# REAL-TIME INDOOR SCENE RECONSTRUCTION WITH RGBD AND INERTIA INPUT

*Zunjie Zhu*

Hangzhou Dianzi University

*Feng Xu*

Tsinghua University

## ABSTRACT

Camera motion estimation is a key technique for 3D scene reconstruction and Simultaneous localization and mapping (SLAM). To make it be feasibly achieved, previous works usually assume slow camera motions, which limits its usage in many real cases. We propose an end-to-end 3D reconstruction system which combines color, depth and inertial measurements to achieve robust reconstruction with fast sensor motions. Our framework extends Kalman filter to fuse the three kinds of information and involve an iterative method to jointly optimize feature correspondences, camera poses and scene geometry. We also propose a novel geometry-aware patch deformation technique to adapt the feature appearance in image domain, leading to a more accurate feature matching under fast camera motions. Experiments show that our patch deformation method improves the accuracy of feature tracking, and our 3D reconstruction outperforms the state-of-the-art solutions under fast camera motions.

## 1. INTRODUCTION

With the rapid development of capture and computation devices, such as depth sensors and GPUs, real-time 3D reconstruction has made big growth. In recent years, a lot of works have focused on indoor scene reconstruction. For example, InfiniTAM[1] only uses depth information to reconstruct 3D models, and estimate camera poses by an iterative closest point(ICP) algorithm[2]. However depth only is extremely brittle in situations such as geometry-less scenes, bright windows and depth sensor noises, and can not eliminates accumulated errors. Drift-free camera tracking have been made breakthrough progress by monocular RGB-based methods, including direct methods[3] and feature point methods[4]. However these approaches can not reconstruct detailed and accurate 3D models. Further more, BundleFusion[5] and ElasticFusion[6] use both color and depth information to estimate camera motions and generate 3D models based on implicit truncated signed distance fields(TSDFs) and surfel representation, respectively.

Although these works exhibit reasonable results[7], they still require strong assumptions, like static scene without dynamic object, sufficient texture and geometrical information, slow camera motions and invariant illumination. However these assumptions can not be satisfied in many applications.
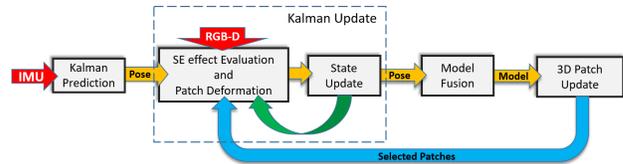


**Fig. 1**. Overview of our pipeline. The red, green and blue arrows represent the input acquired from current frame, iterative operation, and the patches from last frame.

In this paper, we make a step further by handling fast camera motions. For both color and depth, fast camera motion leads to large inter-frame distance, which makes it difficult to perform image feature matching (As images may be blurred and feature appearances may vary a lot.) and ICP based depth aligning. Similar with Tristan et al.[8], we solve the issue by introducing IMU information, gathered by an accelerometer and a gyroscope. Further combining with color and depth information, robust camera pose estimation and geometry fusion of an indoor scene are jointly achieved. The main contributions of our work are as follows:

(1) A *RGB-D-inertial 3D reconstruction system based on extended Kalman filter framework*, which tightly combines the three kinds of information, and jointly achieves camera pose estimation and patch deformation in the kalman update step.

(2) A *geometry-aware feature tracking method* for handling fast camera motion, which utilizes patch features to adapt blurry images and considers the deformation of patches in building feature matchings for images with very different perspectives.

(3) Through experiments on public datasets (including both synthetic and real data) and our data acquired by Intel Realsense ZR300, we see that our approach outperforms the state-of-the-art reconstruction systems under fast camera motions.

