

Dynamic Object Tracking and Masking for Visual SLAM

Jonathan Vincent, Mathieu Labbé, Jean-Samuel Lauzon, François Grondin,
Pier-Marc Comtois-Rivet, François Michaud

Abstract—In dynamic environments, performance of visual SLAM techniques can be impaired by visual features taken from moving objects. One solution is to identify those objects so that their visual features can be removed for localization and mapping. This paper presents a simple and fast pipeline that uses deep neural networks, extended Kalman filters and visual SLAM to improve both localization and mapping in dynamic environments (around 14 fps on a GTX 1080). Results on the dynamic sequences from the TUM dataset using RTAB-Map as visual SLAM suggest that the approach achieves similar localization performance compared to other state-of-the-art methods, while also providing the position of the tracked dynamic objects, a 3D map free of those dynamic objects, better loop closure detection with the whole pipeline able to run on a robot moving at moderate speed.

I. INTRODUCTION

To perform tasks effectively and safely, autonomous mobile robots need accurate and reliable localization from their representation of the environment. Compared to LIDARs (Light Detection And Ranging sensors) and GPS (Global Positioning System), using visual images for Simultaneous Localization and Mapping (SLAM) adds significant information about the environment [1], such as color, textures, surface composition that can be used for semantic interpretation of the environment. Standard visual SLAM (vSLAM) techniques perform well in static environments by being able to extract stable visual features from images. However, in environments with dynamic objects (e.g., people, cars, animals), performance decreases significantly because visual features may come from those objects, making localization less reliable [1]. Deep learning architectures have recently demonstrated interesting capabilities to achieve semantic segmentation from images, outperforming traditional techniques in tasks such as image classification [2]. For instance, Segnet [3] is commonly used for semantic segmentation [4]. It uses an encoder and a decoder to achieve pixel wise semantic segmentation of a scene.

This paper introduces a simple and fast pipeline that uses neural networks, extended Kalman filters and vSLAM algorithm to deal with dynamic objects. Experiments conducted on the TUM dataset demonstrate the robustness of

the proposed method. Our research hypothesis is that a deep learning algorithm can be used to semantically segment object instances in images using *a priori* semantic knowledge of dynamic objects, enabling the identification, tracking and removal of dynamic objects from the scenes using extended Kalman filters to improve both localization and mapping in vSLAM. By doing so, the approach, referred to as Dynamic Object Tracking and Masking for vSLAM (DOTMask)¹ aims at providing six benefits: 1) increased visual odometry performance; 2) increased quality of loop closure detection; 3) produce 3D maps free of dynamic objects; 4) tracking of dynamic objects; 5) modular and fast pipeline.

The paper is organized as follows. Section II presents related work of approaches taking into consideration dynamic objects during localization and during mapping. Section III describes our approach applied as a pre-processing module to RTAB-Map [5], a vSLAM approach. Section IV presents the experimental setup, and Section V provides comparative results on dynamic sequences taken from the TUM dataset.

II. RELATED WORK

Some approaches take into consideration dynamic objects during localization. For instance, BaMVO [6] uses a RGB-D camera to estimate ego-motion. It uses a background model estimator combined with an energy-based dense visual odometry technique to estimate the motion of the camera. Li *et al.* [7] developed a static point weighting method which calculates a weight for each edge point in a keyframe. This weight indicates the likelihood of that specific edge point being part of the static environment. Weights are determined by the movement of a depth edge point between two frames and are added to an Intensity Assisted Iterative Closest Point (IA-ICP) method used to perform the registration task in SLAM. Sun *et al.* [8] present a motion removal approach to increase the localization reliability in dynamic environments. It consists of three steps: 1) detecting moving objects' motion based on ego-motion compensated using image differencing; 2) using a particle filter for tracking; and 3) applying a Maximum-A-Posterior (MAP) estimator on depth images to determine the foreground. This approach is used as the frontend of Dense Visual Odometry (DVO) SLAM [9]. Sun *et al.* [10] uses a similar foreground technique but instead of using a MAP they use a foreground model which is updated on-line. All of these approaches demonstrate good localization results using the Technical University of Munich (TUM) dataset [11], however, mapping is yet to be addressed.

