

# Multi-View Inpainting for RGB-D Sequence

Feiran Li, Gustavo Alfonso Garcia Ricardez, Jun Takamatsu and Tsukasa Ogasawara  
Nara Institute of Science and Technology  
8916-5 Takayama-cho, Ikoma, Nara, Japan  
{li.feiran.kw8, garcia-g, j-taka, ogasawar}@is.naist.jp

## Abstract

*In this work we propose a novel approach to remove undesired objects from RGB-D sequences captured with freely moving cameras, which enables static 3D reconstruction. Our method jointly uses existing information from multiple frames as well as generates new one via inpainting techniques. We use balanced rules to select source frames; local homography based image warping method for alignment and Markov random field (MRF) based approach for combining existing information. For the left holes, we employ exemplar based multi-view inpainting method to deal with the color image and coherently use it as guidance to complete the depth correspondence. Experiments show that our approach is qualified for removing the undesired objects and inpainting the holes.*

## 1. Introduction

3D reconstruction tools [1,2] such as structure from motion (SfM) and visual simultaneous localization and mapping (v-SLAM) serve for many purposes from path planning to scene understanding. Many of such approaches have provided RGB-D versions for the easily achievable distance information. However, with these classical methods undesired objects would be inevitably introduced to the constructed map if the input sequences are captured in dynamic environments. Some methods [3,4] present solutions to deal with the rigidly moving objects by clustering features into either static or dynamic classes. It was not until the recent blossom of image semantic segmentation that demonstrates new insight on how to remove the undesired objects in 2D image level with more flexibility.

However, simply removing the unwanted objects would leave blanks on the images and hence we focus on filling in these holes with inpainting techniques. Within the off-the-shelf approaches, the single image inpainting methods [5,6] are useless to handle large holes since the source region on a single image cannot provide enough information. As for the video completion algorithms [7, 8], they are specifically

designed to deal with color videos and hence not suitable for the RGB-D sequences in our case. We consider our method as one of the multi-view inpainting approaches, which are much more flexible than video completion for being able to handle both color and RGB-D images as well as to make use of information gotten from various kinds of sources, such as images captured with stereo cameras or the discrete ones.

In this paper we propose a unified framework to remove robust classes of clutters from RGB-D sequences. With a slight abuse of notation we inherently define a pair of color and depth images as a “frame”, which is the minimum unit in our algorithm. For a certain frame with masks to fill in, our algorithm searches useful information from the other counterparts. We consider our system as semi-automatic since one can conveniently remove the annoyances via defining their semantic classes.

We summarize the contributions of this paper as follows. First, we introduce a source selection strategy which carefully balances the internal trade-offs among distinct frames to attain expected results. Second, we use the local homography based warping method to achieve more accurate alignments than previous work [7,9]; as well as to prevent information loss during 2D-3D projection in the SfM based methods [10, 11]. Our third contribution is a series of inpainting methods to make use of information from multiple source frames, in which we propose an MRF based approach for combining the candidates; extend the searching region of exemplar based color image inpainting methods to multiple views and coherently inpaint the depth.

The remainder of this paper is organized as follows. Section 2 makes a brief introduction of related inpainting work. Section 3 gives an overview of our entire algorithm. Section 4 introduces the strategy to select source frames. Section 5 and 6 describe the specifics to achieve color and depth inpainting. Section 7 shows the tests and comparison results. Finally, the conclusions are given in section 8.

## 2. Related Work

Image inpainting contains broad research fields and hereby we briefly review those closely related to our work.

Our approach is a combination of multi-view image inpainting and depth inpainting.

## 2.1. Single Image Inpainting

Exemplar based methods are popular in single image inpainting for their texture-protection ability. The beginning of them can be traced to the work of Criminisi *et al.* [12], in which the mask is completed via searching for similar patches from the rest region and inherently copying them. The PatchMatch algorithm proposed by Barnes *et al.* [5] uses random search for quickly finding approximate nearest neighbor matches between patches, which is widely employed as the basis in the follow-up work for its several orders higher time efficiency. Kawai *et al.* [6] extend the energy function by taking into account of brightness changes and spatial locality of texture to deal with unnatural matches. Lee *et al.* [13] propose to take Laplacian pyramid as an error term in patch synthesis in order to protect edges. In summary, the single image inpainting approaches leverage information from the image itself. In contrast, we use the other frames as additional sources.

## 2.2. Video Completion

Video completion aims at dealing with color image sequences. Some work requires manual interaction: Klose *et al.* [8] propose to inpaint a given video by using SfM and manually drawing 3D masks. Other methods either completely copy information from other frames or generate new textures by searching from them: Granados *et al.* [14] enable free movement of camera via using multiple homographies to estimate the geometric registration between frames; whose applications are limited for being required to satisfy the assumption that the missing pixels on the target frame can be completely achieved from the others. On the other hand, Newson *et al.* [15] propose an exemplar based method to search for similar patches on a group of aligned source frames, which is pretty time consuming for minimizing a global energy function. Similarly, Ebdelli *et al.* [7] shrink the searching range by only considering a small number of aligned neighboring frames of the target one. Differently with these methods that handle color videos, our approach is designed for RGB-D sequences; Also we take advantages from both direct copying and multi-view searching and hence more suitable for highly textured scenes.

## 2.3. Multi-View based Inpainting

Multi-view inpainting techniques leverage information from multiple source frames. Hays and Efros [16] gather photos from Internet as a huge database to help with image completion. Similarly, Whyte *et al.* [9] cover an undesired region on the query image with Internet photographs of the same scene, in which multi-homography and photometric registration are used to achieve geometric registration be-

tween the query image and the source ones. Also a Markov random field optimization [17] is employed for selecting the optimal proposals. This kind of methods are pretty unstable since the masked objects are easy to be reintroduced as a result of lacking necessary means to filter the source information.

