# Semantic Nearest Neighbor Fields for Monocular Edge Visual-Odometry

Xiaolong Wu, Assia Benbihi, Antoine Richard, and Cédric Pradalier

*Abstract*— Recent advances in deep learning for edge detection and segmentation opens up a new path for semantic-edge-based ego-motion estimation. In this work, we propose a robust monocular visual odometry (VO) framework using category-aware semantic edges. It can reconstruct large-scale semantic maps in challenging outdoor environments. The core of our approach is a semantic nearest neighbor field that facilitates a robust data association of edges across frames using semantics. This significantly enlarges the convergence radius during tracking phases. The proposed edge registration method can be easily integrated into direct VO frameworks to estimate photometrically, geometrically, and semantically consistent camera motions. Different types of edges are evaluated and extensive experiments demonstrate that our proposed system outperforms state-of-art indirect, direct, and semantic monocular VO systems.

## I. INTRODUCTION

In recent decades, monocular Visual Odometry (VO) and Simultaneous Localization and Mapping (SLAM) systems have shown their full potential to assist various robotic applications in outdoor operations - from autonomous driving in urban scenes to environmental monitoring in natural environments. Among these algorithms, indirect methods [1] [2] [3] are the *de facto* standards since visual features provide considerable robustness to both photometric noise and geometric distortion in images. Recent works have shown that direct methods [4] [5] [6] [7] could provide more accurate and robust motion estimation. However, they present a much smaller convergence basin compared with indirect methods because of loose data association.

Edge-based ego-motion estimation has also gained significant attention for its robustness against illumination changes, motion blur, and occlusion. Edge VO and SLAM can be seen as a crossover of indirect and direct principles. Specifically, edges are binary features extracted from raw images, but edge registration is performed using iterative-closest-point (ICP) based direct alignment. Motion estimation using edges is particularly attractive for outdoor applications for its illumination and convergence robustness. Still, standard edge detection methods (e.g. Canny) used in edge VO and SLAM provide edges with poor repeatability in outdoor settings, which hinders their usage. Recent developments of learning-based edge detection opens up a new path to improve the stability of edge detection for outdoor edge VO and SLAM.

In addition to feature and edge learning, semantic information can also boost motion estimation. [8] integrates a semantic reprojection error into state-of-art indirect and direct VO systems to improve the tracking robustness and accuracy. However, it may be not suitable for edge-based motion estimation because of ambiguous semantic labels of edges around semantic boundaries.
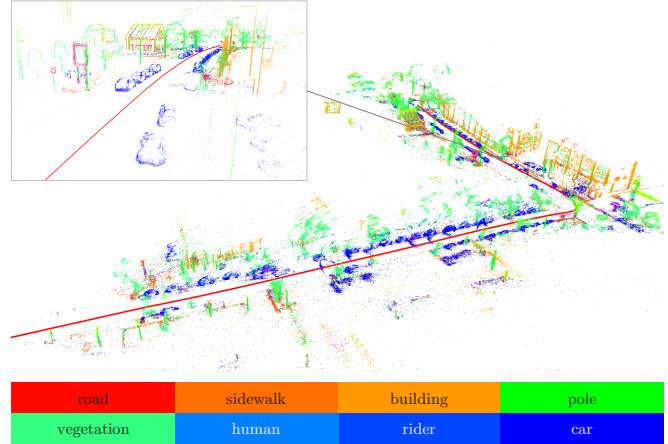


Fig. 1: 3D semantic map with our semantic edge VO system. Semantic edges are learned with CaseNet [9].

In this work, we present an edge-based VO framework that can track and reconstruct semantic edges across frames for outdoor robotic applications. We show that our proposed semantic nearest neighbor fields (SNNFs) offer several advantages over existing edge-based VO algorithms. The learned semantic edges further improve the accuracy and robustness, especially in outdoor environments in Fig. 1. We investigate the influence of several edge detection methods on motion estimation. We conduct extensive evaluation on KITTI, a public autonomous driving dataset, and measure tracking accuracy, robustness and runtime performances. Experimental results show that our proposed system outperforms state-of-art monocular direct, indirect, and semantic VO systems in an outdoor setting. Our main contributions are:

- A novel SNNFs strategy using semantics to improve the efficiency and robustness of ICP-based edge registration.
- A monocular semantic edge VO system that integrates image gradient, edges and semantic information to compute photometrically, geometrically, and semantically consistent motion. It is capable of reconstructing large-scale semantic maps.
- A evaluation study of several edge detection methods for outdoor ego-motion estimation. To our knowledge, this is the first paper evaluating edge-based VO system in an outdoor environment.

## II. RELATED WORK

VO and SLAM are two main solutions for camera tracking and environmental mapping, and can be divided into indirect

or direct methods. Indirect methods [1] [3] minimize the reprojection error between features in two consecutive images. In contrast to feature-based approaches, direct methods [10] [6] [7] minimize the photometric error between all the matching pixels in two successive images. The main drawback is that it presents a much smaller convergence basin than indirect methods because of its loose data association.