¹<https://github.com/introlab/dotmask>

This work was supported by the Institut du véhicule innovant (IVI), Mitacs, InnovÉÉ and NSERC. J. Vincent, M. Labbé, J.-S. Lauzon, F. Grondin and F. Michaud are with the Interdisciplinary Institute for Technological Innovation (3IT), Dept. Elec. Eng. and Comp. Eng., Université de Sherbrooke, 3000 boul. de l'Université, Québec (Canada) J1K 0A5. P.-M. Comtois-Rivet is with the Institut du Véhicule Innovant (IVI), 25, boul. Maisonneuve, Saint-Jérôme, Québec (Canada), J5L 0A1. {Jonathan.Vincent2, Mathieu.m.Labbe, Jean-Samuel.Lauzon, Francois.Grondin2, Francois.Michaud}@USherbrooke.ca, Pmcrivet@ivisolutions.ca

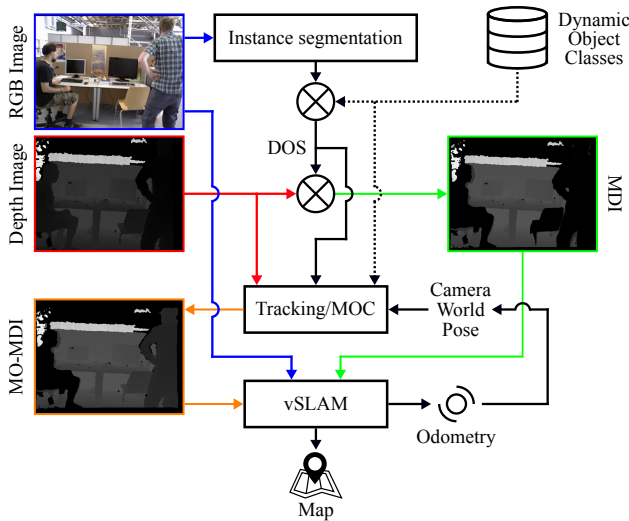


Fig. 1: Architecture of DOTMask

SLAM++ [12] and Semantic Fusion [13] focus on the mapping aspect of SLAM in dynamic environments. SLAM++ [12] is an object-oriented SLAM which achieves efficient semantic scene description using 3D object recognition. SLAM++ defines objects using areas of interest to subsequently locate and map them. However, it needs predefined 3D object models to work. Semantic Fusion [13] creates a semantic segmented 3D map in real time using RGB-CNN [14], a convolutional deep learning neural network, and a dense SLAM algorithm. However, SLAM++ and Semantic Fusion do not address SLAM localization accuracy in dynamic environments, neither do they remove dynamic objects in the 3D map.

Other approaches use deep learning algorithm to provide improved localisation and mapping. Fusion++ [15] and MID-Fusion [16] uses object-level octree-based volumetric representation to estimate both the camera pose and the object positions. They use deep learning techniques to segment object instances. DynaSLAM [17] proposes to combine multi-view geometry models and deep-learning-based algorithms to detect dynamic objects and to remove them from the images prior to a vSLAM algorithm. They also uses inpainting to recreate the image without object occlusion. DynaSLAM achieves impressive results on the TUM dataset. However, these approaches are not optimized for real-time operation.

III. DYNAMIC OBJECT TRACKING AND MASKING FOR vSLAM

The objective of our work is to provide a fast and complete solution for visual SLAM in dynamic environments. Figure 1 illustrates the DOTMask pipeline. As a general overview of the approach, a set of objects of interest (OOI) are defined using *a priori* knowledge and understanding of dynamic objects classes that can be found in the environment. Instance segmentation is done using a neural network trained to identify the object classes from an RGB image. For each dynamic object instance, its bounding box, class type and binary mask are grouped for convenience and referred as the dynamic object state (DOS). The binary mask of the DOS

is then applied to the original depth image, resulting in a masked depth image (MDI). The DOS is also sent to the Tracking module. After computing a 3D centroid for each masked object, the Tracking module predict the position and velocity of the objects. This information is then used by the Moving Object Classification module (MOC) to classify the object as idle or not based on its class, its estimated velocity and its shape deformation. Moving objects are removed from the original depth image, resulting in the Moving Object Masked Depth Image (MO-MDI). The original RGB image, the MDI and the MO-MDI are used by the vSLAM algorithm. It uses the depth images as a mask for feature extraction thus ignoring features from the masked regions. The MO-MDI is used by the visual odometry algorithm of the vSLAM approach while the MDI is used by both its mapping and loop closure algorithms, resulting in a map free of dynamic objects while still being able to use the features of the idle objects for visual odometry. The updated camera pose is then used in the Tracking module to estimate the position and velocity of the dynamic objects resulting in a closed loop.

A. Instance Segmentation

Deep learning algorithms such as Mask R-CNN recently proved to be useful to accomplish instance semantic segmentation [4]. A recent and interesting architecture for fast instance segmentation is the YOLACT [18] and its update YOLACT++ [19]. This network aims at providing similar results as the Mask-RCNN or the Fully Convolutional Instance-aware Semantic Segmentation (FCIS) [20] but at a much lower computational cost. YOLACT and YOLACT++ can achieve real-time instance segmentation. Development in neural networks has been incredibly fast in the past few years and probably will be in the years to come. DOTMask was designed to be modular and can easily change the neural network used in the pipeline. In its current state, DOTMask works with Mask-RCNN, YOLACT and YOLACT++. The YOLACT is much faster than the two others and the loss in precision doesn't impact our results. This is why this architecture is used in our tests. The instance segmentation module takes the input RGB image and outputs the bounding box, class and binary mask for each instance.