## 2. METHOD

The pipline of our system is illustrated in Fig. 1. We introduce our method by following four parts, geometry-aware feature tracking which explores the SE effect and executes patch deformation, filter framework which explain the kalman predic-
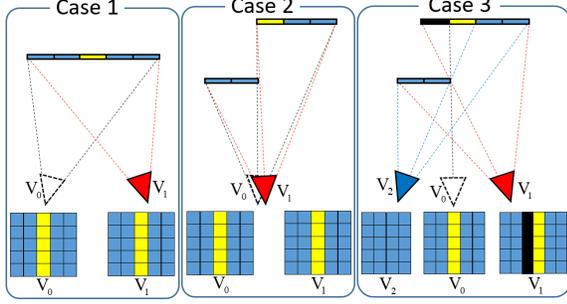
**Fig. 3**. Deformation process of SE patches.

**Fig. 2**. This figure shows the patch SE effects caused by the camera motion and the geometrical shape of patches.

tion and update step, model fusion and patch update.

## 2.1. Geometry-aware Feature Tracking

Point-based feature tracking methods[4] will extract insufficient number of features when images are blurred or with less texture. Thus, patch-based method, which considers larger image regions, is proposed[3] to track features under these situations. However, large patches may contain objects on different depth levels, which causes appearance changes in consecutive frames, especially when the camera motion is fast, leading to inaccurate feature tracking. To address this problem, we combine color and depth information to back project 2D patches into 3D and re-project them to the camera of the next frame by using an initial camera motion. The projection helps us to deform the original patches to model the appearance changes, and the patch tracking can be easily and accurately achieved by the deformed patches.

### 2.1.1. SE effect and Patch Deformation

When camera moves, a feature patch will be seen from different perspectives in different frames, and thus the shape and position of the feature patch in image coordinates vary from different frames. To account for the patch deformation, different from Bloesch et al.[9] which only considers the 2D planar information of a patch, we use the 3D geometry of a patch to determine the 2D shape deformation between images recorded with fast camera motion.

Depending on different geometries of patches and irregularities of camera motions, patches may produce different deformations in consecutive frames. Fig. 2 shows three representative cases of patch deformations:

**Case 1.** If there is no significant difference among pixel depths in a patch, then no matter how aggressive the camera motion is, the general shape of the patch will remain unchanged in two consecutive frames.

**Case 2.** When camera moves slowly, the 2D shape of a patch will still remain unchanged even though there are large depth variances existing in the patch.
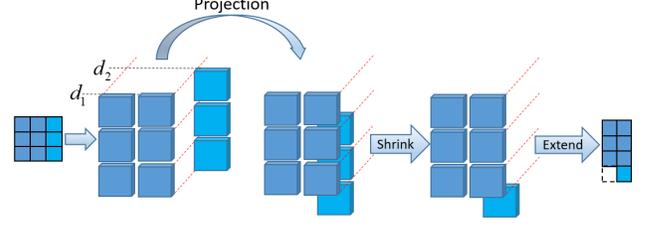
**Case 3.** Different from case 2, if camera moves aggressively, the intensity distribution and shape of the patch will changes. If camera moves from $V_0$ to the viewpoint $V_2$, the yellow region will be occluded and the patch shrinks in the current frame. In addition, the black region which is occluded in $V_0$ is visible once the camera moves to $V_1$, and thus the patch shape extends in the current frame. These phenomenons are called shrink effect and extend effect, so we call them together by SE effect.

We designed a unified deformation method to handle the SE effect. The details of patch deformation process are illustrated in Fig. 3 and formalized as follows:

Each pixel in a patch extracted from last frame $k-1$ are defined as $P_i^{k-1} := (\mathbf{p}_i^{k-1}, I_i^{k-1}, d_i^{k-1}, \mathbf{n}_i^{k-1})$, where $\mathbf{p}_i$ denotes the image coordinates of pixel $i$ in the patch. $I_i, d_i, \mathbf{n}_i$ denote the intensity, depth and 3D normal of pixel $i$. $d_i$ and $\mathbf{n}_i$ are obtained from the depth image, and since we have 3D information, we call our patches 3D patch features. We first back project each patch into 3D world coordinate system:

$$\mathbf{L}_i = \mathbf{T}_{k-1}\pi^{-1}(P_i^{k-1}). \tag{1}$$

Here, $\pi(\cdot)$ is to project a point from the 3D camera coordinate system to 2D pixel domain, and $\pi^{-1}(\cdot)$ represent the inverse operation. $\mathbf{T}_{k-1}$ is the camera pose of frame $k-1$ which transforms a 3D point in the camera coordinate system of frame $k-1$ to the world coordinate system. Thus $\mathbf{L}_i$ is in the world coordinate systme but is indexed from a pixel $i$ in frame $k-1$, so $\mathbf{L}_i := (\mathbf{P}_i, I_i^{k-1}, \mathbf{n}_i)$ where the intensity information does not change in the projection and $\mathbf{P}_i$ indicates the 3D coordinate. Then we project it to the pixel coordinate system of the current frame $k$.

$$P_i^{k'} = \pi(\mathbf{T}_k^{-1}\mathbf{L}_i). \tag{2}$$

Here $'$ means it is projected to this frame, not original from this frame. Note that $\mathbf{T}_k$ is unknown and affects the 2D position of the projection. After projection, if two projected pixels happened to be in the same pixel coordinate, then we consider the shrink effect occurred in this patch. This usually happens when a region close to the camera covers the distant region. Therefore we remove the pixels corresponding to the distant region. Then we evaluate whether the extend effect happens or not. We set the shape of the projected patch to be the bounding box of all projected pixels. Then the extend

effect will be verified if the height (or width) of the projected patch is greater than the original one. No matter which effect happens, we use $P^{k'}$ as the deformed patch in the following feature tracking steps.

### 2.1.2. Objective

In feature tracking, patches affected by SE effect are replaced with the corresponding projected patches. Then We track the projected patch features by both intensity and depth information.

Photometric error for each projected patch is computed as follows. We first extract the patch of the same size at the projected location of the current image as $P^k$, and then calculate the intensity difference between the extracted patch and the projected patch. The photometric error can be formalized as:

$$E_p = \sum_i^Y \left\| I[P_i^k] - I[P_i^{k'}] \right\|_2, \qquad (3)$$

where $Y$ denotes the number of pixels in the patch. $I$ indicates the corresponding intensity information of a patch. Then we compute point-to-plane geometry error as:

$$E_g = \sum_i^Y \left\| \left( \mathbf{P}[\mathbf{T}_k \cdot \pi^{-1}(P_i^k)] - \mathbf{P}[\mathbf{L}_i] \right) \cdot \mathbf{n}[\mathbf{L}_i] \right\|_2. \qquad (4)$$

Given $E_p$ and $E_g$, the cost function for patch tracking is formulated as:

$$E(\mathbf{T}_k) = \lambda \sum_j^M E_p^j + (1 - \lambda) \sum_j^M E_g^j, \qquad (5)$$

where $j$ denotes the patch $j$, and $M$ indicates the number of patches.

## 2.2. Filter Framework

Our EKF framework aims to tightly combine the color, depth and inertial measurements information. To be specific, we model the camera pose of each frame as the state of the EKF, and will solve it in the EKF. The observation of the EKF include both the color and depth images, and the relationship between the state and the observation is measured by the energy defined in Equation 5. If a state fits exactly to an observation, the energy is zero. On the other hand, the inertial information is used in the Kalman prediction step, which serves in building the motion prediction model.