Recent research begins to show interests on using the geometric connections among different views. Baek *et al.* [18] present a multi-view based method to complete the user-defined region by jointly inpaint the color and depth image, which takes advantages from SfM to achieve geometric registration among different views. Similarly, Thonat *et al.* [10] enable free-viewpoint image based rendering with reprojected information from neighboring views. Also a refined method is proposed in the following work [11] that performs inpainting on intermediate, local planes in order to preserve perspective as well as to ensure multi-view coherence. In contrast, we use local homography for achieving pixel-wise correspondences, with which the information loss caused by SfM could be effectively avoided. Also they assume that the input images are of high quality, while in contrast, our approach aims at dealing with more common scenarios, such as those taken with moving cameras and hence cluttered with blurriness.

## 2.4. Depth Inpainting

Depth inpainting is similar to the propagation methods designed for the color images to a certain degree. Result quality may however be limited if the algorithms designed for color images are simply transplanted to the depth counterparts. Therefore popular solutions use color images as guidance to complete the holes on the depth ones. Miao *et al.* [19] introduce a texture assisted inpainting technique via dividing the target area into smooth and edge classes and distribute different partial differential equations (PDE) to each class. Atapour-Abarghouei *et al.* [20] perform semantic segmentation on the color images to get the object edges and the depth value is coherently propagated within every object. Such work targets on assigning value to each unknown pixel. In this work, however, we take the unknown as one of the existing values and only inpaint the mask left by the removed undesired objects.

## 3. Overview

The proposed pipeline can be found in Fig. 1. A masked RGB-D sequence  $F$  and one target frame  $F_t$  are taken as input. The RGB-D sequence serves as source frames from where inpainting information can be gotten and the masks stand for the objects that we would like remove, which are semi-automatically generated in our work with the help of deep learning based semantic segmentation techniques. Our goal is to fill in the masks on  $F_t$  with realistic content by using multi-view information. In order to achieve it, we care-

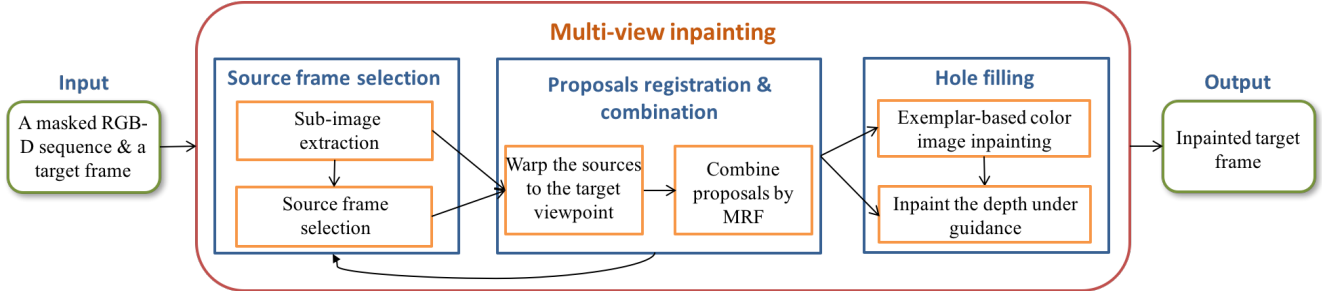


Figure 1: Overview of our method. For every input target frame, the first step is to select satisfactory source frames for inpainting. Then for every sub-images (specifically described in section 4.1), information from the corresponding sources is comprehensively combined. Finally the left holes on the color and depth images are inpainted separately.

fully select a set of source frames from  $F$  for  $F_t$ . Also we take benefits from local homography based image warping method [21] to warp each source  $F_i$  into the same image coordinate with  $F_t$ . Considering that the warpings are not equally accurate, we use a MRF approach similar to those proposed in [9, 17] to reduce bias. After this, we separately use exemplar based multi-view inpainting algorithm to cope with the color image  $I_t$  and coherently inpaint the correlated depth one under guidance.

## 4. Source Frame Selection

In a multi-view inpainting system, we first need to sort out a set of source frames from the input RGB-D sequence that can fill in the blanks in  $F_t$ . This is a quite challenging mission considering the giant quantity and variable qualities of the input frames. In this work, we comprehensively appraise the suitability of each frame to be used as source.

### 4.1. Strategy

Given a potential source frame  $F_i$ , the inpainting accuracy depends on both the image quality itself and the correlations between it and the target frame  $F_t$ . Therefore, we respectively grade the similarity, the inter-frame distance and the image quality of each  $F_i$  to evaluate whether it is suitable or not to inpaint  $F_t$ . It is reasonable to select these three factors because similarity ensures texture consistency; Larger distance would lead to lower warping accuracy by reducing the amount of matched feature points; as well as blurred image would not only weaken the quality of the inpainted image, but may also trigger mismatches among feature points, which would further lead to high-bias warping results.

A certain target frame  $F_t$  can contain several masks. Instead of inpainting them simultaneously, we emphasize local similarity and treat each mask  $M_t$  on  $F_t$  separately. Specifically, we split  $F_t$  into a set of sub-images by extracting the minimum circumscribed rectangle for every  $M_t$ , as shown in Fig. 2. Each  $M_t$  independently gets its own source

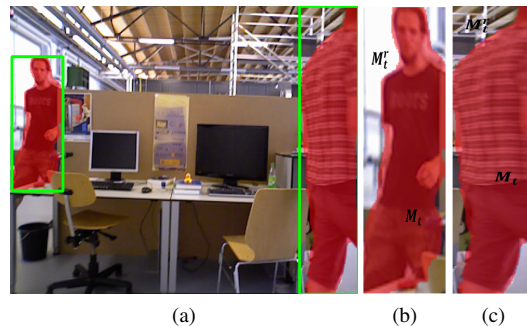


Figure 2: (a) The color image of a target frame. Red color indicates completion region and each green box is a minimum circumscribed rectangle of one mask. (b) and (c) Every box is one sub-image and in each one of them, the red region indicates  $M_t$  and the left part is  $M_t^r$ .

frames and has not influence on the selections of the others.