B. Tracking Using EKF

Using the DOS from the Instance Segmentation module and odometry from vSLAM, the Tracking module predicts the pose and velocity of the objects in the world frame. This is useful when the camera is moving at speed similar to the objects to track (e.g., moving cars on the highway, robot following a pedestrian) or when idle objects have a high amount of features (e.g., person wearing a plaid shirt).

First, the Tracking module receives the DOS and the original depth image as a set, defined as $\mathbf{D}^k = \{\mathbf{d}_1^k, \dots, \mathbf{d}_I^k\}$, where $\mathbf{d}_i^k = \{\mathbf{T}^k, \mathbf{B}_i^k, \zeta_i^k\}$ is the object instance detected by the Instance Segmentation module, with $i \in I$, $I = \{1, \dots, L\}$, L being the total number of object detection in the frame at time k . $\mathbf{T} \in \mathbb{R}^{m \times n}$ is the depth image ,

$\mathbf{B} \in \mathbb{Z}_2^{m \times n}$ is the binary mask and $\zeta \in J$ is the class ID, with $J = \{1, \dots, W\}$, and W is the number of total trained classes in the Instance Segmentation module.

The DOS and the original depth image are used by EKF to estimate the dynamic objects positions and velocities. EKF provides steady tracking of each object instance corresponding to the object type detected by the neural network. An EKF is instantiated for each new object, and *a priori* knowledge from the set of dynamic object classes defines some of the filter's parameters. This instantiation is made using the following parameters: the class of the object, its binary mask and its 3D centroid position. The 3D centroid is defined as the center of the corresponding bounding box. If the tracked object is observed in the DOS, its position is updated accordingly, otherwise its predicted position using EKF is used. If no observations of the object are made for e number of frames, the object is considered removed from the scene and therefore the filter is discarded. The Tracking module outputs the estimated velocity of the objects to the MOC module. The MOC module will classify the objects as idle or not based on the object class, the filter velocity estimation and the object deformation.

To explain further how the Tracking module works, the following subsections presents in more details the Prediction and Update steps of EKF used by DOTMask.

1) *Prediction*: Let us define the hidden state $\mathbf{x} \in \mathbb{R}^{6 \times 1}$ as the 3D position and velocity of an object referenced in the global map in Cartesian coordinates. The *a priori* estimate of the state at time $k \in \mathbb{N}$ is predicted based on the previous state at time $k-1$ as in (1):

$$\hat{\mathbf{x}}^{k|k-1} = \mathbf{F}\hat{\mathbf{x}}^{k-1|k-1} \text{ with } \mathbf{F} = \begin{bmatrix} \mathbf{I}_3 & \Delta t \mathbf{I}_3 \\ \mathbf{0}_3 & \mathbf{I}_3 \end{bmatrix} \quad (1)$$

where $\mathbf{F} \in \mathbb{R}^{6 \times 6}$ is the state transition matrix, $\Delta t \in \mathbb{R}_+$ is the time between each prediction, $\mathbf{0}_3$ is a 3×3 zero matrix and \mathbf{I}_3 is a 3×3 identity matrix. Note that the value of Δt is redefined before each processing cycle.

The *a priori* estimate of the state covariance ($\mathbf{P}^{k|k-1} \in \mathbb{R}^{6 \times 6}$) at time k is predicted based on the previous state at time $k-1$ as given by (2):

$$\mathbf{P}^{k|k-1} = \mathbf{F}\mathbf{P}^{k-1|k-1}\mathbf{F}^T + \mathbf{Q} \quad (2)$$

where $\mathbf{Q} \in \mathbb{R}^{6 \times 6}$ is the process noise covariance matrix defined using the random acceleration model (3):

$$\mathbf{Q} = \mathbf{\Gamma}\mathbf{\Sigma}\mathbf{\Gamma}^T \text{ with } \mathbf{\Gamma} = \begin{bmatrix} \frac{\Delta t^2}{2}\mathbf{I}_{3 \times 3} & \Delta t^2\mathbf{I}_{3 \times 3} \end{bmatrix}^T \quad (3)$$

where $\mathbf{\Gamma} \in \mathbb{R}^{6 \times 3}$ is the mapping between the random acceleration vector $\mathbf{a} \in \mathbb{R}^3$ and the state \mathbf{x} , and $\mathbf{\Sigma} \in \mathbb{R}^{3 \times 3}$ is the covariance matrix of \mathbf{a} . The acceleration components a_x , a_y and a_z are assumed to be uncorrelated.

