

# NGD-SLAM: Towards Real-Time SLAM for Dynamic Environments without GPU

Yuhao Zhang  
yuhao.zhang-7@student.manchester.ac.uk

Department of Computer Science  
University of Manchester, UK

## Abstract

Accurate and robust camera tracking in dynamic environments presents a significant challenge for visual SLAM (Simultaneous Localization and Mapping). Recent progress in this field often involves the use of deep learning techniques to generate mask for dynamic objects, which usually require GPUs to operate in real-time (30 fps). Therefore, this paper proposes a novel visual SLAM system for dynamic environments that obtains real-time performance on CPU by incorporating a mask prediction mechanism, which allows the deep learning method and the camera tracking to run entirely in parallel at different frequencies such that neither waits for the result from the other. Based on this, it further introduces a dual-stage optical flow tracking approach and employs a hybrid usage of optical flow and ORB features, which significantly enhance the efficiency and robustness of the system. Compared with state-of-the-art methods, this system maintains high localization accuracy in dynamic environments while achieving a tracking frame rate of 56 fps on a single laptop CPU without any hardware acceleration, thus proving that deep learning methods are still feasible for dynamic SLAM even without GPU support. Based on the available information, this is the first SLAM system to achieve this.

Project page: <https://github.com/yuhaozhang7/NGD-SLAM>.

## Introduction

Simultaneous Localization and Mapping (SLAM) refers to the technology that leverages sensor's input data to achieve centimeter-level localization of the sensor and to reconstruct surroundings simultaneously. This makes it a fundamental element in areas such as robotics and augmented reality.

A key focus in SLAM research is visual SLAM, which refers to using cameras for localization and mapping. Although numerous visual SLAM algorithms [1, 2, 3, 4, 5] have demonstrated high tracking accuracy across various scenarios, most of them are designed under the assumption of a static environment. Consequently, camera tracking through these methods typically encounter drift in dynamic settings, which leads to a notable reduction in localization accuracy.

With advances in computer vision and deep learning, many SLAM algorithms incorporate deep learning models to identify and filter out potential dynamic objects in the input images based on semantic information [6], which results in remarkable improvement in tracking accuracy under dynamic environments. However, employing deep neural networks often leads to a substantial increase in computational requirements and longer tracking times.

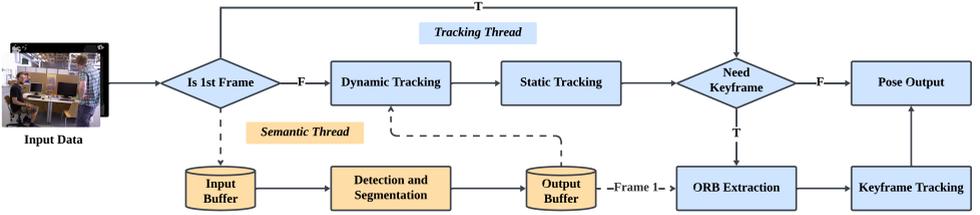


Figure 1: Tracking Pipeline of NGD-SLAM. Based on a mask prediction mechanism, the dynamic tracking process focuses on using optical flow to quickly track dynamic objects identified in the semantic thread and generating a pixel-wise mask to exclude dynamic features. The optical flow method is also employed in the static tracking process to efficiently track the remaining static features across non-keyframes, while ORB feature matching is maintained to establish a strong data association among keyframes.

This is because camera tracking needs to wait for the mask of dynamic objects provided by the deep learning model for each input frame [2]. Consequently, the emphasis in SLAM for dynamic environments has increasingly turned towards enhancing efficiency. Nevertheless, current approaches aiming for real-time performance either yield limited efficiency improvements or entail compromises in localization accuracy, and the reliance on GPU still remains.