Edge-based ego-motion estimation relies on ICP-based optimization. [11] proposes an efficient 2D-3D edge registration framework for real-time motion estimation using distance transform (DT) [12]. However their objective function is neither differentiable nor negative. [13] relieves the first issue with a sub-gradient method. [14] solves the second and enables Gauss-Newton like optimization with a signed residual based on 2-D edge divergence minimization. [15] tackles both problems concurrently and substitutes DT with approximate nearest neighbor fields (ANNFs). The subsequent work [16] extends ANNFs into oriented nearest neighbor fields (ONNFs). It provides a finer data association strategy to reject outliers and enables more robust edge registration. Besides pure edge-based methods, [17] and [18] combine intensity-based photometric and edge-based geometric errors to improve tracking robustness and accuracy. However, performance may degrade with illumination changes contrary to semantic information that is more invariant.

With the recent advances in deep learning, semantic information has become relevant for motion estimation. [19] [20] use semantics to detect moving objects and alleviate their pixels weights in the objective function. Similarly, [21] masks the sky pixels. [8] integrates a semantic reprojection error into existing point-based indirect and direct motion estimation systems. The overall advantage of integrating semantics is to boost the tracking and mapping robustness and accuracy. We pursue these efforts and fuse the advantages from both semantic-SLAM and learned-edges-SLAM to make VO even more robust.

Previous edge-based SLAM methods rely on the standard Canny edge detector [22], but [23] shows that using machine-learned edges [24] [25] gives significant improvements compared to standard edges in indoor settings. Structured edges (SE) [24] generalizes random forests to general structured output spaces to learn edges on patch images. Despite its fast computation time, it relies on hand-crafted features and requires to map the ground-truth edges to low dimensional representation for the training to be scalable. These drawbacks are alleviated with the deep-learning based HED [25]: it processes raw images and outputs an edge probability map over the input in the form of a binary map. In this work, we extend edge-based SLAM and enhance them with semantic information to make them more robust. To do so, we use learned-edges together with their semantic labels as computed in [9] [26]. CaseNet [9] extends HED [25] to multi-label edges and outputs a distinct probability edge map for each semantic class. SEAL [26] continues the efforts of [9] and tackles the challenge of edge-alignment during the CNN training. Label noise in the human-annotated ground-truth edges leads CaseNet to learn thick edges which hides relevant

edges details. However, none of these works investigate the application of semantic edges in VO. In contrast, we provide a method to integrate any semantic edge into VO and run extensive experiments to compare the influence of edge and semantic learning on the VO performances.

## III. SEMANTIC EDGE VISUAL ODOMETRY

We first give a brief introduction on edge-based motion estimation and discuss the main limitation of existing work in Sec. III-A. We then introduce our proposed SNNFs and the corresponding edge registration algorithm in Sec. III-B. The implementation details are provided in Sec. III-C.

### A. Problem Formulation

Let us consider a reference frame, a gray-scale reference image $I_r : \Omega \to \mathbb{R}$ and an inverse depth map $D_r : \Omega \to \mathbb{R}^+$, where $\Omega \subset \mathbb{R}^2$ is the image domain. A 3D scene point $\mathbf{P} = (x, y, z)^T$ is parameterized by its inverse depth $d = z^{-1}$ in the reference frame instead of the conventional 3 unknowns. Thus, each pixel in the reference frame $\mathbf{p} = (u, v)^T \in \Omega$ can be back-projected into 3D world using the back-projection function $\pi^{-1}(\cdot)$ as:

$$\mathbf{P} = \pi^{-1}(\mathbf{p}, d) = \mathbf{K}^{-1}\bar{\mathbf{p}}/d \qquad (1)$$

where $\bar{\mathbf{p}} = (\mathbf{p}^T, 1)^T$ is the homogeneous coordinate of pixel coordinate and $\mathbf{K}$ is the pre-calibrated camera intrinsic matrix. Inversely, the 3D projective warp function $\pi(\cdot)$ can be expressed as:

$$\mathbf{p} = \pi(\mathbf{P}) = \bar{\mathbf{K}}\begin{bmatrix} x/z \\ y/z \end{bmatrix}^T \qquad (2)$$

where $\bar{\mathbf{K}}$ is the first two rows of intrinsic matrix $\mathbf{K}$.