The dynamic of every detected objects may vary greatly depending on its class. For instance, a car does not have the same dynamic as a mug. To better track different types of objects, a covariance matrix is defined for each class to better represent their respective process noise.

2) *Update*: In EKF, the Update step starts by evaluating the innovation $\tilde{\mathbf{y}}^k$ defined as (4):

$$\tilde{\mathbf{y}}^k = \mathbf{z}^k - \hat{\mathbf{h}}^k(\hat{\mathbf{x}}^{k|k-1}) \quad (4)$$

where $\mathbf{z}^k \in \mathbb{R}^3$ is a 3D observation of a masked object in reference to the camera for each object instance, with $\mathbf{z} = [z_x \ z_y \ z_z]^T$, $z_x = (\mu_x - C_x)z_z/f_x$ and $z_y = (\mu_y - C_y)z_z/f_y$, where C_x and C_y are the principal center point coordinate and f_x and f_y are the focal lengths expressed in pixels. z_z is approximated using the average depth from the masked region on the depth image. The expressions μ_x and μ_y stand for the center of the bounding box.

To simplify the following equations, (s , c) represent respectively the sine and cosine operations of the the Euler angles ϕ , θ , ψ (roll, pitch, yaw). $h(\mathbf{x}^k) \in \mathbb{R}^4$ is the observation function which maps the true state space \mathbf{x}^k to the observed state space \mathbf{z}^k . $\hat{\mathbf{h}}(\mathbf{x}^k)$ is the three first terms of $h(\mathbf{x}^k)$. However, in our case, the transform between those spaces is not linear, justifying the use of EKF. The non-linear rotation matrix used to transform the estimate state $\hat{\mathbf{x}}^k$ in the observed state \mathbf{z}^k follows the (x , y , z) Tait-Bryan convention and is given by $h(\hat{\mathbf{x}}^k) = [h_\phi \ h_\theta \ h_\psi \ 1]$, where:

$$\begin{aligned} h_\phi &= (c_\phi c_\theta)\hat{x}_x + (c_\phi s_\theta s_\psi - c_\psi s_\phi)\hat{x}_y + (s_\phi s_\psi + c_\phi c_\psi s_\theta)\hat{x}_z + c_x \\ h_\theta &= (c_\theta s_\phi)\hat{x}_x + (c_\theta c_\psi + s_\phi s_\theta s_\psi)\hat{x}_y + (c_\psi s_\phi s_\theta - c_\phi s_\psi)\hat{x}_z + c_y \\ h_\psi &= -(s_\theta)\hat{x}_x + (c_\theta s_\psi)\hat{x}_y + (c_\theta c_\psi)\hat{x}_z + c_z \end{aligned} \quad (5)$$

and c_x , c_y and c_z are the coordinate of the camera referenced to the world, which is derived using vSLAM odometry.

The innovation covariance $\mathbf{S}_k \in \mathbb{R}^{3 \times 3}$ is defined as follows, where the expression $\mathbf{H}^k \in \mathbb{R}^{3 \times 6}$ stands for the Jacobian of $h(\hat{\mathbf{x}}^k)$:

$$\mathbf{S}_k = \mathbf{H}^k \mathbf{P}^{k|k-1} (\mathbf{H}^k)^T + \mathbf{R}^k \quad (6)$$

where $\mathbf{R}^k \in \mathbb{R}^{3 \times 3}$ is the covariance of the observation noise, its diagonal terms stand for the imprecision of the RGB-D camera. The near optimal Kalman gain $\mathbf{K}^k \in \mathbb{R}^{3 \times 3}$ is defined as follows:

$$\mathbf{K}^k = \mathbf{P}^{k|k-1} (\mathbf{H}^k)^T (\mathbf{S}_k)^{-1} \quad (7)$$

Finally, the updated state estimate $\hat{\mathbf{x}}^{k|k}$ and the covariance estimate are given respectively by (8) and (9).

$$\hat{\mathbf{x}}^{k|k} = \hat{\mathbf{x}}^{k|k-1} + \mathbf{K}^k \tilde{\mathbf{y}}^k \quad (8)$$

$$\mathbf{P}^{k|k} = (\mathbf{I}_6 - \mathbf{K}^k \mathbf{H}^k) \mathbf{P}^{k|k-1} \quad (9)$$