We evaluate the suitability  $S(F_i, M_t)$  for inpainting  $M_t$  with potential source frame  $F_i$  by Eq. 1

$$S(F_i, M_t) = \frac{(w_1 s(F_i, M_t) - w_2 d(F_i, M_t)) q(F_i)}{q(F_t)}, \quad (1)$$

where  $s(F_i, M_t)$  stands for the similarity between  $F_i$  and  $M_t$ ;  $d(F_i, M_t)$  represents the distance between  $F_i$  and  $M_t$ ; as well as  $q(F_i)$  and  $q(F_t)$  shows the image quality of  $F_i$  and  $F_t$  respectively.  $w_1$  and  $w_2$  are positive numbers for weight parameters. We use linear normalization to normalize all the factors to  $[0, 1]$  range.

In practicing, it is pretty burdensome to calculate the distance term for all the  $F_i$  in  $F$  and hence we employ a two-step selection strategy. Namely, we first select the most similar  $n$  source frames by only using the similarity term  $s(F_i, M_t)$  and then the  $S(F_i, M_t)$  is calculated within these  $n$  frames.

## 4.2. Factors

We consider the similarity evaluation as a typical image retrieval problem and in this work we use the bag of visual word (BoVW) model to solve it. In practice we use the SIFT feature [22] and vocabulary tree [23] for description and searching. The feature points are extracted from unmasked area of the sub-image  $M_i^r$  as shown in Fig. 2 (b) and Fig. 2 (c).

For the distance term, we take advantages from the extrinsic parameters of the camera. Since the distance relation between  $F_i$  and  $M_t$  is the same as that between  $F_i$  and  $F_t$ , we define the distance  $d(F_s, M_t)$  as follows

$$d(F_i, M_t) = \|r(F_i, F_t)\| + \|t(F_i, F_t)\|, \quad (2)$$

where  $r(F_i, F_t)$  and  $t(F_i, F_t)$  represent the rotation and translation distances respectively between the potential source frame  $F_i$  and the target one  $F_t$ . We set  $d(F_i, M_t) = +\infty$  for unsolvable conditions, which are discarded when doing normalization.

We use gradient based methods [24–26] for image quality assessment, which base their evaluation criteria on the proportion of relatively large gradients in all of them. A higher proportion indicates a larger amount of explicit edges and hence stands for higher image quality. This kind of methods is suitable for our case since we consider edges consistency as one of the factors in MRF (specifically described in section 5.2) and coherently use gradients in the energy function of it.

## 5. Proposals Registration and Combination

So far the source frames have been retrieved, the first thing to do is warping them into the same viewpoint with the target frame. For RGB-D sequence, the classical method to achieve pixel-wise correspondence is to use depth information for calculating the SE(3) transformation. However, a pixel with unknown depth value cannot be transformed and the information loss triggered by such invalid transformation would significantly depress the inpainting accuracy. Hence instead of it we propose to use local homography based warping method in this work.

Given the registered source frames, we need to decide from which one of them the mask pixel should get its value. A naive method is to project all the source frames into the target one and then blending them. However it is easy to trigger blurriness. In this work we propose to use MRF to make optimal choices among sources by considering it as a multi-label problem, similar to [9, 10, 17].

We also carry out post refinements for more natural and preciser results after solving the MRF. The entire framework of the proposed warping and combination approach is summarized in algorithm 1.

---

### Algorithm 1 Warp and combine multiple proposals

---

```

for Each Mask  $M_t$  in the target frame do
  for Each selected source frame  $F_i$  do
    Warp  $F_i$  into the same viewpoint with the target frame  $F_t$ ;
  end for
  for All the color images  $I$  of  $F_i$  do
    Calculate the average  $SSD$  between each  $I_i$  and  $I_t$ ;

    Create the median image by weighted overlaying each  $I_i$ ;
  end for
  Minimize the energy function of MRF by graph-cut;
  for Each source frame  $F_i$  do
    Compute the transformation matrix  $T_{ir}$  between  $F_i$  and  $F_t$ ;
    for Each source frame  $F_j (j \neq i)$  do
      Compute the transformation matrix  $T_{ij}$  between  $F_i$  and  $F_j$ ;
    end for
    end for
    Pose graph initialization;
    for Each vertexes  $V_i$  in the graph do
      Assign  $T_{ir}$  to  $V_i$ ;
      for Each vertexes  $V_j (j \neq i)$  in the graph do
        if  $T_{ij}$  then
          Assign  $T_{ij}$  to the edge  $E_{ij}$  that connects  $V_i$  and  $V_j$ ;
        else
          Disconnect the two vertexes;
        end if
      end for
    end for
    Do graph optimization;
  end for
  Do Poisson image editing;

```

---

### 5.1. Image warping

A single global homography is the simplest solution to describe pixel-wise correspondences between images. However it is restricted to subject to planar scene or pure rotation motion assumptions. Multi-homography methods [9, 27] have been employed in various kinds of research to deal with scenes that contain multiple planes, in which the images are divided into several planes and for each of them, an independent homography matrix is calculated. In this work we use the grid based local homography method for warping images taken with freely moving cameras. Whether a local homography can be calculated or not depends on the abundance of matched feature points, hence we use the affine SIFT [28] for points extraction.