We follow the traditional Kalman filter to define the variables. A nonlinear discrete time system with state $\mathbf{x}$, observation term $z$, process noise $\boldsymbol{\omega} \sim N(0, \mathbf{Q})$, and update noise $\boldsymbol{\mu} \sim N(0, \mathbf{U})$ in $k$th frame can be written as

$$\mathbf{x}_k = f(\mathbf{x}_{k-1}, \boldsymbol{\omega}_k), \qquad (6)$$

$$\mathbf{z}_k = h(\mathbf{x}_k, \boldsymbol{\mu}_k). \qquad (7)$$

In our framework, the state of the filter is composed of the following elements $\mathbf{x} : \mathbf{T} = (\mathbf{R}, \mathbf{t})$, with a $3 \times 3$ camera rotation matrix $\mathbf{R}$ and a $3 \times 1$ camera translation vector $\mathbf{t}$ related to the world coordinate system which is assigned to be the camera coordinate system of the first frame.In the following, the superscript symbol '+' denotes a-posteriori estimate of a variable calculated from the Kalman update step and '−' denotes a-prior estimate from the Kalman prediction step.

### 2.2.1. Kalman Prediction and State Propagation

Given an a-posteriori estimate $\mathbf{x}_{k-1}^+$ with covariance $\mathbf{P}_{k-1}^+$, the prediction step of the EKF yields a-priori estimate at the next frame:

$$\mathbf{x}_k^- = f(\mathbf{x}_{k-1}^+, 0), \qquad (8)$$

$$\mathbf{P}_k^- = \mathbf{F}_k \mathbf{P}_{k-1}^+ \mathbf{F}_k^T + \mathbf{Q}_k, \qquad (9)$$

with the Jacobians:

$$\mathbf{F}_k = \left. \frac{\partial f}{\partial \mathbf{x}} \right|_{\mathbf{x}_{k-1}^+, 0}. \qquad (10)$$

The key in the Kalman prediction step is to define the function $f$. In our EKF framework, the inertial measurements are employed in the definition. Following [10], we get the actual sensor acceleration $\mathbf{a}$ and angular velocity $\mathbf{w}$ from inertial measurements. We assume that the IMU is synchronized with the camera and acquires measurements with time interval $\tau$ which is much smaller than that of the camera. Hence, we denote $N$ as the number of inertial measurements acquired in two consecutive camera frames, and then merge them together by pre-integration method[11] to predict camera rotation $\Delta \mathbf{R}$ and translation $\Delta \mathbf{t}$ between two consecutive camera frames:

$$\Delta \mathbf{R} = \Phi \cdot \prod_{n=1}^N \mathrm{Exp}(\mathbf{w}_n \cdot \tau) \qquad (11)$$

$$\Delta \mathbf{v} = \sum_{n=1}^N \Delta \mathbf{R}_n \cdot \mathbf{a}_n \cdot \tau \qquad (12)$$

$$\Delta \mathbf{t} = \Phi \cdot \sum_{n=1}^N (\Delta \mathbf{v}_n \cdot \tau + \frac{1}{2} \mathbf{g} \cdot \tau^2 + \frac{1}{2} \Delta \mathbf{R}_n \cdot \mathbf{a}_n \cdot \tau^2). \qquad (13)$$

In the above three equations, the subscript '$n$' denotes the corresponding variable at the $n$th IMU input in consecutive camera frames. Besides, $\Delta \mathbf{v}$ is the accumulated IMU linear velocity from the last camera frame to the current camera frame, and $\Phi$ is the extrinsic matrix from the IMU coordinate to the camera coordinate. $\mathbf{g}$ is the gravity acceleration, and $\mathrm{Exp}(\cdot)$ denotes the exponential map from Lie-algebra to Lie-group. Details about this predition step can be found in [10]. Finally, the states predicted in current frame $k$ can be formulated as:

$$\begin{bmatrix} \mathbf{R}_k^- & \mathbf{t}_k^- \\ \mathbf{0} & 1 \end{bmatrix} = \begin{bmatrix} \Delta \mathbf{R}_k \cdot \mathbf{R}_{k-1}^+ & \Delta \mathbf{t}_k + \mathbf{t}_{k-1}^+ \\ \mathbf{0} & 1 \end{bmatrix}. \qquad (14)$$