C. Moving Object Classification

The MOC module classify dynamic objects as either moving or idle. It takes as inputs the dynamic objects class, velocity and mask. The object velocity comes from the tracking module estimation. The object class and mask are directly obtained from the DOS. The object class defines if the object is rigid or not. The deformation of non-rigid object is computed using the intersection over union (IoU) of the masks of the object at time k and $k-1$. The IoU algorithm takes two arbitrary convex shape \mathbf{M}^{k-1} , \mathbf{M}^k and is defined as $\text{IoU} = |\mathbf{M}^k \cap \mathbf{M}^{k-1}| / |\mathbf{M}^k \cup \mathbf{M}^{k-1}|$, where

TABLE I: Experimental Parameters

Description	Value
Frame to terminate object tracking	10
Score threshold (s)	0.1
Maximum number of observations (m)	5
Velocity threshold for a person	0.01 m/sec
Velocity threshold for the other objects	0.1 m/sec
Random acceleration for a person	0.62 m/s ²
Random acceleration for other objects	1.0 m/s ²

$|\dots|$ is the cardinality of the set. A dynamic object is classified as moving if its velocity is higher than a predefined threshold or if it is a non-rigid object with an IoU above another predefined threshold. The original depth image is then updated resulting in the MO-MDI. The MO-MDI is sent to the vSLAM odometry to update the camera pose.

IV. EXPERIMENTAL SETUP

To test our DOTMask approach, we chose to use the TUM dataset because it presents challenging indoor dynamic RGB-D sequences with ground truth to evaluate visual odometry techniques. Also, TUM is commonly used to compare with other state-of-the-art techniques. We used sequences in low dynamic and highly dynamic environments.

For our experimental setup, ROS is used as a middleware to make the interconnections between the input images, segmentation network, EKF and RTAB-Map. The deep learning library PyTorch is used for the instance segmentation algorithm. The ResNet-50-FPN backbone is used for the YOLACT architecture because this configuration achieves the best results at a higher framerate [18]. Our Instance segmentation module is based on the implementation of YOLACT by dbolya² and its pre-trained weights. The network is trained on all 91 classes of the COCO dataset. The COCO dataset is often used to compare state-of-the-art instance segmentation approaches, which is why we chose to use it in our trials. In our tests, person, chair, cup and bottle are the the OOI used because of their presence in the TUM dataset and in our in-house tests. The RTAB-Map library [5] is also used, which includes various state-of-the-art visual odometry algorithms, a loop closure detection approach and a 3D map render.

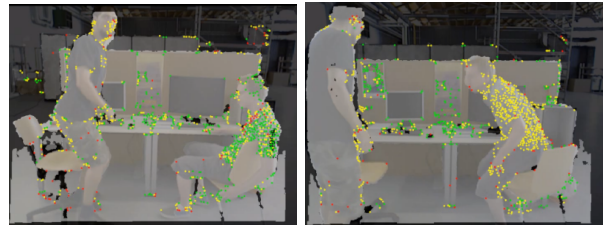
Table I presents the parameters used for DOTMask in our trials, based on empirical observations in the evaluated TUM sequences and our understanding of the nature of the objects. A probability threshold p and a maximum instance number m are used to reduce the number of object instances to feed into the pipeline. Only detections with a score above p are used and at maximum, m objects detections are processed. This provides faster and more robust tracking.

V. RESULTS

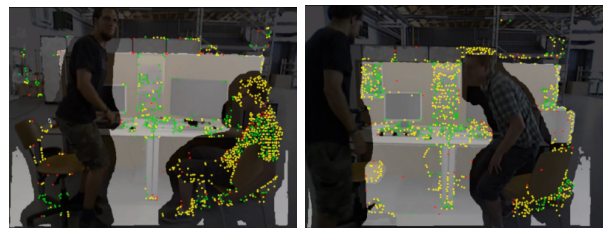
Trials were conducted in comparison with approaches by Kim and Kim [6], Sun *et al.* [8], Bescos *et al.* [17] and RTAB-Map, the latter being also used with DOTMask. Figure 2a shows two original RGB frames in the TUM



(a) Original RGB Image



(b) RGB and depth image superposed without DOTMask



(c) RGB and depth image superposed with DOTMask

Fig. 2: RTAB-Map features (colored dots) not appearing on moving objects with DOTMask