## 5.2. Multiple Proposals Combination

Each source image (for convenience, without ambiguity we hereafter use source image to refer to the color image in the source frame) is considered as a label  $l$ . Our goal is to assign a label  $l_p$  for each pixel  $p$  in the mask area. We use the data cost term to represent the cost of assigning  $l_p$  to  $p$  and the smooth cost term to encourage avoiding explicit boundaries among distinct proposals. The energy function we would like to minimize it formulated in the form of

$$E(l) = \sum_{p \in \zeta} (\lambda_1 T_1(p, l_p) + \lambda_2 T_2(p, l_p)) + \sum_{(p, q) \in \zeta} \lambda_3 W(p, q, l_p, l_q), \quad (3)$$

where the sum of  $T_1(p, l_p)$  and  $T_2(p, l_p)$  indicate the data cost and  $\zeta$  represents all the pixels in the sub-image.  $W(p, q, l_p, l_q)$  stands for the smooth cost;  $(p, q)$  is a pair of neighboring pixels of which we use the 4-neighbor system and  $\zeta$  is the set of all such pairs in the sub-image region.  $\lambda_1$ ,  $\lambda_2$  and  $\lambda_3$  are weight parameters. It is important to notice that the unmasked pixels are also taken into account as contributions of the total energy since it can serve as constraint for the masked ones via the smooth cost term.

$T_1(p, l_p)$  is set to infinity if  $l_p(p)$  belongs to mask, representing waiting to be inpainted by in the hole filling section later on. For other cases, we choose  $T_1(p, l_p)$  to have the form

$$T_1(p, l_p) = \begin{cases} \|I_{l_p}(p) - I_t(p)\| & p \in M_t^r \\ \|I_{l_p}(p) - I_m(p)\| & p \in M_t \end{cases}, \quad (4)$$

where  $I_{l_p}(p)$  indicates the value of pixel  $p$  on the proposal  $l_p$ ;  $I_t(p)$  is that on the target image and  $I_m(p)$  represents the one on the median image, which is considered as a weighted combination of all source images with the rule

$$I_m(p) = \sum_{i \in N} w_i I_{l_i}(p), \quad (5)$$

where  $N$  is the set of source images;  $I_{l_i}(p)$  is the value of pixel  $p$  on the proposal  $l_i$  and  $w_i$  is weight. The weight is imported basing on the fact that the warping accuracies of all the proposals are not equal. Therefore we use the  $w_i$  for describing the relative accuracy and define it as

$$w_i = 1 - \frac{SSD(I_{l_i}(p), I_t(p))/n_i}{\sum_{j \in N} SSD(I_{l_j}(p), I_t(p))/n_j}, \quad (6)$$

where  $SSD$  is the sum of square difference between two pixel values; and  $n_i$  and  $n_j$  are the amounts of overlapped unmasked pixels between the source and target, which is used for normalization. Also we linearly normalize the RGB value to avoid influence caused by illumination changes.

As described earlier, the warping accuracy drops notably when distance is increased between viewpoints. Therefore we choose  $T_2(p, l_p)$  to have the form

$$T_2(p, l_p) = \exp(\|t(I_{l_i}, I_t)\|) - 1, \quad (7)$$

where  $t(I_{l_i}, I_t)$  is the translation distance between  $I_{l_i}$  and  $I_t$ ; different from Eq. 2, we ignore rotation distance since, as mentioned above, in principle, pure rotation would not introduce errors in homography computation. Also this function can provide approximately linear and distinguishable error term.

For the smooth cost term, we use the same gradient cost function proposed in [9], where  $W(p, q, l_p, l_q) = 0$  if  $l_p = l_q$  and otherwise

$$W(p, q, l_p, l_q) = \|\nabla I_{l_p}(p) - \nabla I_{l_q}(p)\| + \|\nabla I_{l_p}(q) - \nabla I_{l_q}(q)\|, \quad (8)$$

where  $\nabla I_{l_p}(p)$  indicates the gradient of  $I_{l_p}$  at pixel  $p$  and so on. We use the Sobel operator to get it.

Finally, to minimize the energy function described in Eq. 3, we use the graph-cut algorithm proposed in [29–31].

## 5.3. Color Adjustment

The brightness of the same scene may vary among distinct source images because of the different illumination conditions, which will cause obvious contrasts among the boundaries of sources. Therefore, in practice the mask is first expanded to a few more pixels with a dilation filter before solving the MRF and then the Poisson image editing [32] technique is used for blending the varied lighting.

## 5.4. Depth Transformation

Within previous procedures we have established pixel-wise correspondence between the target frame and the source ones. Now the task is to transform the depth from sources to the target. In principle we achieve it via following the classical "inverse projection, SE(3) transformation and projection" steps. For those source pixels whose depth values are unknown, we set them to blank to represent waiting to be inpainted in the hole filling section later on. Namely, for each source frame a transformation matrix  $T \in \mathbb{R}^{4 \times 4}$  should be calculated.

However, this is only part of the story. Another problem is that the estimation of pairwise transformation is not such accurate. Also considering that information for inpainting is collected from multiple proposals, it is easy to trigger inconsistencies on the inpainted depth image. Therefore, we implement a global optimization process to get the preciser transformation matrices by modeling it into a graph optimization problem. Specifically, we take the pose of the target frame as the origin and its camera coordinate as the world one. Then the vertexes in the graph can be set to the

transformation matrices  $T_{sr}$  from each source frame to the target. Coherently, the edges connecting pairwise vertexes are conditionally either set to the transformation matrices  $T_{ij}$  between the two source frames if the  $T_{ij}$  can be computed or left as blank to represent disconnections. We employ the g2o framework [33] for implementation.