Given arbitrary edge detector $E(\cdot)$, the group of edge pixels $\mathscr{E}_r$ extracted from the reference image can be expressed as:

$$\mathscr{E}_r = \{\mathbf{p}_r\} = E(I_r) \qquad (3)$$

The detected edge pixels $\mathscr{E}_r$ are subsequently projected to current frame $k$. The transformed edge pixels $\mathscr{E}_{kr}$ can be computed as:

$$\mathscr{E}_{kr} = \{\mathbf{p}_{kr}\} = \left\{\pi(\mathbf{R}_{kr}\pi^{-1}(\mathbf{p}_r, D_r(\mathbf{p}_r)) + \mathbf{t}_{kr})\right\} \qquad (4)$$

where $\mathbf{R}_{kr} \in SO(3)$ and $\mathbf{t}_{kr} \in \mathbb{R}^3$ are the 3D rigid body rotation and translation from reference frame to current frame $k$, respectively.

Let us define a function to find the nearest neighbor of the projected edge pixel $\mathbf{p}_{kr}$ in current frame $k$ using the Euclidean distance metric as:

$$n(\mathbf{p}_{kr}) = \underset{\mathbf{p}_k \in \mathscr{E}_k}{\arg\min} \|\mathbf{p}_k - \mathbf{p}_{kr}\| \qquad (5)$$

As the result, the total energy $E_{kr}^{\mathscr{E}}$ sums up all edge distance errors from the reference frame to the current frame $k$ expressed as:

$$E_{kr}^{\mathscr{E}} := \sum_{\mathbf{p}_r \in \mathscr{E}_r} w_{\mathbf{p}_r}^{\mathscr{E}} \|\mathbf{p}_{kr} - n(\mathbf{p}_{kr})\|_\gamma \qquad (6)$$

where $w_{\mathbf{p}_r}$ is the weight assigned for each edge pixel in the reference frame and $\|\cdot\|_\gamma$ is the Huber norm.
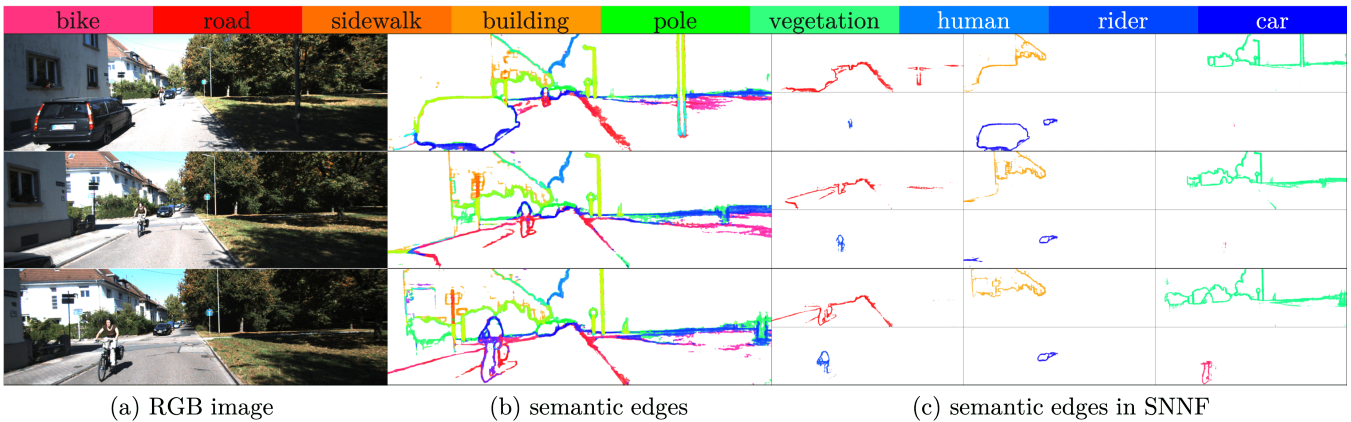
| bike | road | sidewalk | building | pole | vegetation | human | rider | car |

(a) RGB image        (b) semantic edges        (c) semantic edges in SNNF

Fig. 2: Example sequential RGB images, fused semantic edge maps, and their SNNFs.

To find the optimal camera transformation $\mathbf{R}$ and $\mathbf{t}$ using Eqn. 5, we use a 2D-3D ICP-based optimization [11] that alternates between finding approximate nearest neighbors and register the putative correspondences using an iteratively reweighted Gauss-Newton algorithm. Following the theory of optimization under unitary constraints [27], we optimize the energy function defined in Eqn. 5 on Lie-manifolds for better convergence.

### B. Semantic Nearest Neighbor Fields

ANNFs-base edge VO is extensively evaluated in [15] on an indoor dataset, and presents impressive performance. However, when it comes to outdoor environments, the results of both ANNFs and ONNFs deteriorate because of unrepeatable edges and large camera motions. This comes from the weak data association strategy of edge-based VO algorithms: it makes them sensitive to outliers and motion initialization. To solve this issue, we incorporate semantic information into existing edge VO frameworks and facilitate a robust data association of edges across frames.