dataset, along with their superimposed RGB and depth images with features used by RTAB-Map (Fig. 2b) and with DOTMask (Fig. 2c). Using the depth image as a mask to filter outlying features, dynamic objects (i.e., humans and chairs in this case) are filtered out because the MDI includes the semantic mask. The MO-MDI is used by RTAB-Map to compute visual odometry, keeping only the features from static objects as seen in Fig. 2c (left vs right) with the colored dots representing visual features used for visual odometry. In the left image of Fig. 2c, the man on the left is classified by the Tracking module as moving, while the man on the right is classified as being idle, resulting in keeping his visual features. In the right image of Fig. 2c, the man on the right is also classified as moving because he is standing up, masking his visual features. Figure 3 illustrates the influence of MDI, which contains the depth mask of all the dynamic objects, either idle or not, to generate a map free of dynamic objects. This has two benefits: it creates a more visually accurate 3D rendered map, and it improves loop closure detection. The differences in the 3D generated maps between RTAB-Map without and with DOTMask are very apparent: there are less artifacts of dynamic objects and less drifting. The *fr3/walking_static* sequence shows improved quality in the map, while the *fr3/walking_rpy* sequence presents some undesirable artifacts. These artifacts are caused either by the mask failing to identify dynamic objects that are tilted or upside down or by the time delay between the RGB image and its corresponding depth image. The *fr3/sitting_static*

²<https://github.com/dbolya/yolact>

TABLE II: Absolute Transitional Error (ATE) RMSE in cm

TUM Seqs	BaMVO	Sun et al.	DynaSLAM	RTAB-Map	DOTMask	Impr. (%)
fr3/sit_static	2.48	-	-	1.70	0.60	64.71
fr3/sit_xyz	4.82	3.17	1.5	1.60	1.80	-12.50
fr3/wlk_static	13.39	0.60	2.61	10.7	0.80	92.52
fr3/wlk_xyz	23.26	9.32	1.50	24.50	2.10	91.42
fr3/wlk_rpy	35.84	13.33	3.50	22.80	5.30	76.75
fr3/wlk_halfsph	17.38	12.52	2.50	14.50	4.00	72.41

TABLE III: Loop Closure Analysis

TUM Seqs	RTAB-Map			DOTMask		
	Nb loop	T_{err} (cm)	R_{err} (deg)	Nb loop	T_{err} (cm)	R_{err} (deg)
fr3/sit_static	33	1.80	0.26	1246	0.60	0.21
fr3/sit_xyz	288	2.10	0.42	1486	2.50	0.45
fr3/wlk_static	105	9.00	0.18	1260	7.00	0.15
fr3/wlk_xyz	55	6.5	0.99	1516	2.9	0.45
fr3/wlk_halfs.	121	5.90	0.84	964	4.90	0.79
fr3/wlk_rpy	94	6.7	1.06	965	6.00	1.04

shows the result when masking idle object, resulting in completely removing the dynamic objects from the scene.

Table II characterizes the overall SLAM quality in terms of absolute trajectory error (ATE). In almost all cases, DOTMask improves the ATE compared to RTAB-Map alone (as seen in the last column of the table). Table II characterizes the overall SLAM quality in terms of absolute trajectory error (ATE). While DynaSLAM is better in almost every sequences, DOTMask is not far off with closer values compared to the other techniques.

Table III presents the number of loop closure detections, the mean translation error (T_{err}) and the mean rotational error (R_{err}) on each sequences both with and without DOTMask. In all sequences, DOTMask helps RTAB-Map to make more loop closures while also lowering both mean errors. Since loop closure features are computed from the depth image (MDI), using DOTMask forces RTAB-Map to use only features from static object hence providing better loop closures.

On the *fr3/sitting_xyz* sequence, RTAB-Map alone provides better performance in both ATE and loop closure detection. In this entire sequence, the dynamic objects do not move. While the MO-MDI enables features from idle dynamic objects to be used by the odometry algorithm, the MDI does not enables those same features for the loop closure algorithm. Since nothing is moving in this particular sequence, all features will help to provide a better localization. However, this case is not representative of dynamic environments.

Table IV presents the average computation time to process a frame for each approach without vSLAM and odometry algorithms. Results are processed on a computer equipped with a GTX 1080 GPU and a I5-8600k CPU. DOTMask was also tested on a laptop with a GTX 1050 where it achieved an average of 8 frames per second. At 70 ms, it can run on

TABLE IV: Timing Analysis

Approach	Img. Res.	Avg. Time	CPU	GPU
BaMVO.	320×240	42.6 ms	i7 3.3GHz	-
Sun et al.	640×480	500 ms	i5	-
DynaSLAM	640×480	500 ms	-	-
DOTMask	640×480	70 ms	i5-8600K	GTX1080
DOTMask	640×480	125 ms	i7-8750H	GTX1050

a mobile robot operating at a moderate speed. The fastest method is BaMVO with only 42 ms cycle time.

Figure 4 shows the tracked dynamic objects in the ROS visualizer RViz. DOTMask generates ROS transforms to track the position of the objects. Those transforms could easily be used in other ROS applications. Figure 5 shows the difference between RTAB-Map and DOTMask in a real scene where a robot moves at a similar speed as dynamic objects (chairs and humans). The pink and blue lines represent the odometry of RTAB-Map without and with DOTMask. These results suggest qualitatively that DOTMask improves the odometry and the 3D map.