## 6. Hole filling

Haven finished the multiple proposals combination, we can now proceed to the final step to cope with the still existing holes on the target frame. For the color image, a multi-view exemplar based inpainting method is proposed. Also, the inpainted color image will serve as guidance for completing the corresponding depth one.

### 6.1. Color Image Inpainting

Exemplar based inpainting methods basically synthesize values for the mask from the source by minimizing an energy function describing the similarity between them. In this work, we base our multi-view exemplar inpainting approach on the method proposed in [6, 34] and define the energy function in the form of

$$E = \sum_{p_i \in \phi, p_j \in \Phi} SSD(p_i, p_j), \quad (9)$$

where, the same as the general definitions in image inpainting work,  $\phi$  is the boundary area of the mask;  $\Phi$  is the source area from where the information is gotten;  $p_i$  and  $p_j$  respectively indicate pixels in  $\phi$  and  $\Phi$ ; and  $SSD$  represents the patch similarity in the form of

$$SSD(p_i, p_j) = \sum_{s \in \omega} \|I(p_{i+s}) - \alpha_{p_i p_j} I(p_{j+s})\| + \|\nabla I(p_{i+s}) - \nabla I(p_{j+s})\|, \quad (10)$$

where  $\omega$  is the patch size;  $s$  is a shift vector used to traverse all the pixels within  $\omega$ ;  $I(p_{i+s})$  and  $I(p_{j+s})$  indicate values of pixel  $p_{i+s}$  and  $p_{j+s}$ ;  $\nabla I(p_{i+s})$  and  $\nabla I(p_{j+s})$  are gradients used as constraints for preserving texture consistency.  $\alpha_{p_i p_j}$  is a parameter used for dealing with brightness changes [6], which is defined as

$$\alpha_{p_i p_j} = \sqrt{\frac{\sum_{s \in \omega} I^2(p_{i+s})}{\sum_{s \in \omega} I^2(p_{j+s})}}, \quad (11)$$

For multi-view inpainting, we expand the definition of  $\Phi$  as the unmasked area in the single image to that in both the target image and the source ones used for proposals combination. It is worth mentioning that the source images used in our approaches are the warped ones as described in section 5.2 to ensure texture consistency between the inpainted area and the original unmasked one. It is reasonable to do

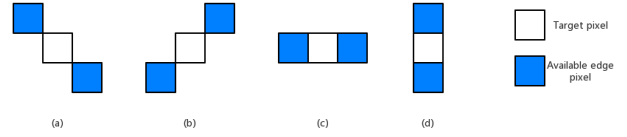


Figure 3: Creditable distribution samples for edge class

so since the warped images are in the same viewpoint with the target one, which means that we are seeking information from a spatially consistent region rather than several independent counterparts.

### 6.2. Guided Depth Image Inpainting

Depth images are usually low textured and propagation based methods are hence applicable. Similar to [19], we extract edges from the color image as guidance and coherently divide the masked pixels of the depth image into either smooth or edge classes. Specifically, a pixels would be classified into edge class if its color correspondence is edge and smooth otherwise. Propagations are processed separately on each class.

A given pixel can reckon its value from propagation only if its neighbors can provide enough information. In this work we use the 8-neighbor system and define respective rules for smooth and edge classes to evaluate the conditions of their neighbors. For the smooth class, a credible neighbor should has no less than 4 available pixels that also belong to the smooth; and for the edge one, a neighbor region is reliable if the distribution of available edge pixels within it satisfies one of the conditions shown in Fig. 3.

The mask area is iteratively inpainted until convergence. In each iteration, a pixel is either set to the propagated value provided satisfying the rules above or skipped otherwise. An edge pixel would be consider as mis-masked and move to the smooth class if it still cannot be assigned value after several loops. For propagation we use the Laplace equation in its discrete form

$$I(p, t + 1) = \frac{\sum_{p' \in \mu_8(p)} \kappa(p') I(p', t)}{\sum_{p' \in \mu_8(p)} \kappa(p')}, \quad (12)$$

where  $I(p, t + 1)$  is the value of pixel  $p$  at time  $t + 1$ ;  $\mu_8(p')$  indicates the 8-neighbor pixels of  $p$ ; and  $\kappa(p')$  is an indicator function which is set to 1 if the pixel is in the same class with  $p$  and 0 otherwise.

## 7. Results and Comparisons

We test our approach on different datasets from the TUM RGB-D benchmark [35] (freiburg3-walking-rpy, freiburg3-walking-xyz and freiburg3-walking-halosphere), in which



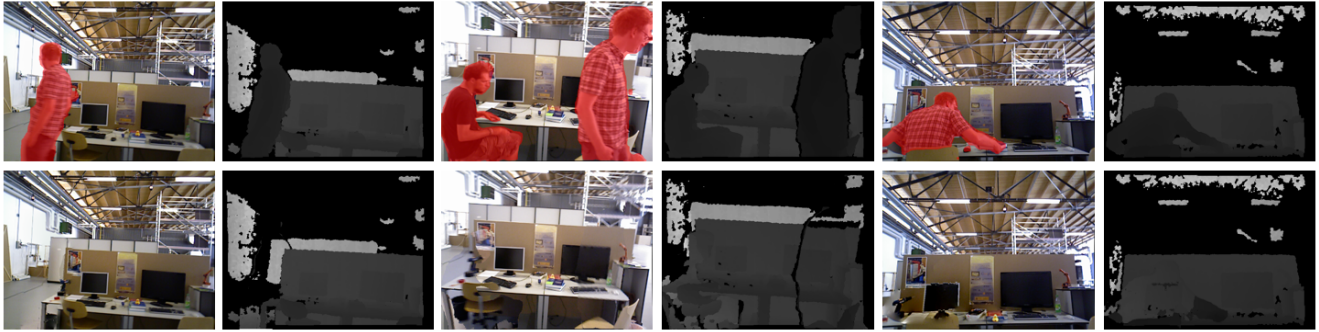


Figure 4: Testing results from different datasets. The first row indicates the original frames and the second one is the inpainted results with our method. The first two columns come from the freiburg3-walking-xyz dataset; middle two columns: freiburg3-walking-rpy; last two columns: freiburg3-walking-halosphere.

