# OpenStreetMap-based LiDAR Global Localization in Urban Environment without a Prior LiDAR Map

Younghun Cho[1], Giseop Kim[2], Sangmin Lee[1] and Jee-Hwan Ryu[1*]

arXiv:2202.07516v1 [cs.RO] 15 Feb 2022

*Abstract*—Using publicly accessible maps, we propose a novel vehicle localization method that can be applied without using prior light detection and ranging (LiDAR) maps. Our method generates OSM descriptors by calculating the distances to buildings from a location in OpenStreetMap at a regular angle, and LiDAR descriptors by calculating the shortest distances to building points from the current location at a regular angle. Comparing the OSM descriptors and LiDAR descriptors yields a highly accurate vehicle localization result. Compared to methods that use prior LiDAR maps, our method presents two main advantages: (1) vehicle localization is not limited to only places with previously acquired LiDAR maps, and (2) our method is comparable to LiDAR map-based methods, and especially outperforms the other methods with respect to the top one candidate at KITTI dataset sequence 00.

*Index Terms*—Localization, Range Sensing, Mapping

## I. INTRODUCTION

Vehicle localization, or perceiving vehicle position, is an essential prerequisite for autonomous driving technologies. Autonomous driving algorithms, including path planning and following, cannot be accomplished without accurate vehicle position information. Owing to their high accuracy, Light Detection and Ranging (LiDAR) sensors have been actively studied for this purpose. LiDAR-based localization is usually performed by comparing real-time LiDAR scans with prior LiDAR maps created using LiDAR scans.

However, building LiDAR maps is a challenging: To build a LiDAR map, a vehicle equipped with a LiDAR sensor must visit every road in the target area, scanning every road that an autonomous vehicle might be using, which is a very challenging goal. Additionally, the LiDAR maps would need to be continuously updated every time new roads are built. Furthermore, the LiDAR scans must also be appropriately collated to create the LiDAR map. Dynamic objects and changing environments, such as passing cars and trees, make this process even more difficult, because these obstacles may hide static objects such as buildings. Thus, efforts have been

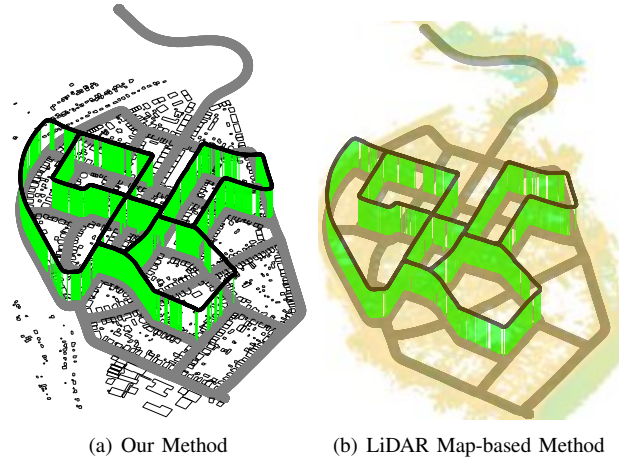|           |           |
|-----------|-----------|
| (a) Our Method | (b) LiDAR Map-based Method |

Fig. 1. This paper proposes a method that performs LiDAR-based global localization without a prior LiDAR map. Based on this method, we can obtain a comparable global localization result against PointNetVLAD. In particular, our method outperforms PointNetVLAD in terms of the top one accuracy at KITTI dataset sequence 00. (a) Localization results on OpenStreetMap. (b) Localization results on the LiDAR map. Black: query LiDAR scans trajectory. Gray: reference LiDAR scans trajectory (only for (b)). Green: successful localization within 5m.

made to replace LiDAR maps by finding ways to localize vehicles without visiting their real-world locations.