In response, this paper propose NGD-SLAM (**No GPU Dynamic SLAM**), a real-time SLAM system designed for dynamic environments that operates without GPU support. This system is built upon ORB-SLAM3 framework [3] and includes an additional semantic thread for dynamic object detection using YOLO [20]. One core innovation of this paper is the tracking pipeline shown in Figure 1. Rather than focusing on deploying lighter deep learning models [11, 25, 28] or using them less frequently during the tracking process [12, 15, 23], it approach the challenge of efficiency as a task of utilizing past information to quickly generate the desired current outcomes. Specifically, it introduces a straightforward yet highly effective and *framework-independent mask prediction mechanism* (Figure 2) that utilizes previous segmentation result from the semantic thread to predict the mask of dynamic objects in the current frame. This mechanism can be integrated into any visual SLAM framework to enable concurrent operation of camera tracking and dynamic object segmentation, such that neither process needs to wait for the other’s output. To further enhance the robustness and efficiency, the system employs a **dual-stage tracking** that uses the optical flow method to continuously track identified dynamic and static features across images, along with a **hybrid strategy** that combines the strengths of optical flow and ORB features for camera tracking.

Through comprehensive experiments, NGD-SLAM demonstrates localization accuracy similar to that of state-of-the-art methods and stands as the most accurate one among all benchmarked real-time methods. Notably, it achieves a tracking frame rate of 56 fps on a single laptop CPU without any hardware acceleration.

## 2 Related Work

To address the challenges in dynamic environment, identifying pixels or keypoints originating from dynamic objects is crucial. Initial studies, such as [24], utilizes epipolar geometry constraints and a Bayesian approach to identify moving object. Klappstein et al. [13] in-

investigate detection of moving objects using motion metrics derived from optical flow. Using these traditional approaches can preserve the real-time nature of SLAM, but they also exhibit limited adaptability across various scenarios and result in reduced accuracy in localization.

With advancements in deep learning, more sophisticated methods have been developed for detecting or segmenting dynamic instances across input images using semantic information. Building on ORB-SLAM2 [14], DS-SLAM [26] improves localization accuracy in dynamic environment by incorporating SegNet [10] for dynamic instances segmentation and also maintains the epipolar geometry constraints for moving consistency check. Similarly, DynaSLAM [9] employs Mask R-CNN [10] and multi-view geometry to eliminate dynamic ORB features. Instead of using segmentation model, Detect-SLAM [28] leverages YOLO model [20] for dynamic object detection, along with moving probability propagation to remove dynamic objects. More recently, CFP-SLAM [11] proposes a coarse-to-fine static probability approach to enhance localization accuracy by calculating and updating static probabilities as weights for pose optimization.

Although deep learning-driven methods demonstrate remarkable improvement in accuracy, they typically rely on powerful GPUs to achieve (near) real-time performance due to the computationally intensive nature of deep neural networks. One strategy to boost the efficiency is incorporating a separate semantic thread for deep learning model. Based on this idea, common approaches involve parallelizing the execution of feature extraction and dynamic object identification [11, 26], but they only yield limited improvements in efficiency as the subsequent tracking process must await the output from the deep neural network. On the other hand, some algorithms apply the deep learning method only to keyframes [11, 15, 28], which indeed greatly reduce the average tracking time but result in compromised accuracy.

While NGD-SLAM also includes a separate semantic thread for deep neural network, it incorporates a mask prediction mechanism that allows the tracking thread to operate completely independently without waiting for the result from the semantic thread. Moreover, distinct from the methods that apply optical flow or scene flow for differentiating dynamic points from static ones [24, 26, 27], NGD-SLAM leverages optical flow method only for sparse features tracking and applies it across two different tasks designed to enhance the system’s accuracy and efficiency.

## 3 Methodology

### 3.1 Tracking Pipeline

Figure 1 illustrates the tracking pipeline of NGD-SLAM. Initially, the image data of the first frame is forwarded to the semantic thread, and this frame is set as a keyframe. Its ORB features are extracted with a mask from the semantic thread (Section 3.2) to filter out dynamic features; only the first frame will wait for the mask given by the semantic thread. The system then initializes the map as part of the keyframe tracking. For subsequent frames, the system employs dual-stage tracking (Section 3.4), which involves tracking potential dynamic objects to generate mask for them based on the mask prediction mechanism (Section 3.3) and tracking the remaining static points from the last frame to estimate the camera’s pose. The decision to set a new keyframe is then made. If a new keyframe is required, the system extracts static ORB features from the current frame by filtering out features within the mask generated by dynamic tracking process, and continues with the keyframe tracking process using ORB features (Section 3.5).