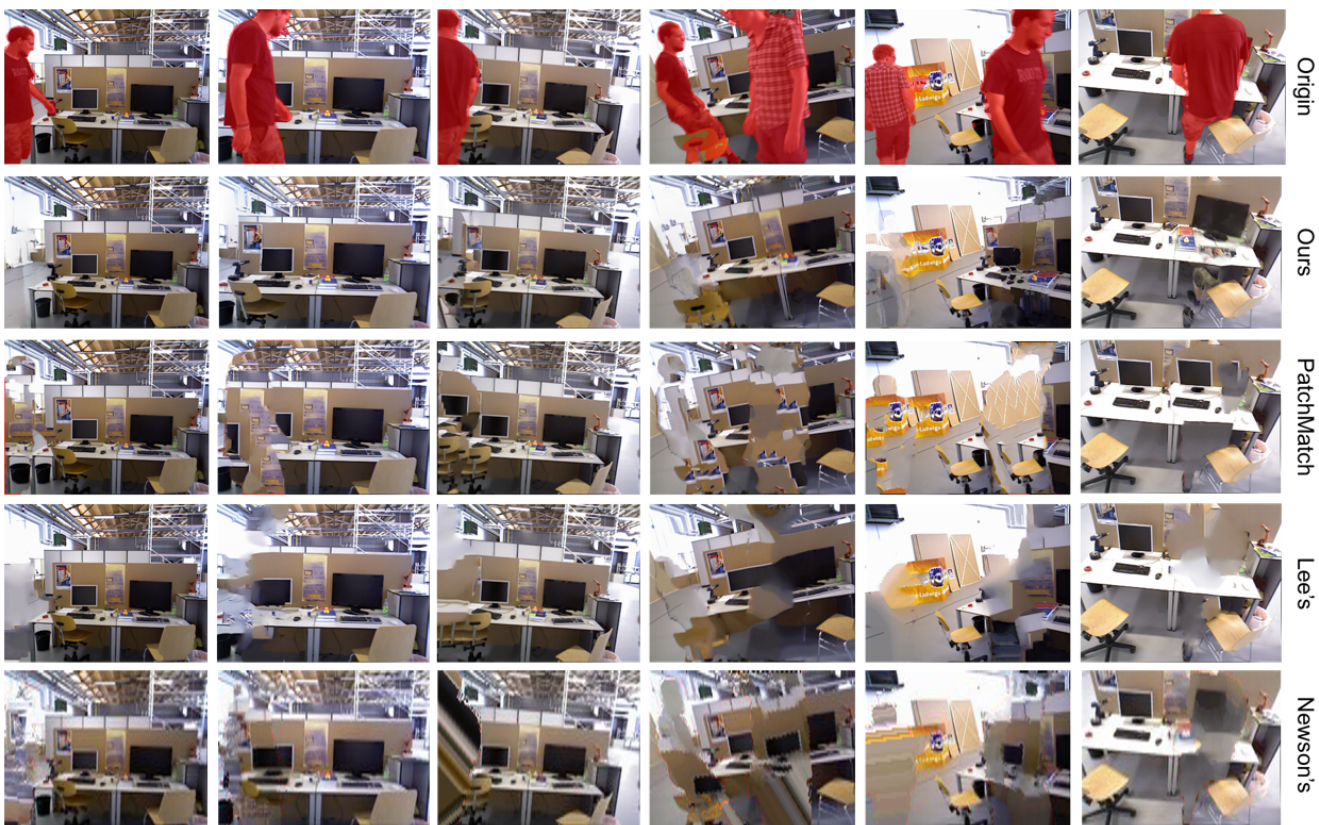


Figure 5: Compare with other methods. Each column shows one image inpainted by different methods. The first and second rows show the origin images and inpainting results with our method. The third and fourth rows present the results of single image inpainting methods respectively proposed by Barnes *et al.* [5] and Lee *et al.* [13]. The last row is the results of the video inpainting algorithm proposed by Newson *et al.* [15].

two persons randomly walk across the scenes and the cameras are also on moving. The "human" class is defined as undesired. It is important to notice that these three datasets are totally different because of the distinct camera motions although they present similar scenes. These datasets

are pretty challenging for being highly textured. For each dataset we use the first 800 frames and one target frame within it as input. The results are shown in Fig. 4.

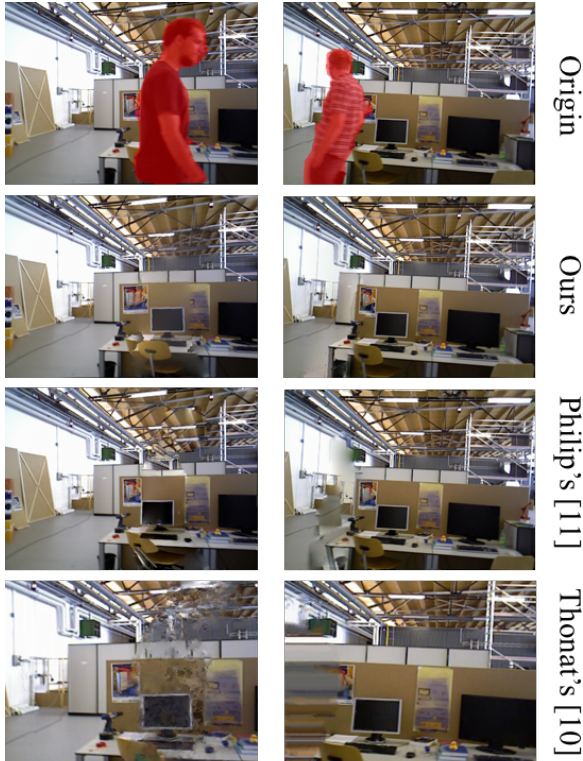


Figure 6: Compare with other multi-view inpainting approaches. Our method can effectively preserve texture consistencies (e.g. the upper bound of the white board and the ceiling).