One such alternative approach was to use OpenStreetMap, which provides geographic, including roads and buildings, information. By utilizing static building information that does not vary in the short-term, OpenStreetMap-based descriptors can be generated at every road point to replace a LiDAR map. Several studies have utilized OpenStreetMap for localization problems; however, most of them perform visual localization using images [1, 2, 3]. Only a few studies have attempted to utilize OpenStreetMap for the vehicle localization problem, but prior attempts could not provide reliable localization results for autonomous driving [4]. The research and its limitattions are reviewed in Section II-B.

This paper proposes a 1D descriptor-based vehicle localization method using OpenStreetMap without prior LiDAR maps. Our method generates 1D descriptors, named OSM descriptors, from both OpenStreetMap information and LiDAR scans acquired from vehicles, to determine the vehicle's position by comparing those descriptors. In this process, a set of descriptors generated from OpenStreetMap replaces the conventional LiDAR maps. Fig. 1 shows that the OSM descriptor performs comparably with PointNetVLAD on the KITTI dataset sequence 00. Moreover, the OSM descriptor is complicated enough to discriminate between places without the adapting additional frameworks, such as Monte Carlo

Localization.

The contribution of this paper is as follows:

- To the best of our knowledge, our method is the first attempt for LiDAR kidnapped place recognition using OpenStreetMap and shows localization performance comparable to other methods using prior LiDAR maps.
- As opposed to conventional methods, our method is not limited to only places that are covered by LiDAR maps.
- Unlike conventional OpenStreetMap-based LiDAR localization methods, our method does not require additional processes such as a particle filter.

## II. RELATED WORKS

### A. Descriptor-based LiDAR Localization

LiDAR is the most widely used sensor for vehicle localization research. There are two forks of research that utilized LiDAR sensor: registration-based methods and feature-based methods. Registration-based methods, or dense methods, utilize all points in the point clouds for matching. In contrast, feature-based methods, or sparse methods, extract features from point clouds and match them with prior LiDAR maps. Methods using descriptors can also be considered as feature-based methods.

Himstedt et al., He et al., and Kim et al. developed 2D descriptors called GLARE, M2DP, and Scan context respectively to convert point clouds into descriptors [5, 6, 7, 8, 9]. GLARE transforms 2D LiDAR scans into histograms and uses the bag-of-words framework [5]. M2DP projects 3D point clouds into multiple 2D planes and uses their densities as descriptors, which are more robust even with the presence of noisy point clouds [6]. The Scan context extracts the highest point for every angle and calculates the distance from the point clouds, outperforming other descriptor-based methods [8, 9].

Recently, deep learning-based methods have been widely applied to LiDAR localization. The most representative works include L3-Net and PointNetVLAD [10, 11]. L3-Net was the first learning-based LiDAR localization method and was developed into the DeepVCP [12]. As can be inferred from the title of the paper, PointNetVLAD is a combination of two methods: PointNet and NetVLAD [13, 14]. This method achieves permutation invariance, because NetVLAD is a symmetric function. Additionally, the researchers adopted a lazy triplet loss function that maximizes the difference between descriptors.

However, these LiDAR localization methods rely on prior LiDAR maps, because they compare point clouds or features with prior LiDAR maps, which are very limited.

### B. OpenStreetMap-based Localization

Therefore, several researchers sought to localize autonomous vehicles using OpenStreetMap. Most OpenStreetMap-based localization methods utilized visual sensors rather than LiDARs, or overhead images rather than OSM data, for their works [1, 2, 3]. Furthermore, as mentioned earlier, there have been a few attempts to achieve LiDAR localization using OpenStreetMap. Floros et al. fitted trajectories generated from visual odometry to OpenStreetMap road information using Chamfer matching [15]. Ruchti et al. also fitted the calculated trajectory to OpenStreetMap road information, but it used the LiDAR point clouds instead of visual sensors. Road and non-road road cells from LiDAR point clouds were classified, and the output was used as the weights of the Monte Carlo Localization [16]. However, these methods will not function accurately when a vehicle is driven straight for a long time, because they only use the topological information of road networks.