The key idea behind SNNFs is to classify the extracted edges as shown in Fig.2. Then each group of edges can only be registered to their counterparts in subsequent frames. Specifically, for a dataset with $C$ classes, we generate $C$ edge-class probability maps. An edge pixel can belong to one or more class: for example, an edge pixel lying between a car and the road has both labels. Then, we compute distance fields with region growing algorithm on the seeded region of each map. Upon registration, each edge pixels is only registered with edges with the same semantic label instead of spatially nearest ones. This way, we avoid ambiguous associations and enlarge the convergence basin during registration: SNNFs constrains edge registration with semantic consistency so the distance to the region of attraction that is semantically and geometrically consistent is much larger than for ANNFs. Note that our proposed SNNFs implements a 'soft' data association strategy: each pixel can be classified into one or multiple semantic classes. This allows our method to be robust to edge classification errors.

SNNFs is, in essence, an extension of the ANNFs, so it

can be seamlessly integrated to the optimization formulation in Eqn. 6. The new energy to optimize is:

$$E_{kr}^{\mathscr{E}} := \sum_{i=1}^{C} \sum_{\mathbf{p}_r \in \mathscr{C}_r^i} w_{\mathbf{p}_r}^{\mathscr{E}} \|\mathbf{p}_{kr} - n_{\mathscr{C}_r^i}(\mathbf{p}_{kr})\|_\gamma \qquad (7)$$

where $\mathscr{C}_r^i$ denotes the subset of edge pixels belonging to $i^{th}$ semantic class, $n_{\mathscr{C}_r^i}(\cdot)$ is the function returns the nearest neighbor from the $i^{th}$ semantic edge group in frame $k$. The relationship between $\{\mathscr{C}_r^i\}$ and $\mathscr{E}_r$ can be expressed as $\mathscr{C}_r^1 \cup \cdots \cup \mathscr{C}_r^C = \mathscr{E}_r$.

Following [15], we implement a point-to-tangent residual i.e. we project the original pixel-wise residual onto its local gradient direction to obtain additional robustness against outliers. It should be noted that this formulation makes the underlying assumption that the camera motion is free of large inter-frame rotations. In reality, this assumption is valid for the autonomous driving applications considered in this paper.

### C. Implementation Detail

We integrate the semantic edge constraints (Eqn. 7) into both tracking and mapping. In the tracking phase, we put more weights on the edge residuals to enforce a better convergence basin. In the mapping phase, we set smaller weights for edge terms and use a depth regularization component to penalize large inverse depth updates: the inverse depth of some edge pixels may be unobservable as the epipolar line are perpendicular to edge normals.

We use image gradient magnitude instead of intensity as it is more robust to illumination changes and preserves the high-frequency photometric information. Edge-based VO is robust to light variations but relying on image intensity may jeopardize this property. This requires only a slight code modification.

As detailed in Sec. IV-C.3, we implement a flexible point selection strategy to boost the tracking performances in vegetation dominated environments. Specifically, our proposed semantic edge VO algorithm not only selects edge pixels with semantic labels but also randomly select unlabeled pixels as support pixels to boost the robustness of motion estimation.

When the number of edge pixels is large enough, the pixel selector merely samples the minimum number of supportive pixels. However, as the firm edges get less or only occupied in small areas, we select more supportive points to keep well-distributed pixels into the optimization layer. We use the same sampling strategy described in [7] to choose the supportive pixels. Note that the supportive pixels only have photometric constraint while edge pixels hold both edge and photometric constraints. Since the supportive points don't have any labels, we don't push them into the semantic maps.

## IV. EXPERIMENTS

We first describe the experimental setup and the dataset. Qualitative results on semantic mapping results are presented in Sec. IV-B. Sec. IV-C evaluates the influence of several edge detection and registration algorithms on tracking accuracy and robustness. Our proposed semantic edge VO is compared to the state-of-art edge and semantic VO systems on KITTI in Sec. IV-D. Finally, the runtime performance is discussed in Sec. IV-E.

### A. Dataset and Metric

We test our method on sequences 00-10 of the KITTI odometry dataset [28]. We use the rectified color image of the left camera only with intrinsic calibration and ground truth camera poses.

TABLE I: KITTI Dataset Description

| Scene | Sequence No. | Description |
|---|---|---|
| city | 00, 05, 06, 07 | buildings, cars with few vegetation |
| village | 02, 03, 04, 08, 09 | vegetation with few buildings, cars |
| highway | 01 | roads, cars, and signs |

For fair comparison, we do not use loop-closure and all methods use the ground truth poses to recover the scale of motion every 200 frames. During evaluation, we find that the pose estimation for the first frames are unstable and vary for each method. So we discard the first ten pose estimates for all experiments. We choose the absolute trajectory error (ATE) as our metric: it measures the absolute difference between camera positions of two trajectories.