### 7.1. Application Example

Hereby we present an application example of our work. For a set of selected frames, the fine-tuned PSPNet [36] is used to detect and draw masks on the "human" class. Also we expand the masks to a bit more pixels by dilation in order to cope with unmasked edges [37]. We use MeshLab [38] to present the final results as shown in Fig. 7.

### 7.2. Comparisons

As far as we know, similar work that semi-autonomously inpainting RGB-D sequences barely exists and therefore we could only compare our algorithm with other color image inpainting techniques. Specifically, we compare our method with other three approaches on all the three datasets, two of which are single image inpainting methods [5, 13] and the other one is proposed for video inpainting [15]. The results are presented in Fig. 5. Evidently the single image inpainting approaches cannot perfectly handle large blanks. As for the video inpainting method, we have to down-sample the images and use the neighboring 100 frames of the target as input because of its extremely high time cost. We also present comparisons among our methods and those pro-



Figure 7: Application example of our proposed algorithm. Above: point clouds projected from original images. Below: the counterpart after inpainting with our method.

posed in [10, 11], as shown in Fig. 6.

## 8. Conclusions and Discussion

We have introduced a multi-view based method for inpainting RGB-D sequences. Experiment shows the improvements of our method over the existing ones. However, like other homography based approaches, our method is easy to suffer from the sparseness of feature points when handling large baseline conditions. Also the segmentation accuracy would significantly infect the inpainting quality since the undesired objects might be re-introduced into the target frame provided poor segmentation.

For future work, other methods like those employing grid optimization can be explored to deal with large baseline conditions [39]. Also more suitable source selection strategy could be designed to avoid the current demand on weights adjustment. Another line of interest is, besides segment out the movable objects themselves (like the human), other objects which are passively moved (like the chair) can also be taken into account.

### Acknowledgement

We thank Norihiko Kawai for providing the source code of [34]; and the authors of [10, 11, 15] for their kind help on making comparisons.