In contrast, Vysotska and Stachniss utilized the building information instead of the road information to localize vehicles [17]. LiDAR point clouds were fitted to the OpenStreetMap building information, and the results were used as constraints for graph-based SLAM. Suger and Burgard classified the semantic terrain information of the surrounding environment using LiDAR point clouds, and used the result for the Monte Carlo Localization's weighting function [18]. Yan et al. suggested 4-bit semantic descriptors, classifying the existence of buildings and roads in four directions, and performing Monte Carlo Localization using their descriptors [4]. However, the results of this method are not sufficiently accurate for autonomous driving, because the resulting trajectories pass through buildings.

Another major constraint of these localization methods is that they cannot be used independently, but must be combined with other methods such as Monte Carlo Localization or graph-based SLAM. Additionally, because of the fundamental problem of a Monte Carlo Localization, time is required to initialize the vehicle's position. On the contrary, our method is a stand-alone method for kidnapped localization problems that can be used independently from other methods.

## III. VEHICLE LOCALIZATION USING OSM CONTEXT

The pipeline of the proposed method is visualized in Fig. 2. This chapter states the problem definition and describes the methods used in this paper, which consist of four modules: OSM descriptor generation, LiDAR descriptor generation, rotation invariant key generation, and finding the most similar descriptor.

### A. Problem Definition

This paper focuses on the problem of LiDAR localization without prior visits to the target area. Specifically, we build a reference map consisting of OSM descriptors at each position in the target area, without visiting the actual location. Because LiDAR scans contain comprehensive information about the surrounding environment, while OSM contains only building information, only the building points should be extracted from LiDAR scans to warrant an equal comparison. We converted each extracted building point from the LiDAR scans into 1D descriptors, which have the same dimensions as the OSM descriptor, and found the descriptor most similar to the reference map. We assume that the ground is locally flat and buildings are not necessarily rectangular but structured and have straight vertical walls.

Denoting the function that converts OpenStreetMap building information, $b_n$, into OSM descriptor as $f(.)$, the reference map is built as a set of OSM descriptors $M =$
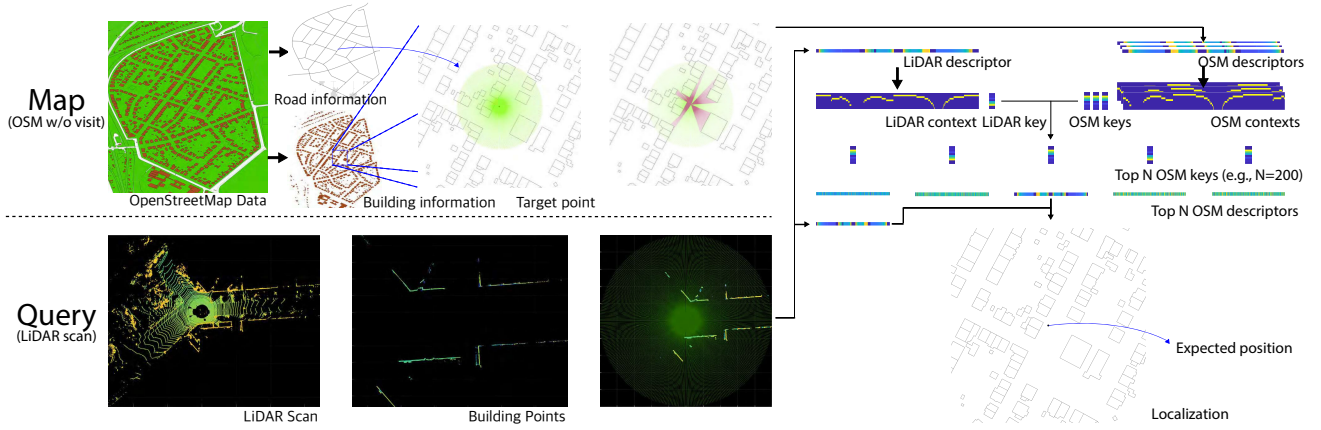
Fig. 2. Overall pipeline of the proposed method. We generated a reference map using the OpenStreetMap data by converting the building information into OSM descriptors and OSM keys. For every query LiDAR scan, we converted it into a LiDAR descriptor and LiDAR key. Comparing them in a 2-stage method, we can find the most similar OSM descriptor to localize the vehicle's position.
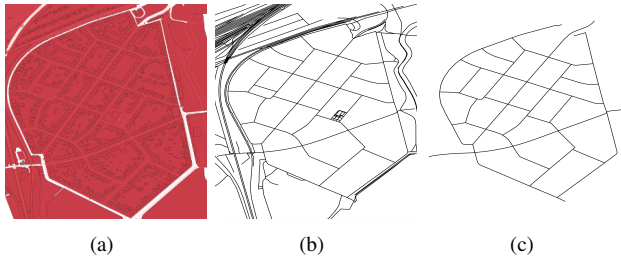