### B. Semantic Mapping

Fig. 4 shows large-scale semantic edge maps with pixel resolution generated with our approach on several environments. It should be noted that our proposed monocular VO system can only recover locally consistent 3D semantic edge maps rather than global ones since our algorithm also suffers from scale drift like all monocular methods.

### C. Edge-learning evaluation

We investigate the influence of the edge learning and registration methods on the robustness and accuracy of SNNFs (Fig. 3).
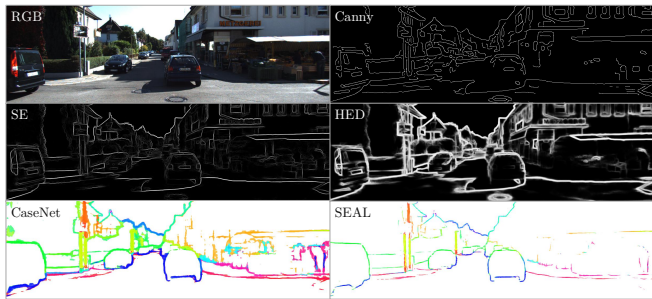


Fig. 3: Visualization of the evaluated edge methods.

*1) Edge Repeatability:* For outdoor VO applications, the choice of edge detector $E(\cdot)$ in Eqn. 3 is still an open question. [23] observes that the performance of edge-based VO highly depends on the *repeatability* of the edges: it computes the number of re-detected edge-pixels over the number of edge-pixels that should reappear. For fair comparison, we randomly choose 9000 edge pixels for each edge detectors. We run evaluation on vKITTI [29] since the ground truth depth is available. Fig. 5 shows the repeatability analysis: learned edges significantly outperform the conventional Canny detector.

*2) Edge Registration Evaluation:* One of the main advantage of our method is to improve tracking robustness and accuracy by integrating semantic information in the optimization. Experiments show that the convergence basin of our method is indeed larger than state-of-art ANNFs- and ONNFs-based methods [16].

Fig. 5 shows ATE for ANNFs, ONNFs, and SNNFs based edge registration approaches as a function of initial displacement using vKITTI dataset. We use the ground-truth depth maps to rule out the bias of depth reconstruction. We first use the ground truth poses as initial camera poses and control the initial displacements in a range of $[0,5]$m. The ATE is obtained by averaging all inter-frame tracking errors to compute the ATE. We observe that the convergence basin of our proposed SNNFs-based method is about two times larger than that of ONNFs and ANNFs. Note that the convergence test is not implementing any multi-level optimization techniques to improve the convergence basin. But our proposed VO system implements a pyramid-based tracking strategy to boost the tracking robustness further.

*3) Edge-Learning Evaluation:* We evaluate the influence of the edge-learning method on the VO performances. Fig 5 (bottom-right) shows the ATE on KITTI for ANNFs, ONNFs, and SNNFs using conventional edges, fused-semantic-edges and learned-semantic-edges.

We distinguish two ways to generate semantic edges: one is end-to-end learned semantic edges such as CaseNet and SEAL trained on Cityscapes [30]. The second one is the fusion of the learned-edges, HED or SE trained on BSDS [31], and learned-semantics trained separately from state-of-the-art DeepLabV3 [32]. We use the Xception-65 [33] model pretrained on Cityscapes and finetuned on KITTI [1]. We

---
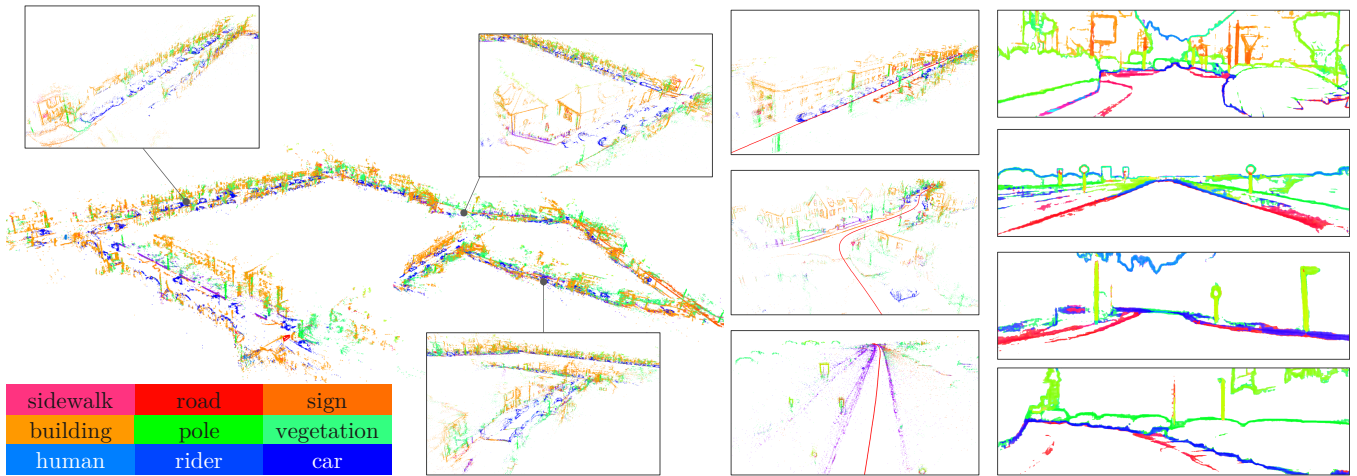
[1] https://github.com/hiwad-aziz/kitti_deeplab