## References

- [1] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohi, J. Shotton, S. Hodges, and A. Fitzgibbon, "Kinectfusion: Real-time dense surface mapping and tracking," in *Mixed and augmented reality (ISMAR), 2011 10th IEEE international symposium on*, pp. 127–136, IEEE, 2011.
- [2] C. Kerl, J. Sturm, and D. Cremers, "Dense visual slam for rgb-d cameras," in *Intelligent Robots and Systems (IROS), 2013 IEEE/RSJ International Conference on*, pp. 2100–2106, IEEE, 2013.
- [3] S. Caccamo, E. Ataer-Cansizoglu, and Y. Taguchi, "Joint 3d reconstruction of a static scene and moving objects," *arXiv preprint arXiv:1802.04738*, 2018.
- [4] E. Ataer-Cansizoglu and Y. Taguchi, "Object detection and tracking in rgb-d slam via hierarchical feature grouping," in *Intelligent Robots and Systems (IROS), 2016 IEEE/RSJ International Conference on*, pp. 4164–4171, IEEE, 2016.
- [5] C. Barnes, E. Shechtman, A. Finkelstein, and D. B. Goldman, "Patchmatch: A randomized correspondence algorithm for structural image editing," *ACM Transactions on Graphics-TOG*, vol. 28, no. 3, p. 24, 2009.
- [6] N. Kawai, T. Sato, and N. Yokoya, "Image inpainting considering brightness change and spatial locality of textures and its evaluation," in *Pacific-Rim Symposium on Image and Video Technology*, pp. 271–282, Springer, 2009.
- [7] M. Ebdelli, O. Le Meur, and C. Guillemot, "Video inpainting with short-term windows: application to object removal and error concealment," *IEEE Transactions on Image Processing*, vol. 24, no. 10, pp. 3034–3047, 2015.
- [8] F. Klose, O. Wang, J.-C. Bazin, M. Magnor, and A. Sorkine-Hornung, "Sampling based scene-space video processing," *ACM Transactions on Graphics (TOG)*, vol. 34, no. 4, p. 67, 2015.
- [9] O. Whyte, J. Sivic, and A. Zisserman, "Get out of my picture! internet-based inpainting," in *BMVC*, vol. 2, p. 5, 2009.
- [10] T. Thonat, E. Shechtman, S. Paris, and G. Drettakis, "Multi-view inpainting for image-based scene editing and rendering," in *3D Vision (3DV), 2016 Fourth International Conference on*, pp. 351–359, IEEE, 2016.
- [11] J. Philip and G. Drettakis, "Plane-based multi-view inpainting for image-based rendering in large scenes," in *13D 2018-ACM SIGGRAPH Symposium on Interactive 3D Graphics and Games*, pp. 1–11, 2018.
- [12] A. Criminisi, P. Pérez, and K. Toyama, "Region filling and object removal by exemplar-based image inpainting," *IEEE Transactions on image processing*, vol. 13, no. 9, pp. 1200–1212, 2004.
- [13] J. H. Lee, I. Choi, and M. H. Kim, "Laplacian patch-based image synthesis.," in *CVPR*, pp. 2727–2735, 2016.
- [14] M. Granados, K. I. Kim, J. Tompkin, J. Kautz, and C. Theobalt, "Background inpainting for videos with dynamic objects and a free-moving camera," in *European Conference on Computer Vision*, pp. 682–695, Springer, 2012.
- [15] A. Newson, A. Almansa, M. Fradet, Y. Gousseau, and P. Pérez, "Video inpainting of complex scenes," *SIAM Journal on Imaging Sciences*, vol. 7, no. 4, pp. 1993–2019, 2014.
- [16] J. Hays and A. A. Efros, "Scene completion using millions of photographs," in *ACM Transactions on Graphics (TOG)*, vol. 26, p. 4, ACM, 2007.
- [17] A. Agarwala, M. Dontcheva, M. Agrawala, S. Drucker, A. Colburn, B. Curless, D. Salesin, and M. Cohen, "Interactive digital photomontage," *ACM Transactions on Graphics (ToG)*, vol. 23, no. 3, pp. 294–302, 2004.
- [18] S.-H. Baek, I. Choi, and M. H. Kim, "Multiview image completion with space structure propagation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 488–496, 2016.
- [19] D. Miao, J. Fu, Y. Lu, S. Li, and C. W. Chen, "Texture-assisted kinect depth inpainting," in *Circuits and Systems (IS-CAS), 2012 IEEE International Symposium on*, pp. 604–607, IEEE, 2012.
- [20] A. Atapour-Abarghouei and T. P. Breckon, "Depthcomp: real-time depth image completion based on prior semantic scene segmentation.," 2017.
- [21] J. Zaragoza, T.-J. Chin, M. S. Brown, and D. Suter, "As-projective-as-possible image stitching with moving dlt," in *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, pp. 2339–2346, IEEE, 2013.
- [22] P. C. Ng and S. Henikoff, "Sift: Predicting amino acid changes that affect protein function," *Nucleic acids research*, vol. 31, no. 13, pp. 3812–3814, 2003.
- [23] D. Gálvez-López and J. D. Tardós, "Bags of binary words for fast place recognition in image sequences," *IEEE Transactions on Robotics*, vol. 28, pp. 1188–1197, October 2012.
- [24] A. Mittal, A. K. Moorthy, and A. C. Bovik, "No-reference image quality assessment in the spatial domain," *IEEE Transactions on Image Processing*, vol. 21, no. 12, pp. 4695–4708, 2012.
- [25] K. Bahrami and A. C. Kot, "A fast approach for no-reference image sharpness assessment based on maximum local variation," *IEEE Signal Processing Letters*, vol. 21, no. 6, pp. 751–755, 2014.
- [26] L. Ying, Z. Li, and C. Zhang, "No-reference sharpness assessment with fusion of gradient information hvs filter (in chinese)," *Journal of Image and Graphics*, vol. 20, no. 11, pp. 1446–1452, 2015.
- [27] H. Liu, G. Zhang, and H. Bao, "Robust keyframe-based monocular slam for augmented reality," in *Mixed and Augmented Reality (ISMAR), 2016 IEEE International Symposium on*, pp. 1–10, IEEE, 2016.
- [28] J.-M. Morel and G. Yu, "Asift: A new framework for fully affine invariant image comparison," *SIAM journal on imaging sciences*, vol. 2, no. 2, pp. 438–469, 2009.
- [29] Y. Boykov, O. Veksler, and R. Zabih, "Fast approximate energy minimization via graph cuts," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 23, no. 11, pp. 1222–1239, 2001.

- [30] V. Kolmogorov and R. Zabini, "What energy functions can be minimized via graph cuts?," *IEEE transactions on pattern analysis and machine intelligence*, vol. 26, no. 2, pp. 147–159, 2004.
- [31] Y. Boykov and V. Kolmogorov, "An experimental comparison of min-cut/max-flow algorithms for energy minimization in vision," *IEEE transactions on pattern analysis and machine intelligence*, vol. 26, no. 9, pp. 1124–1137, 2004.
- [32] P. Pérez, M. Gangnet, and A. Blake, "Poisson image editing," *ACM Transactions on graphics (TOG)*, vol. 22, no. 3, pp. 313–318, 2003.
- [33] R. Kümmerle, G. Grisetti, H. Strasdat, K. Konolige, and W. Burgard, "g2o: A general framework for graph optimization," in *Robotics and Automation (ICRA), 2011 IEEE International Conference on*, pp. 3607–3613, IEEE, 2011.
- [34] N. Kawai and N. Yokoya, "Image inpainting considering symmetric patterns," in *Pattern Recognition (ICPR), 2012 21st International Conference on*, pp. 2744–2747, IEEE, 2012.
- [35] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers, "A benchmark for the evaluation of rgb-d slam systems," in *Proc. of the International Conference on Intelligent Robot Systems (IROS)*, Oct. 2012.
- [36] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pp. 2881–2890, 2017.
- [37] L. Feiran, D. Ming, T. Jun, and O. Tsukasa, "Static 3d map reconstruction based on image semantic segmentation," in *The 15th International Conference on Ubiquitous Robots*, 2018.
- [38] P. Cignoni, M. Callieri, M. Corsini, M. Dellepiane, F. Ganovelli, and G. Ranzuglia, "Meshlab: an open-source mesh processing tool.," in *Eurographics Italian Chapter Conference*, vol. 2008, pp. 129–136, 2008.
- [39] C.-C. Lin, S. U. Pankanti, K. Natesan Ramamurthy, and A. Y. Aravkin, "Adaptive as-natural-as-possible image stitching," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1155–1163, 2015.