### 2.2.2. Kalman Update and Iteration

In traditional extended Kalman update step, the measurement residual is modeled as:

$$\mathbf{y}_k = \mathbf{z}_k - h(\mathbf{x}_k^-, \mathbf{0}). \tag{15}$$

Here, $\mathbf{0}$ means we directly use $\mathbf{x}_k^-$ to calculate the residual without adding any Gausian noise. The updated state is formulated as:

$$\mathbf{x}_k^+ = \mathbf{x}_k^- + \mathbf{K}_k \cdot \mathbf{y}_k, \tag{16}$$

where $\mathbf{K}_k$ is the Kalman gain. In our method, we defined the residual as the photometric and geometric error of patches (equation 5), thus the residual can be formulated as:

$$\mathbf{y}_k = 0 - E(\mathbf{T}_k^-) = -E\left(\begin{bmatrix} \mathbf{R}_k^- & \mathbf{t}_k^- \\ \mathbf{0} & 1 \end{bmatrix}\right). \tag{17}$$

Notice that the deformations of the SE patches, which are used in calculating $\mathbf{y}_k$ by Equation 17, is heavily affected by the camera poses. So after we obtained an updated camera pose by Equation 16, we use the newly updated camera pose to iteratively calculate the deformations of the SE patches and refine the camera poses by Equation 16 again. In this manner, we can estimate a more accurate $\mathbf{x}_k^+$. To be more specific, we use $m$ to denote the iterations, and thus we have:

$$h(\mathbf{x}_{k,m}^+, 0) = E\left(\begin{bmatrix} \mathbf{R}_{k,m}^+ & \mathbf{t}_{k,m}^+ \\ \mathbf{0} & 1 \end{bmatrix}\right), \tag{18}$$

and the Kalman gain respect to each iteration is:

$$\mathbf{K}_{k,m} = \mathbf{P}_k^- \mathbf{H}_{k,m}^T \mathbf{S}_{k,m}^{-1} \tag{19}$$

$$\mathbf{S}_{k,m} = \mathbf{H}_{k,m} \mathbf{P}_k^- \mathbf{H}_{k,m}^T + \mathbf{U}_k. \tag{20}$$

As defined in the begin of section 2.2, $\mathbf{U}_k$ is the covariance matrix of noise $\boldsymbol{\mu}_k$. And the Jacobians updated in every iteration are formulated as:

$$\mathbf{H}_{k,m} = \left.\frac{\partial h}{\partial \mathbf{x}}\right|_{\mathbf{x}_{k,m}^+, 0} \tag{21}$$

Then the updated state of each iteration is calculated as follows:

$$\mathbf{x}_{k,m+1}^+ = \mathbf{x}_{k,m}^+ - \mathbf{K}_{k,m} \cdot h(\mathbf{x}_{k,m}^+, \mathbf{0}). \tag{22}$$

Notice that $\mathbf{x}_{k,0}^+$ is set to be $\mathbf{x}_k^-$. Finally, the iteration is terminated when the absolute value of $\mathbf{K}_{k,m} \cdot h(\mathbf{T}_{k,m}^+, \mathbf{0})$ is below a certain threshold and the covariance matrix is only updated once the process has converged after $\eta$ iterations:

$$\mathbf{P}_k^+ = (\mathbf{I} - \mathbf{K}_{k,\eta} \mathbf{H}_{k,\eta}) \mathbf{P}_k^- \tag{23}$$
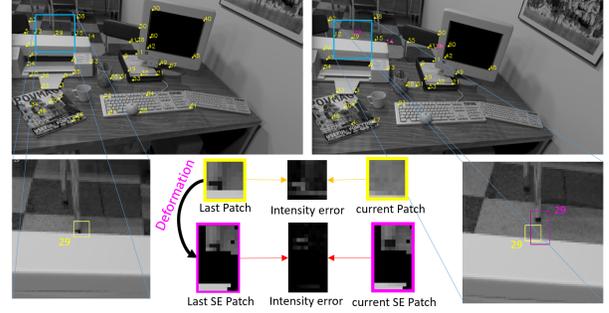


**Fig. 4**. This figure shows the feature tracking results of ours and direct method.

### 2.3. Model Fusion and Patch Update

We use the volumetric truncated signed distance function (TSDF)[12] to incrementally fuse each consecutive depth frame $\mathbf{D}_k$ into one 3D geometry model $\mathbf{M}_k(\mathbf{X})$, with the associated camera pose from Kalman update $\mathbf{R}_k^+, \mathbf{t}_k^+$. Details about depth fusion can be found in[13].