Fig. 4: Reconstructed semantic edge maps for KITTI. Left: semantic edge maps recovered from city, village, and highway sequences. Right: semantic edge images generated using CaseNet [9].
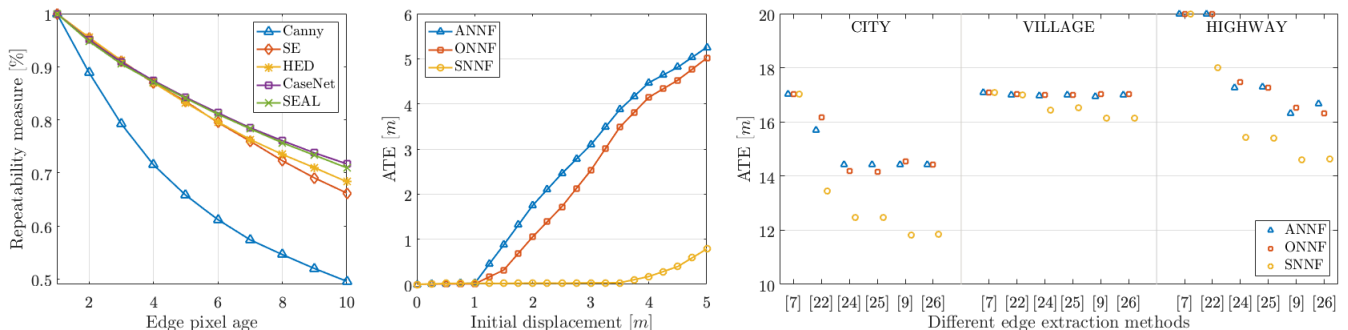


Fig. 5: Left: Repeatability analysis on vKITTI. Middle: convergence basins of ANNFs, ONNFs, and our proposed SNNFs edge registration on vKITTI. Right: tracking errors for different registration methods and different edge generation methods (city, village, highway scenes). We compare conventional edge detector (Canny [22]), learned edges (SE [24], HED [25]), and semantic edges (CaseNet [9], SEAL [26]). DSO [7] serves as a baseline VO.

assume that edge and segmentation probability are independent and multiply them to compute the fused semantic-edge probabilities. We observe that CaseNet and SEAL generalize well from Cityscapes to KITTI without further finetuning contrary to dense segmentation. This is a strong advantage of end-to-end methods over fusion ones. It should be noted that the edge VO algorithms implemented in this evaluation merely use edge pixels rather than the more flexible ones described in Sec. III-C, so that we can observe the pros and cons of integrating edge and semantic constraints in outdoor VO applications.

This evaluation allows us to draw five significant observations: (1) Introducing edge constraints into VO system is advantageous for operation in city and highway scenes, where there are many easy-to-track edges. However, it shows up little improvements or marginally worse performance for village datasets due to the poor repeatability and the few edges in vegetation-dominated images. (2) The VO systems using learned edges show a slightly better tracking accuracy compared with the one uses conventional edge detectors. This can be explained with the better repeatability of learned

edges in outdoor environments. (3) Semantic information benefits camera tracking more in city and highway scenes, where there are more semantic elements than in non urban areas. The performances are barely improved for the village images due to the edge sampling strategy. (4) End-to-end semantic edges have higher tracking accuracy than the fused ones which suggests that the end-to-end learning is more stable. (5) Introducing edge constraints cannot significantly improves the overall tracking precision in village image sequences, which indicates that merely relying on the edges in vegetation dominated environments could potentially jeopardize the overall tracking robustness.