VI. CONCLUSION

This paper presents DOTMask, a fast and modular pipeline that uses a deep learning algorithm to semantically segment images, enabling the tracking and masking of dynamic objects in scenes to improve both localization and mapping in vSLAM. Our approach aims at providing a simple and complete pipeline to allow mobile robots to operate in dynamic environments. Results on the TUM dataset suggest that using DOTMask with RTAB-Map provides similar performance compared to other state-of-the-art localization approaches while providing an improved 3D map, dynamic objects tracking and higher loop closure detection. While DOTMask does not outperform DynaSLAM on the TUM dataset or outrun BaMVO, it reveals to be a good compromise for robotic applications. Because DOTMask pipeline is highly modular, it can also evolve with future improvements of deep learning architectures and new sets of dynamic object classes. In future work, we want to use the tracked dynamic objects to create a global 3D map with object permanence, and explore more complex neural networks³ to add body keypoint tracking, which could significantly improve human feature extraction. We would also like to explore techniques to detect outlier segmentations from the neural network to improve robustness.

REFERENCES

- [1] J. Fuentes-Pacheco, J. Ruiz-Ascencio, and J. M. Rendón-Mancha, "Visual simultaneous localization and mapping: A survey," *Artificial Intelligence Review*, vol. 43, no. 1, pp. 55–81, 2015.
- [2] D. Ciregan, U. Meier, and J. Schmidhuber, "Multi-column deep neural networks for image classification," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2012, pp. 3642–3649.
- [3] V. Badrinarayanan, A. Kendall, and R. Cipolla, "Segnet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 39, no. 12, pp. 2481–2495, 2017.

³https://github.com/dajucug/Mask-RCNN-TF_detection-human_segment-body_keypoint-regression

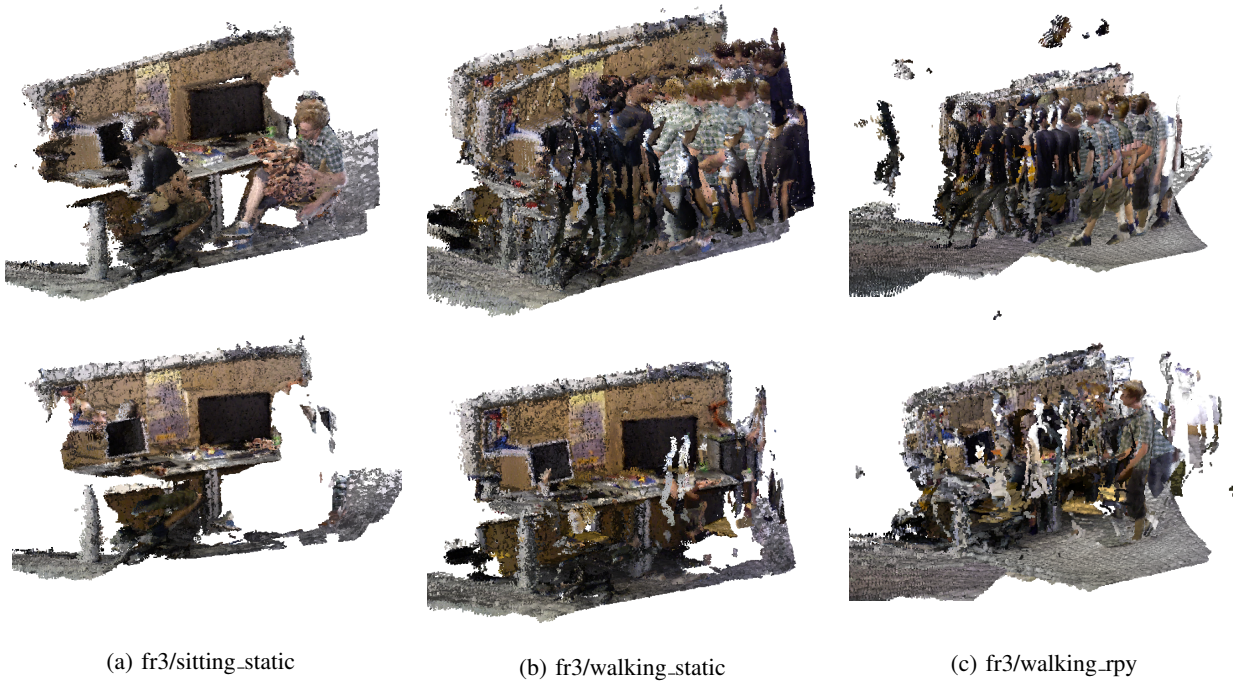
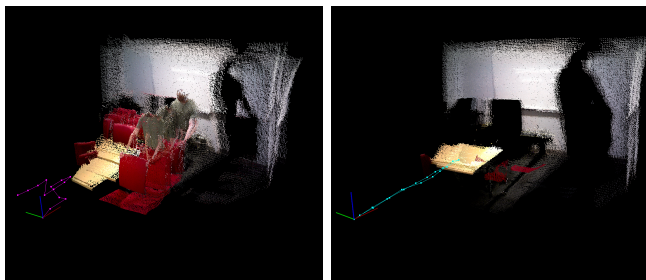