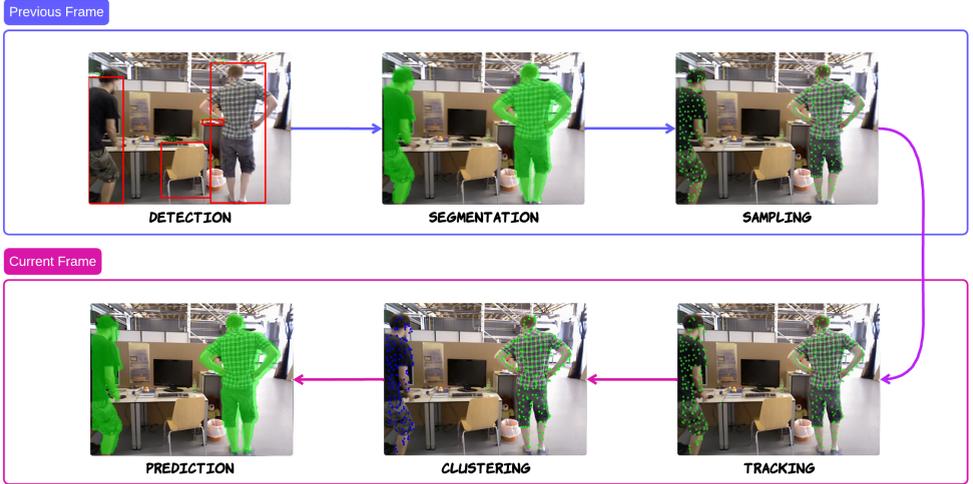


Figure 2: Mask Prediction Pipeline. The detection and segmentation are processed inside semantic thread, while tracking thread uses this result to predict mask in current frame.

## 3.2 Semantic Thread

NGD-SLAM enhances real-time performance by adding a semantic thread that detects dynamic objects using the YOLO-fastest model [6, 20]. This thread works with an input and output buffer, each holding one frame, in which a new input or output frame will replace the existing one in the buffer. Frames that passed from the tracking thread to the input buffer, are processed by the network, and then move to the output buffer for the tracking thread to access. This setup ensures that the semantic thread processes the most recent input frame and the tracking thread does not have to wait for the result from the semantic thread, instead taking the latest output directly from the buffer. Although there is a temporal mismatch because two threads run at different frequencies, meaning the result taken from the output buffer corresponds to one of the previous frames, NGD-SLAM compensates for this by using this past result to predict mask for dynamic object in the current frame (Section 3.3).

## 3.3 Mask Prediction

One important module of NGD-SLAM is the *framework-independent* mask prediction mechanism, which also serves as a foundational component of dynamic tracking (Section 3.4). Compared to the forward process of a deep neural network, this prediction is much faster, which takes only a few milliseconds to process once on a CPU. It can be directly adapted to any visual SLAM framework and maintains its real-time performance. Figure 2 shows the pipeline of the mask prediction, which consists of six components.

**Detection.** The YOLO-fastest model, employed in the semantic thread, is utilized to detect multiple objects in the frame and identify potential dynamic ones based on semantic information. It is integrated using OpenCV DNN for C++ compatibility. NGD-SLAM employs this model for its efficiency but also allows for easy substitutions with other models for object detection or segmentation as needed.

**Segmentation.** Using an object detection model is more efficient than a segmentation network but it only provides bounding boxes for detected objects. To improve this, NGD-SLAM uses depth information to segment the objects within these boxes. After identifying dynamic objects such as people, it then sorts pixel depths in each bounding box, using the median as the estimated depth for that detected object. The bounding box is expanded and all pixels close to this depth value are masked. This approach is preferred over clustering methods for its computational efficiency. Subsequently, the connected component labeling algorithm [9] is used to exclude mask of objects that share a similar depth but are not part of the main object, and a dilation process further refines the mask before it is passed into the output buffer.

**Sampling.** Given the mask of identified dynamic objects, dynamic points are extracted within the mask. The image is divided into cells with size of  $15 \times 15$ . Within each cell, pixels that are masked and meet the FAST keypoint criteria [20] are identified, and the one with the highest threshold value is selected. If no suitable FAST keypoints present, a random pixel is chosen within that cell. This grid structure allows for uniform extraction of dynamic keypoints, and the tracking stability is ensured due to the distinctiveness of keypoints and their association with high-gradient regions.