After the reconstruction, we should update patch features for subsequent tracking. We get rid of bad features based on average pixel intensity error, and re-extract squared patch feature for those with non-square shapes affected by the SE effect. Then, we add new features with distinct intensity gradient and sufficient depth information evaluated by FAST corner detector[14] and the number of pixels with available depth information. Finally, patch intensity information is updated by current color and depth information is acquired from the 3D geometry model which has better quality than the current depth image.

## 3. EXPERIMENTS

We first demonstrate the effectiveness of our geometry-aware feature tracking method, which evaluates SE effect and deforms patches for accurate feature tracking in sequences with fast camera motion. Then, we evaluate the benefits of inertial information by comparing our system with and without IMU. Finally, our 3D reconstruction method is compared against state-of-the-art systems in datasets with fast sensor motion.

### 3.1. Evaluation

**Feature tracking.** We compare our feature tracking method against traditional direct method which does not take the SE effect under consideration. In order to achieve equitable comparison, we use a patch-size $10 \times 10$ in both methods and extract no more than 100 patches in each frame. Fig. 4 shows the tracking results of a patch feature in two consecutive frames. The tracking result of traditional method is severely influenced by SE effect and got bad intensity error,

**Table 1**. Comparison of Patch Feature Tracking

| Type | Dataset | AIE | |
|------|---------|-----|-----|
| | | DM | Ours |
| slow | TUM_freiburg1_desk | 13.3756 | **9.53** |
| | ICL_NUIM_lr_kt2 | 4.8981 | 4.0825 |
| | Dorm_slow | 8.1312 | 7.8934 |
| fast | ICL Fast Motion | 17.219 | **7.8328** |
| | Dorm_fast | 13.9011 | **7.9325** |

while our method deforms the patch and eliminates the influence caused by SE effect to get lower intensity error.

We compare on several datasets, contains ICL datasets[15, 16], TUM datasets[17] and our datasets gathered by a handheld sensor. The average intensity error(AIE) of patches are listed in Table 1. All datasets are divided into slow and fast depending on the qualities of recorded images. To be more specific, as there is no explicit criteria for dividing camera speed, thus we empirically set, based on the unified characteristics of most public datasets, the motion without creating image motion blur as slow camera motion, and the motion which creates severe image blur as fast camera motion. From the table, we find that our method gets lower AIE in all datasets, especially in datasets with fast camera motion.

**IMU evaluation.** To verify whether the integration of IMU helps to reconstruct the scene geometry during fast camera motion, we compare the results with and without IMU on two datasets with slow and fast camera motions, respectively. As shown in Fig. 5, on the dataset with slow camera motion, the system without IMU works on-par with the complete system, while it fails to reconstruct the model for fast camera motions. Fig. 6 demonstrates the details of camera motions in the two datasets. From the figure, we find that in the fast dataset, there exist some subsequences with large linear and angular velocities of camera, which cause the system without IMU fails to track camera poses. Notice that other fast datasets used in our experiments also contain this kinds of subsequences.

### 3.2. Comparison

We compare our 3D reconstruction systems with InfiniTAM[1], a typical voxel based scene reconstruction method, Bundlefusion[5] which proposed an efficient global pose optimization algorithm, and a surfel based method ElasticFusion[6] which contains loop closure and executes model refinement through non-rigid surface deformations.