Fig. 3: RTAB-Map 3D rendered map from the TUM sequences, without (top) and with (bottom) DOTMask



Fig. 4: Position of tracked dynamic objects shown in RVIZ



(a) RTAB-Map alone (b) RTAB-Map with DOTMask

Fig. 5: 3D map and odometry improved with DOTMask

[4] A. Garcia-Garcia, S. Orts-Escolano, S. Oprea, V. Villena-Martinez, and J. Garcia-Rodriguez, "A review on deep learning techniques applied to semantic segmentation," *arXiv preprint arXiv:1704.06857*, 2017.

[5] M. Labbé and F. Michaud, "Online global loop closure detection for large-scale multi-session graph-based SLAM," in *Proc. IEEE/RSJ Int. Conf. on Intelligent Robots and Systems*, 2014, pp. 2661–2666.

[6] D.-H. Kim and J.-H. Kim, "Effective background model-based RGB-D dense visual odometry in a dynamic environment," *IEEE Trans. Robotics*, vol. 32, no. 6, pp. 1565–1573, 2016.

[7] S. Li and D. Lee, "RGB-D SLAM in dynamic environments using

static point weighting," *IEEE Robotics and Automation Letters*, vol. 2, no. 4, pp. 2263–2270, 2017.

[8] Y. Sun, M. Liu, and M. Q.-H. Meng, "Improving RGB-D SLAM in dynamic environments: A motion removal approach," *Robotics and Autonomous Systems*, vol. 89, pp. 110–122, 2017.

[9] C. Kerl, J. Sturm, and D. Cremers, "Dense visual SLAM for RGB-D cameras," in *Proc. IEEE/RSJ Int. Conf. Intelligent Robots and Systems*, 2013, pp. 2100–2106.

[10] Y. Sun, M. Liu, and M. Q.-H. Meng, "Motion removal for reliable rgb-d slam in dynamic environments," *Robotics and Autonomous Systems*, vol. 108, pp. 115–128, 2018.

[11] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers, "A benchmark for the evaluation of RGB-D SLAM systems," in *Proc. IEEE/RSJ Int. Conf. Intelligent Robots and Systems*, Oct. 2012.

[12] R. F. Salas-Moreno, R. A. Newcombe, H. Strasdat, P. H. Kelly, and A. J. Davison, "SLAM++: Simultaneous localisation and mapping at the level of objects," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2013, pp. 1352–1359.

[13] J. McCormac, A. Handa, A. Davison, and S. Leutenegger, "SemanticFusion: Dense 3D semantic mapping with convolutional neural networks," in *Proc. IEEE Int. Conf. Robotics and Automation*, 2017, pp. 4628–4635.

[14] H. Noh, S. Hong, and B. Han, "Learning deconvolution network for semantic segmentation," in *Proc. IEEE Int. Conf. Computer Vision*, 2015, pp. 1520–1528.

[15] J. McCormac, R. Clark, M. Bloesch, A. Davison, and S. Leutenegger, "Fusion++: Volumetric object-level SLAM," in *2018 international conference on 3D vision (3DV)*. IEEE, 2018, pp. 32–41.

[16] B. Xu, W. Li, D. Tzoumanikas, M. Bloesch, A. Davison, and S. Leutenegger, "Mid-fusion: Octree-based object-level multi-instance dynamic slam," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 5231–5237.

[17] B. Bescos, J. M. Fácil, J. Civera, and J. Neira, "DynaSLAM: Tracking, mapping, and inpainting in dynamic scenes," *IEEE Robotics and Automation Letters*, vol. 3, no. 4, pp. 4076–4083, 2018.

[18] D. Bolya, C. Zhou, F. Xiao, and Y. J. Lee, "Yolact: Real-time instance segmentation," in *ICCV*, 2019.

[19] —, "Yolact++: Better real-time instance segmentation," 2019.

[20] Y. Li, H. Qi, J. Dai, X. Ji, and Y. Wei, "Fully convolutional instance-aware semantic segmentation," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2017, pp. 2359–2367.