**Tracking.** The system utilizes the Lucas-Kanade optical flow method to track sampled dynamic points, which assumes that the grayscale value of a keypoint remains constant across different frames [9, 16]. Mathematically, this is represented as:

$$I(x + dx, y + dy, t + dt) = I(x, y, t)$$

where  $I(x, y, t)$  is the grayscale intensity value of a pixel at coordinates  $(x, y)$  with time  $t$ . This equation can be further expanded using a Taylor series as:

$$I(x + dx, y + dy, t + dt) \approx I(x, y, t) + \frac{\partial I}{\partial x} \cdot dx + \frac{\partial I}{\partial y} \cdot dy + \frac{\partial I}{\partial t} \cdot dt$$

Subsequently, assuming that pixels move uniformly within a small window, an optimization problem can be formulated, which uses the pixel gradients within the window to estimate pixel motion  $(dx, dy)$  by minimizing the total grayscale error of matched pixels.

**Clustering.** For tracked dynamic points in the current frame, the DBSCAN (Density-Based Spatial Clustering of Applications with Noise) algorithm [8] is applied to cluster points that are close to each other while distinguishing those significantly distant from any cluster as noise. Employing this method allows for the effective distinction between tracked points originating from distinct dynamic objects. Additionally, by considering the 2D position and depth of each tracked point, the algorithm effectively excludes outliers, such as points that significantly deviate from clusters or those that seem close in the 2D image but have notable depth differences.

**Prediction.** After deriving the clusters, the system becomes aware of the dynamic entity's location within the current frame. Therefore, it initially generates a rough mask, typically a rectangle, for each cluster to cover all points within it. Subsequently, it uses the depth information of points in each cluster to refine the corresponding mask into precise shape of the dynamic object. As shown in Figure 2, the mask predicted for the current frame is remarkably precise.

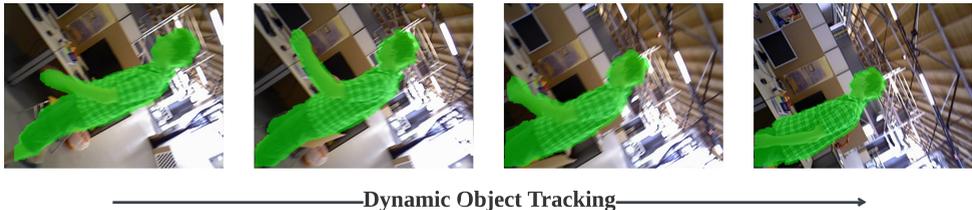


Figure 3: Dynamic Object Tracking. The system continuously tracks the dynamic object over 30 frames, even when the deep learning model completely fails.

### 3.4 Dual-Stage Tracking

**Dynamic Tracking.** Although the mask prediction mechanism is already kind of dynamic tracking, its dependence on YOLO for object detection can sometimes result in dynamic object tracking failures. For instance, YOLO may struggle to detect objects when the camera is rotated. Therefore, when the semantic thread fails to provide a mask for tracking, NGD-SLAM implements a straightforward solution: it uses the predicted mask from the previous frame to predict the mask for the current frame. This method allows for the continuous tracking of dynamic points until the dynamic object is no longer visible to the camera, even if the deep neural network experiences issues. By continually sampling new points from the previous mask for tracking, rather than tracking the same points across multiple frames, it ensures that the keypoints remain up-to-date and enhances the robustness of the tracking (Figure 3). While alternative approaches involve deriving the *roll* value from the quaternion expression  $(w, x, y, z)$  of camera pose’s rotational part using:



Figure 4: Rotate the image as alternative solution.

$$roll = \arctan 2(2 \times (w \times x + y \times z), 1 - 2 \times (x^2 + y^2)) \times \frac{180}{\pi}$$

and use this value to rotate the image for detection (Figure 4), the purpose of dynamic tracking is to provide robust tracking in general cases, as camera rotation is not the only factor that can lead to detection failure.