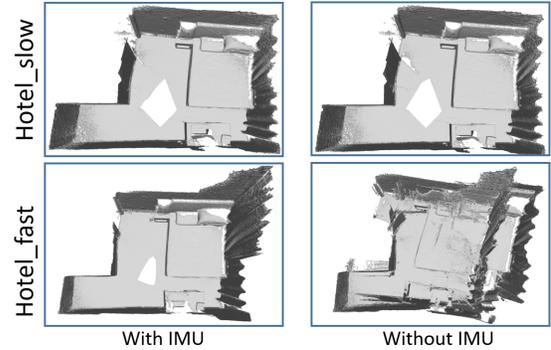The results of sequence Dorm_fast which reconstruct the entire scene are exhibited in Fig. 7. As Bundle-



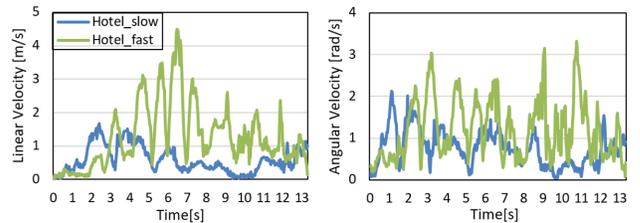**Fig. 5**. The reconstruction results of a hotel under slow and fast camera motion.



**Fig. 6**. Camera linear and angular velocity of the two sequences used in IMU evaluation.
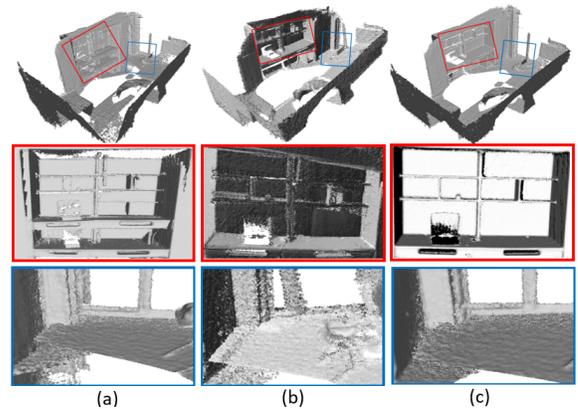


**Fig. 7**. Comparison of reconstruction with fast camera motion in (a)InfiniTAM (b)ElasticFusion and (c)Ours.

Fusion fails once the camera speeds up and subsequently restarts when the camera slows down, thus we only show its reconstruction precess in our supplementary video https://www.youtube.com/watch?v=Jy3SGqWuhp8. From the Fig. 7 we find that InfiniTAM can not maintain consistency of the reconstructed geometry, which is mainly caused by the inaccurate camera pose estimation and the large accumulated error. Meanwhile, the loop closure function of ElasticFusion, aiming to eliminate accumulated error, is always invalid in fast camera motion, and thus leads to the fail reconstruction of the parts shown with red and blue bounding

boxes. In the opposite, our system reconstructs a good geometry of the scene even without loop closure.

We encourage the reader to watch our video for a better visualization of comparison results.

## 4. CONCLUSION AND FUTURE WORK

We present a real-time system for indoor scene reconstruction by tightly-coupled RGB-D-Inertial information with an extended Kalman filter. The key feature of our method is that it can estimate camera pose and reconstruct 3D scene model with fast camera motion. In addition, we explore the SE effect and propose a geometry-aware patch deformation method to eliminate the influence during feature tracking. However, our system has not achieved loop closure with fast camera motion. The reason is that the degraded image information caused by fast camera motion, such as image motion blur, results in the difficulties in loop detection(or feature association) of loop closure method. In future work, we wish to address the problem of loop closure under fast camera motion.

## 5. REFERENCES

[1] Olaf Kähler, Victor Adrian Prisacariu, Carl Yuheng Ren, Xin Sun, Philip Torr, and David Murray, "Very high frame rate volumetric integration of depth images on mobile devices," *IEEE transactions on visualization and computer graphics*, vol. 21, no. 11, pp. 1241–1250, 2015.

