cloth Simulation. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* pp. 5499–5509, (2019).
- [ZDI\*15] ZOLLHOFER M., DAI A., INNMANN M., WU C., STAMMINGER M., THEOBALT C., NIESSNER M.: Shading-based

refinement on volumetric signed distance functions. *ACM Transactions on Graphics* 34, 4 (2015), 96.

[ZNI\*14] ZOLLHOFER M., NIESSNER M., IZADI S., REHMANN C., ZACH C., FISHER M., WU C., FITZGIBBON A., LOOP C., THEOBALT C., et al.: Real-time non-rigid reconstruction using an RGB-D camera. *ACM Transactions on Graphics* 33, 4 (2014), 156.

[ZYL\*18] ZHENG Z., YU T., LI H., GUO K., DAI Q., FANG L., LIU Y.: Hybridfusion: Real-time performance capture using a single

depth sensor and sparse imus. In *European Conference on Computer Vision (ECCV)* (Munich, Germany, September 2018).

### Supporting Information

Additional supporting information may be found online in the Supporting Information section at the end of the article.

Data S1

Data Video S2