**Static Tracking.** In contrast to the original ORB-SLAM3, which depends on ORB feature matching for every input frame, NGD-SLAM uses the same optical flow method [16] described in Section 3.3 to track the static keypoints from the last frame. Geometric constraints are formed by linking 3D map points from the last frame to tracked 2D keypoints in the current frame, and the initial pose of the current frame is calculated based on the motion model below [8]:

$$T_t = (T_{t-1}T_{t-2}^{-1})T_{t-1}$$

where  $T_{t-1}$  is the pose of the last frame. Based on this information, the RANSAC (Random Sample Consensus) algorithm [7] is applied to iteratively estimate the actual pose and filter out outliers.

This is designed based on two considerations. Firstly, the mask prediction mechanism typically fails when a dynamic object is presented in the current frame but doesn’t appear

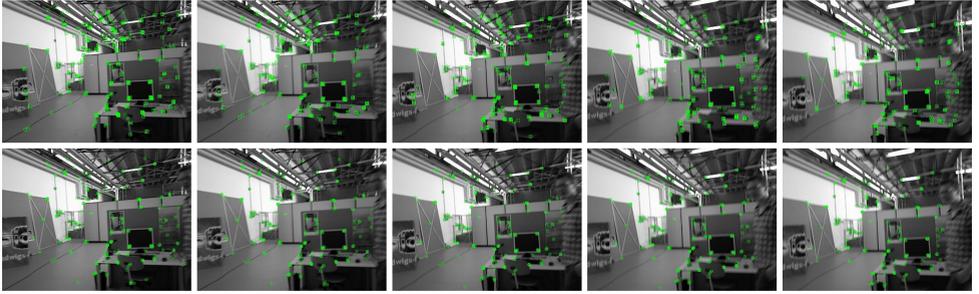


Figure 5: Static Points Tracking. These are frames where mask prediction fails because a dynamic object newly appears. (Top) Extracting and matching ORB features for each frame leads to matches with dynamic points. (Bottom) Using the optical flow method to track the same static points across multiple frames.

in the frame used for prediction. Thus, extracting new ORB features and matching based on descriptor value differences for each frame can result in mismatches or matches between dynamic points (Figure 5). In contrast, the static tracking process, which track the same static points across multiple frames based on image gradient and outliers are filtered out for future tracking, tends to be more robust. Although new ORB features will be extracted for keyframe tracking, this occurs when there are not enough tracked static points, indicating that a dynamic object has appeared for a while across multiple frames. Secondly, the static tracking process skips the computationally expensive ORB feature extraction for non-keyframes, making the tracking process much more efficient.

### 3.5 Keyframe Tracking

ORB features are retained for keyframes to establish connections between them. It allows the maintenance of essential ORB-SLAM3 components like co-visible graphs and local maps, which are crucial for robust tracking and efficient data retrieval. It also serves as a backup tracking method when static tracking underperforms, as it is more robust to illumination changes and motion blur. Keyframe tracking is achieved by using the pose estimated in the static tracking phase to project the map points from the last keyframe into the current frame to get matches and perform optimizations. Similar to original ORB-SLAM3 [3], new map points are added, and the frame is passed to the local mapping thread for further optimization.

## 4 Experiments

### 4.1 Experimental Setup

The state-of-the-art algorithms and several other baseline methods for SLAM in dynamic environments are chosen for comparison, including DynaSLAM [2], DS-SLAM [27], CFP-SLAM [11], RDS-SLAM [15], and TeteSLAM [12]. NGD-SLAM without dual-stage tracking (keeping mask prediction only) is also evaluated for ablation study.

Four sequences from the TUM dataset [23] are utilized, which captures scenarios of two individuals moving around a desk, one wearing a plaid shirt that is rich in texture. These se-

	DynaSLAM	DS-SLAM	CFP-SLAM	RDS-SLAM	TeteSLAM	NGD-SLAM	w/o DST
f3/w_xyz	0.015 (2)	0.025 (5)	<b>0.014</b> (1)	0.057 (6)	0.019 (4)	0.015 (2)	0.015 -
f3/w_rpy	0.036 (2)	0.444 (6)	0.037 (3)	0.160 (5)	0.037 (3)	<b>0.034</b> (1)	0.040 ↑
f3/w_half	0.027 (3)	0.030 (5)	<b>0.024</b> (1)	0.081 (6)	0.029 (4)	<b>0.024</b> (1)	0.028 ↑
f3/w_static	<b>0.007</b> (1)	0.008 (4)	<b>0.007</b> (1)	0.008 (4)	0.011 (6)	<b>0.007</b> (1)	0.007 -

Table 1: Comparison of RMSE of absolute trajectory error (m).