Fig. 3. OSM data visualization. (a) Entire information provided by OSM. (b) 'lines' layer of OSM. (c) Extracted road information and interpolated (x,y) points.
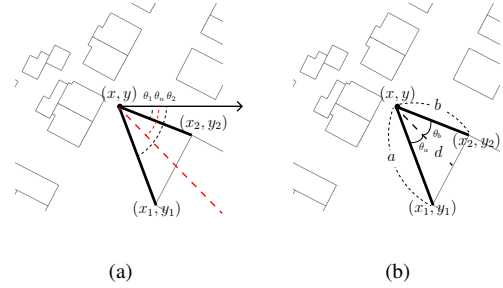


Fig. 4. (a) Distance is calculated only if the target angle is in between two vertices of the building edge. (b) Distance to building is calculated using the fact that a sum of two small triangle areas is equal to the area of a large triangle.

$\{m_1, m_2, ..., m_N\}$ where $m_n = f(b_n)$. Given a query point cloud $q$, our aim is to find an OSM descriptor, $m_*$, with the minimum distance.

To solve this problem, we extracted a building point cloud, $q' = R(q)$, where $R(.)$ is a sematic segmentation network using RangeNet++ [19]. Using the function $f(.)$, we generated a LiDAR descriptor $\bar{q} = f(q')$. By comparing $\bar{q}$ and $M$, we found the OSM descriptor with the minimum distance using the distance function $d(.)$, where $d(\bar{q}, m_*) \leq d(\bar{q}, m_n), n \leq N$. In this paper, we used the L1-norm as the distance function $d(.)$.

### B. OSM Descriptor Generation

The OSM descriptor is a 1D descriptor that is designed to replace LiDAR maps in the vehicle localization scheme. The OSM descriptor mimics a 2D LiDAR scan using the OpenStreetMap building information, and the set of OSM descriptors in the target area can replace the LiDAR map.

An OSM descriptor is generated by calculating the distance to the buildings using OpenStreetMap data. OpenStreetMap provides five layers of information: 'points', 'lines', 'multiline strings, 'multi polygons', and 'other relations'. 'multi polygons' layer provides building information as tuples that represent building edges, and 'lines' layer provides roads, rivers, and other line information. Among the 'lines' layer, road information is labeled as 'highway' and provides sparse points that are just enough to express the road shape as shown in Fig. 3. Therefore, we interpolated road points with a 1-meter

interval and generated an OSM descriptor for every interpolated point $(x, y)$ in UTM converted coordinate. Denoting a tuple of the building edges $t = (X_1, X_2) = ((x_1, y_1), (x_2, y_2))$ and the target position in the UTM coordinate as $(x, y)$, we find the shortest distance to a building $d_n^*$ for each search angle $\theta_n$ using (1).