These evaluations show that edge and semantic constraints significantly improve the tracking performance for environments with semantic elements and edges such as urban areas. However, it shows poor advantages in scenes with weak edges. In comparison, DSO implements a flexible sampling strategy to achieve similar tracking accuracy without incorporating edge and semantic constraints. We can find that the main limitation of pure edge VO is the deficiency of well-distributed edge pixels all over the image area. This
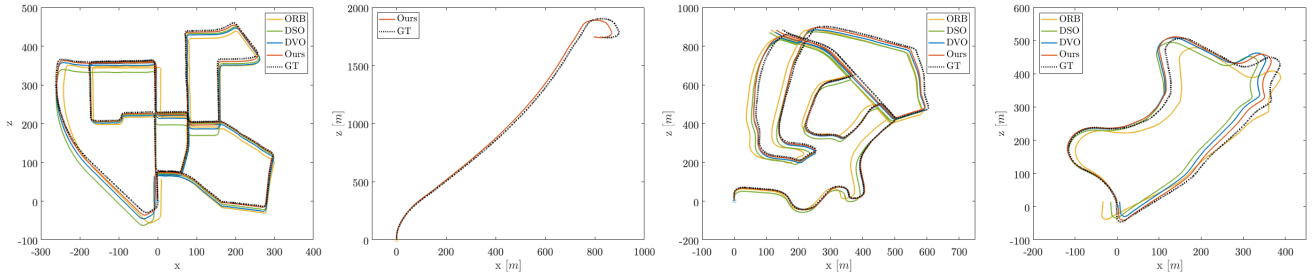
Fig. 6: Trajectories from our method, indirect ORBSLAM2, direct DSO, and semantic VSO systems on KITTI. Left to Right: KITTI-seq00, 01, 02, and 09. Note that seq01 merely shows our trajectory and the ground truth bcause other methods cannot generate the whole trajectory.

TABLE II: Tracking Error Comparison

| KITTI | city | | | | village | | | | | highway |
|---|---|---|---|---|---|---|---|---|---|---|
| | 00 | 05 | 06 | 07 | 02 | 03 | 04 | 08 | 09 | 01 |
| ORBSLAM2 | 16.14 | 15.96 | 13.35 | 10.63 | 15.58 | **3.44** | 3.05 | 15.43 | 12.88 | 36.32 |
| DSO | 16.83 | 13.64 | 16.83 | 9.55 | 17.08 | 3.71 | **3.01** | 18.31 | 13.05 | - |
| SVO | 15.31 | 10.08 | 14.10 | 8.39 | 14.57 | 3.76 | 3.09 | 15.29 | 13.12 | - |
| Proposed | **11.82** | **8.39** | **10.92** | **6.11** | **14.15** | 3.72 | 3.03 | **15.07** | **12.63** | **14.59** |

justifies our flexible sampling strategy to introduce pixels without labels to stabilize the motion estimation in vegetation dominated environment in Sec. III-C.

### D. Evaluation using Different Visual Odometry Methods

We assess the robustness and accuracy of the whole system on KITTI. We compare to mono-ORBSLAM2 [3], DSO [7] and VSO [8] as state-of-the-art monocular indirect, direct, and semantic VO algorithms for comparison. Since there is no released code for VSO, we implement it by introducing the semantic constraint energy into DSO for both tracking and mapping. For fair comparison, we choose 4000 active points for all approaches.

Table II shows the quantitative results and Fig. 6 shows trajectory comparisons on KITTI 00, 01, 02, and 09. Our methods outperforms the state-of-art in terms of tracking accuracy. Significant improvement can be observed on the highway trajectories, where only our proposed method could recover the full trajectory for such environment. Similar improvement are obtained on urban sequences and marginally better or equivalent performance for village sequences.

### E. Runtime

Runtime depends on the semantic edge pixels and the image size. In our experiments, we keep the original KITTI and vKITTI image sizes (around $(1242,375)px$), and use 3000 edge pixels and 1000 supportive pixels to balance accuracy and robustness. The processing time for both tracking and mapping is $1.34\times$ longer on average than DSO using 4000 active pixels. The semantic edge generations runs at $0.7s/image$ on an NVIDIA GTX1080Ti which makes our approach quasi-online.

## V. CONCLUSION

In this work, we present a monocular semantic edge VO framework capable of reconstructing 3D semantic edge maps in unstructured outdoor environments. Our proposed semantic nearest neighbor fields (SNNFs) offers several advantages over existing edge VO algorithms by using deep-learned semantic as a robust data association strategy. We analyse the influence of edge-learning and alignment methods on edge-based motion estimation and overcome the primary limitations of edge VO for outdoor application. An extensive evaluation of the accuracy and the robustness is conducted on KITTI. Results show that our method outperforms the state-of-art monocular direct, indirect, and semantic VO systems.

### REFERENCES

[1] G. Klein and D. Murray, "Parallel tracking and mapping for small ar workspaces," in *Proceedings of the 2007 6th IEEE and ACM International Symposium on Mixed and Augmented Reality*. IEEE Computer Society, 2007, pp. 1–10.

[2] H. Strasdat, J. Montiel, and A. J. Davison, "Scale drift-aware large scale monocular slam," *Robotics: Science and Systems VI*, vol. 2, no. 3, p. 7, 2010.

[3] R. Mur-Artal and J. D. Tardós, "Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras," *IEEE Transactions on Robotics*, vol. 33, no. 5, pp. 1255–1262, 2017.