	DynaSLAM	DS-SLAM	CFP-SLAM	RDS-SLAM	TeteSLAM	NGD-SLAM	w/o DST
f3/w_xyz	0.021 (3)	0.033 (5)	<b>0.019</b> (1)	0.042 (6)	0.023 (4)	0.020 (2)	0.020 -
f3/w_rpy	0.045 (2)	0.150 (6)	0.050 (4)	0.132 (5)	0.047 (3)	<b>0.044</b> (1)	0.052 ↑
f3/w_half	0.028 (3)	0.030 (4)	0.026 (2)	0.048 (6)	0.042 (5)	<b>0.025</b> (1)	0.029 ↑
f3/w_static	<b>0.008</b> (1)	0.010 (4)	0.009 (2)	0.022 (6)	0.011 (5)	0.009 (2)	0.010 ↑

Table 2: Comparison of RMSE of relative pose error in translation (m/s).

	DynaSLAM	DS-SLAM	CFP-SLAM	RDS-SLAM	TeteSLAM	NGD-SLAM	w/o DST
f3/w_xyz	<b>0.452</b> (1)	0.826 (5)	0.602 (3)	0.922 (6)	0.636 (4)	0.470 (2)	0.482 ↑
f3/w_rpy	0.902 (2)	3.004 (5)	1.108 (4)	13.169 (6)	1.058 (3)	<b>0.889</b> (1)	1.012 ↑
f3/w_half	0.737 (2)	0.814 (4)	0.757 (3)	1.883 (6)	0.965 (5)	<b>0.695</b> (1)	0.779 ↑
f3/w_static	0.258 (2)	0.269 (4)	<b>0.253</b> (1)	0.494 (6)	0.287 (5)	0.262 (3)	0.283 ↑

Table 3: Comparison of RMSE of relative pose error in rotation ( $^{\circ}$ /s).

	DynaSLAM	DS-SLAM	CFP-SLAM	RDS-SLAM	TeteSLAM	NGD-SLAM	w/o DST
Ratio	15.94	2.07	2.21	1.30	1.06	0.89	1.51
Real-Time	<b>X</b>	<b>X</b>	<b>X</b>	✓	✓	✓	✓
Device	CPU + GPU	CPU	CPU				

Table 4: Comparison of average tracking time per frame.

quences exhibit diverse camera movements, including translations along xyz axes, rotational movements (significant roll, pitch, yaw changes), a trajectory resembling a half-sphere, and one that is nearly stationary. For accuracy, ATE (Absolute Trajectory Error) and RPE (Relative Trajectory Error) are chosen as evaluation metrics, with the interval for RPE measurements set to 1 second (30 frames). For efficiency, it measures the time taken in milliseconds to process a given function.

All experiments are conducted on a laptop equipped with an AMD Ryzen 7 4800H CPU and an NVIDIA GeForce RTX 2060 GPU.

## 4.2 Experimental Results

**Accuracy.** Table 1 presents a comparison of NGD-SLAM against various baseline methods regarding RMSE of ATE. The minimal error value is highlighted, with the corresponding rank indicated alongside (except for NGD-SLAM without DST, which indicates whether the error is increased or decreased compared to its base implementation). It is clear that DynaSLAM, CFP-SLAM, and NGD-SLAM demonstrate high precision throughout all sequences. However, DynaSLAM generally gives slightly higher error, and CFP-SLAM is less impressive in the *f3/w\_rpy* sequence because the epipolar constraints they employ do not work well in the scenario with rotation [10]. Moreover, Tables 2 and 3 display the RPE RMSE comparison, where NGD-SLAM continues to uphold the state-of-the-art outcomes.

	Dynamic Tracking	Static Tracking	Keyframe Tracking	Total Tracking
Processing Time	9.74	3.18	19.92	32.84

Table 5: Average tracking time per keyframe of NGD-SLAM (ms)

It is also important to note that the dual-stage tracking plays a crucial role here, as both ATE and RPE of the system generally increase when dual-stage tracking is disabled. This is particularly evident in more challenging sequences like  $f3/w\_rpy$  and  $f3/w\_half$ .