[2] Szymon Rusinkiewicz and Marc Levoy, "Efficient variants of the icp algorithm," in *3-D Digital Imaging and Modeling, 2001. Proceedings. Third International Conference on*. IEEE, 2001, pp. 145–152.

[3] Christian Forster, Zichao Zhang, Michael Gassner, Manuel Werlberger, and Davide Scaramuzza, "Svo: Semidirect visual odometry for monocular and multi-camera systems," *IEEE Transactions on Robotics*, vol. 33, no. 2, pp. 249–265, 2017.

[4] Raul Mur-Artal, Jose Maria Martinez Montiel, and Juan D Tardos, "Orb-slam: a versatile and accurate monocular slam system," *IEEE Transactions on Robotics*, vol. 31, no. 5, pp. 1147–1163, 2015.

[5] Angela Dai, Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Christian Theobalt, "Bundlefusion: Real-time globally consistent 3d reconstruction using on-the-fly surface reintegration," *ACM Transactions on Graphics (TOG)*, vol. 36, no. 4, pp. 76a, 2017.

[6] Thomas Whelan, Renato F Salas-Moreno, Ben Glocker, Andrew J Davison, and Stefan Leutenegger, "Elasticfusion: Real-time dense slam and light source estimation,"

*The International Journal of Robotics Research*, vol. 35, no. 14, pp. 1697–1716, 2016.

[7] Chenggang Yan, Yongdong Zhang, Jizheng Xu, Feng Dai, Liang Li, Qionghai Dai, and Feng Wu, "A highly parallel framework for hevc coding unit partitioning tree decision on many-core processors," *IEEE Signal Processing Letters*, vol. 21, no. 5, pp. 573–576, 2014.

[8] Tristan Laidlow, Michael Bloesch, Wenbin Li, and Stefan Leutenegger, "Dense rgb-d-inertial slam with map deformations," in *Intelligent Robots and Systems (IROS), 2017 IEEE/RSJ International Conference on*. IEEE, 2017, pp. 6741–6748.

[9] Michael Bloesch, Michael Burri, Sammy Omari, Marco Hutter, and Roland Siegwart, "Iterated extended kalman filter based visual-inertial odometry using direct photometric feedback," *The International Journal of Robotics Research*, vol. 36, no. 10, pp. 1053–1072, 2017.

[10] Anastasios I Mourikis and Stergios I Roumeliotis, "A multi-state constraint kalman filter for vision-aided inertial navigation," in *Robotics and automation, 2007 IEEE international conference on*. IEEE, 2007, pp. 3565–3572.

[11] Christian Forster, Luca Carlone, Frank Dellaert, and Davide Scaramuzza, "On-manifold preintegration for real-time visual–inertial odometry," *IEEE Transactions on Robotics*, vol. 33, no. 1, pp. 1–21, 2017.

[12] Brian Curless and Marc Levoy, "A volumetric method for building complex models from range images," in *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*. ACM, 1996, pp. 303–312.

[13] Hao Zhang and Feng Xu, "Mixedfusion: Real-time reconstruction of an indoor scene with dynamic objects," *IEEE Transactions on Visualization and Computer Graphics*, 2017.

[14] Edward Rosten and Tom Drummond, "Machine learning for high-speed corner detection," in *European conference on computer vision*. Springer, 2006, pp. 430–443.

[15] Ankur Handa, Richard A Newcombe, Adrien Angeli, and Andrew J Davison, "Real-time camera tracking: When is high frame-rate best?," in *European Conference on Computer Vision*. Springer, 2012, pp. 222–235.

[16] A. Handa, T. Whelan, J.B. McDonald, and A.J. Davison, "A benchmark for RGB-D visual odometry, 3D reconstruction and SLAM," in *IEEE Intl. Conf. on Robotics and Automation, ICRA*, Hong Kong, China, May 2014.

[17] 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.