[4] R. A. Newcombe, S. J. Lovegrove, and A. J. Davison, "Dtam: Dense tracking and mapping in real-time," in *2011 international conference on computer vision*. IEEE, 2011, pp. 2320–2327.

[5] M. Pizzoli, C. Forster, and D. Scaramuzza, "Remode: Probabilistic, monocular dense reconstruction in real time," in *2014 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2014, pp. 2609–2616.

[6] J. Engel, T. Schöps, and D. Cremers, "Lsd-slam: Large-scale direct monocular slam," in *European conference on computer vision*. Springer, 2014, pp. 834–849.

[7] J. Engel, V. Koltun, and D. Cremers, "Direct sparse odometry," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 3, pp. 611–625, 2018.

[8] K.-N. Lianos, J. L. Schonberger, M. Pollefeys, and T. Sattler, "Vso: Visual semantic odometry," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 234–250.

[9] Z. Yu, C. Feng, M.-Y. Liu, and S. Ramalingam, "Casenet: Deep category-aware semantic edge detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 5964–5973.

[10] J. Engel, J. Sturm, and D. Cremers, "Semi-dense visual odometry for a monocular camera," in *Proceedings of the IEEE international conference on computer vision*, 2013, pp. 1449–1456.

[11] L. Kneip, Z. Yi, and H. Li, "Sdicp: Semi-dense tracking based on iterative closest points." in *Bmvc*, 2015, pp. 100–1.

[12] P. F. Felzenszwalb and D. P. Huttenlocher, "Distance transforms of sampled functions," *Theory of computing*, vol. 8, no. 1, pp. 415–428, 2012.

[13] M. Kuse and S. Shen, "Robust camera motion estimation using direct edge alignment and sub-gradient method," in *2016 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2016, pp. 573–579.

[14] C. Kim, P. Kim, S. Lee, and H. J. Kim, "Edge-based robust rgb-d visual odometry using 2-d edge divergence minimization," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 1–9.

[15] Y. Zhou, L. Kneip, and H. Li, "Semi-dense visual odometry for rgb-d cameras using approximate nearest neighbour fields," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2017, pp. 6261–6268.

[16] Y. Zhou, H. Li, and L. Kneip, "Canny-vo: Visual odometry with rgb-d cameras based on geometric 3-d–2-d edge alignment," *IEEE Transactions on Robotics*, vol. 35, no. 1, pp. 184–199, 2019.

[17] X. Wang, W. Dong, M. Zhou, R. Li, and H. Zha, "Edge enhanced direct visual odometry." in *BMVC*, 2016.

[18] S. Li and D. Lee, "Fast visual odometry using intensity-assisted iterative closest point," *IEEE Robotics and Automation Letters*, vol. 1, no. 2, pp. 992–999, 2016.

[19] C. Yu, Z. Liu, X.-J. Liu, F. Xie, Y. Yang, Q. Wei, and Q. Fei, "Ds-slam: A semantic visual slam towards dynamic environments," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 1168–1174.

[20] H. Mahé, D. Marraud, and A. I. Comport, "Semantic-only visual odometry based on dense class-level segmentation," in *2018 24th International Conference on Pattern Recognition (ICPR)*. IEEE, 2018, pp. 1989–1995.

[21] M. Kaneko, K. Iwami, T. Ogawa, T. Yamasaki, and K. Aizawa, "Mask-slam: Robust feature-based monocular slam by masking using semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2018, pp. 258–266.

[22] J. Canny, "A computational approach to edge detection," in *Readings in computer vision*. Elsevier, 1987, pp. 184–203.

[23] F. Schenk and F. Fraundorfer, "Robust edge-based visual odometry using machine-learned edges," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2017, pp. 1297–1304.

[24] P. Dollár and C. L. Zitnick, "Fast edge detection using structured forests," *IEEE transactions on pattern analysis and machine intelligence*, vol. 37, no. 8, pp. 1558–1570, 2015.

[25] S. Xie and Z. Tu, "Holistically-nested edge detection," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1395–1403.

[26] Z. Yu, W. Liu, Y. Zou, C. Feng, S. Ramalingam, B. Vijaya Kumar, and J. Kautz, "Simultaneous edge alignment and learning," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 388–404.

[27] J. H. Manton, "Optimization algorithms exploiting unitary constraints," *IEEE Transactions on Signal Processing*, vol. 50, no. 3, pp. 635–650, 2002.

[28] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *2012 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2012, pp. 3354–3361.

[29] A. Gaidon, Q. Wang, Y. Cabon, and E. Vig, "Virtual worlds as proxy for multi-object tracking analysis," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4340–4349.

[30] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 3213–3223.

[31] P. Arbelaez, M. Maire, C. Fowlkes, and J. Malik, "Contour detection and hierarchical image segmentation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 33, no. 5, pp. 898–916, 2011.

[32] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *ECCV*, 2018.

[33] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1251–1258.