**Efficiency.** Table 4 presents a comparison of the average processing time per frame across these dynamic sequences. As some algorithms are not open sourced and are tested on various devices, a fair comparison is ensured by displaying the ratio of the average processing time per frame of the algorithm to the average processing time of the original ORB-SLAM2/3 as reported in their source literature. The average tracking time of ORB-SLAM2/3 is around 20 milliseconds; therefore, a value within 1.65 can be considered real-time performance.

While DynaSLAM and CFP-SLAM demonstrate high accuracy, neither achieves real-time performance even with GPU support. RDS-SLAM and TeteSLAM run in real-time but compromise accuracy as they applying deep learning models only to keyframes [14, 15]. In contrast, both NGD-SLAM and it without DST (keeping mask prediction only) achieve real-time performance and only require a CPU.

**Real-Time Analysis.** The average processing time per frame for the system is 17.88 milliseconds (56 fps), attributed to the efficiency of optical flow tracking in non-keyframes. Since only the tracking of a keyframe involves all components of NGD-SLAM, the experiment is further segmented to evaluate the average time required for tracking a keyframe using static, dynamic, and keyframe tracking methods, as shown in Table 5. It is evident that static tracking with optical flow is highly efficient, and although dynamic tracking requires slightly more time, the combined tracking process falls within 33.3 milliseconds (30 fps), which still meets the criteria for real-time performance. Notably, while YOLO detection takes approximately 30 to 40 milliseconds per frame, it is not factored into this tracking time since the tracking and semantic threads operate concurrently at different frequencies.

## 5 Conclusion

This paper introduces a real-time visual SLAM system designed for dynamic environments that operates on a CPU. It incorporates a *framework-independent* mask prediction mechanism to mitigate the inefficiencies caused by the use of deep learning models while preserving their high accuracy in dynamic object identification. Additionally, a dual-stage tracking method is developed to complement the limitations of the mask prediction mechanism and to further enhance the system’s efficiency. Experimental evaluations on dynamic environment datasets and comparisons with baseline methods reveal that NGD-SLAM achieves accuracy on par with top-performing algorithms and obtains real-time tracking on a single laptop CPU without any hardware acceleration. This progress underscores the potential of deep learning-based methods in enhancing the effectiveness of SLAM systems in dynamic environments, even without GPU support.

## Acknowledgement

Huge thanks to Professor Mikel Luján (University of Manchester) for valuable discussions.

## References

- [1] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 39(12):2481–2495, 2017.
- [2] Berta Bescos, José M Fácil, Javier Civera, and José Neira. Dynaslam: Tracking, mapping, and inpainting in dynamic scenes. *IEEE Robotics and Automation Letters*, 3(4): 4076–4083, 2018.
- [3] Carlos Campos, Richard Elvira, Juan J Gómez Rodríguez, José MM Montiel, and Juan D Tardós. Orb-slam3: An accurate open-source library for visual, visual–inertial, and multimap slam. *IEEE Transactions on Robotics*, 37(6):1874–1890, 2021.
- [4] Luigi Di Stefano and Andrea Bulgarelli. A simple and efficient connected components labeling algorithm. In *Proceedings 10th international conference on image analysis and processing*, pages 322–327. IEEE, 1999.
- [5] Dogquiqui. Yolo-fastest: yolo-fastest-v1.1.0, 2021. URL <https://doi.org/10.5281/zenodo.5131532>.
- [6] Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. A density-based algorithm for discovering clusters in large spatial databases with noise. In *kdd*, volume 96, pages 226–231, 1996.
- [7] Martin A Fischler and Robert C Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 24(6):381–395, 1981.
- [8] Christian Forster, Matia Pizzoli, and Davide Scaramuzza. Svo: Fast semi-direct monocular visual odometry. In *2014 IEEE international conference on robotics and automation (ICRA)*, pages 15–22. IEEE, 2014.
- [9] Xiang Gao, Tao Zhang, Yi Liu, and Qinrui Yan. *14 Lectures on Visual SLAM: From Theory to Practice*. Publishing House of Electronics Industry, 2017.
- [10] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017.
- [11] Xinggang Hu, Yunzhou Zhang, Zhenzhong Cao, Rong Ma, Yanmin Wu, Zhiqiang Deng, and Wenkai Sun. Cfp-slam: A real-time visual slam based on coarse-to-fine probability in dynamic environments. In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 4399–4406. IEEE, 2022.
- [12] Tete Ji, Chen Wang, and Lihua Xie. Towards real-time semantic rgb-d slam in dynamic environments. In *2021 IEEE international conference on robotics and automation (ICRA)*, pages 11175–11181. IEEE, 2021.

- [13] Jens Klappstein, Tobi Vaudrey, Clemens Rabe, Andreas Wedel, and Reinhard Klette. Moving object segmentation using optical flow and depth information. In *Advances in Image and Video Technology: Third Pacific Rim Symposium, PSIVT 2009, Tokyo, Japan, January 13-16, 2009. Proceedings 3*, pages 611–623. Springer, 2009.
- [14] Abhijit Kundu, K Madhava Krishna, and Jayanthi Sivaswamy. Moving object detection by multi-view geometric techniques from a single camera mounted robot. In *2009 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 4306–4312. IEEE, 2009.
- [15] Yubao Liu and Jun Miura. Rds-slam: Real-time dynamic slam using semantic segmentation methods. *Ieee Access*, 9:23772–23785, 2021.
- [16] Bruce D Lucas and Takeo Kanade. An iterative image registration technique with an application to stereo vision. In *IJCAI’81: 7th international joint conference on Artificial intelligence*, volume 2, pages 674–679, 1981.
- [17] Raul Mur-Artal and Juan D Tardós. Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras. *IEEE transactions on robotics*, 33(5):1255–1262, 2017.
- [18] Raul Mur-Artal, Jose Maria Martinez Montiel, and Juan D Tardos. Orb-slam: a versatile and accurate monocular slam system. *IEEE transactions on robotics*, 31(5):1147–1163, 2015.
- [19] Tong Qin, Peiliang Li, and Shaojie Shen. Vins-mono: A robust and versatile monocular visual-inertial state estimator. *IEEE Transactions on Robotics*, 34(4):1004–1020, 2018.
- [20] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016.
- [21] Edward Rosten and Tom Drummond. Fusing points and lines for high performance tracking. In *Tenth IEEE International Conference on Computer Vision (ICCV’05) Volume 1*, volume 2, pages 1508–1515. Ieee, 2005.
- [22] Muhamad Risqi U Saputra, Andrew Markham, and Niki Trigoni. Visual slam and structure from motion in dynamic environments: A survey. *ACM Computing Surveys (CSUR)*, 51(2):1–36, 2018.
- [23] Jürgen Sturm, Nikolas Engelhard, Felix Endres, Wolfram Burgard, and Daniel Cremers. A benchmark for the evaluation of rgb-d slam systems. In *2012 IEEE/RSJ international conference on intelligent robots and systems*, pages 573–580. IEEE, 2012.
- [24] Youbing Wang and Shoudong Huang. Towards dense moving object segmentation based robust dense rgb-d slam in dynamic scenarios. In *2014 13th International Conference on Control Automation Robotics & Vision (ICARCV)*, pages 1841–1846. IEEE, 2014.
- [25] Wenxin Wu, Liang Guo, Hongli Gao, Zhichao You, Yuekai Liu, and Zhiqiang Chen. Yolo-slam: A semantic slam system towards dynamic environment with geometric constraint. *Neural Computing and Applications*, pages 1–16, 2022.

- 
- [26] Chao Yu, Zuxin Liu, Xin-Jun Liu, Fugui Xie, Yi Yang, Qi Wei, and Qiao Fei. Ds-slam: A semantic visual slam towards dynamic environments. In *2018 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages 1168–1174. IEEE, 2018.
- [27] Tianwei Zhang, Huayan Zhang, Yang Li, Yoshihiko Nakamura, and Lei Zhang. Flow-fusion: Dynamic dense rgb-d slam based on optical flow. In *2020 IEEE international conference on robotics and automation (ICRA)*, pages 7322–7328. IEEE, 2020.
- [28] Fangwei Zhong, Sheng Wang, Ziqi Zhang, and Yizhou Wang. Detect-slam: Making object detection and slam mutually beneficial. In *2018 IEEE winter conference on applications of computer vision (WACV)*, pages 1001–1010. IEEE, 2018.