$$d_n^* \leq d_n = \frac{ab \sin(\theta_a + \theta_b)}{a \sin(\theta_a) + b \sin(\theta_b)}, \quad \theta_1 \leq \theta_n \leq \theta_2, \quad (1)$$

where $a$ and $b$ denote the distance between the current position and $t$. $\theta_1$ and $\theta_2$, respectively, denote the smaller and larger angle between the horizontal axis and the line segments joining the target point $(x, y)$ and building vertices $(x_1, y_1)$ and $(x_2, y_2)$, and $\theta_a = \theta_n - \theta_1$ and $\theta_b = \theta_2 - \theta_n$ denote the angles between the target angle and the angle to $t$, as shown in Fig. 4. If there are no buildings at the target angle, $d_*$ is considered as zero. $d_n$ is calculated when $\theta_1 \leq \theta_n \leq \theta_2$, which means that the target angle is in between the two vertices of the building edge. By stacking the calculated $d_n^*$ from all angles, the OSM descriptor was defined as $m_n = \{d_1^*, d_2^*, ..., d_{360}^*\}$.

### C. LiDAR Descriptor Generation

To compare the LiDAR scans with the OSM descriptors, LiDAR scans should be converted into the same format as the OSM descriptors. First, we extracted the building points in a

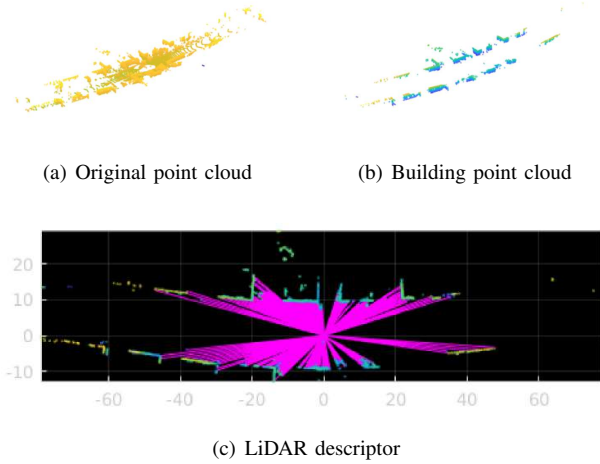(a) Original point cloud    (b) Building point cloud



(c) LiDAR descriptor

Fig. 5. Visualization of LiDAR descriptor generation process.

LiDAR scan using RangeNet++ [19], because the OSM descriptors are made using building information only. However, LiDAR scans consist of points, whereas OpenStreetMap building information consists of the vertices of the buildings only. Therefore, a query LiDAR descriptor $\bar{q} = \{d_1^*, d_2^*, ..., d_N^*\}$ can be generated by calculating the shortest distance $d_n^* \leq d_n$ to building points $P = \{X_1, X_2, ..., X_N\}$ at each search angle, where $d_n$ denotes the distance to the point $X_n$. The entire process for Section III-C is shown in Fig. 5.

### D. Rotation-Invariant Descriptor Generation

Because vehicles do not always travel in one direction only, there can be multiple scenarios of cars in a spot: cars can come in from two opposite roads, or even from multiple roads, in the case of an intersection. To account for localization in all of these cases, the building rotation-invariant descriptors are pivotal. Here, we introduce a method for making a previously described descriptor invariant to rotation.

This can be accomplished simply by summing or averaging the features of all the angles. However, adding all 1-dimensional vectors yield only one value, which can diminish the discriminatory power as a descriptor. Therefore, to avoid the loss, we transform the OSM descriptors and LiDAR descriptors that are 1-dimensional vectors into 2-dimensional context $\hat{q}'$ by setting $\lceil \bar{q}(i)/l_b \rceil$-th element of every integer angle $i$-th column to be one and others to be zero where $l_b$ denotes the length of each bin of the descriptor as shown in (2). In this paper, we used $l_b = 5\,m$ because 5 m bin performed better compared to 2 m and 10 m bin as shown in Table II.

$$\hat{q}'(k,i) = \begin{cases} 1, & k = \lceil \bar{q}(i)/l_b \rceil \\ 0, & others \end{cases}, \quad 1 \leq i \leq 360, \quad i \in Z,$$
(2)

where $\bar{q}(i)$ denotes $i$-th value of the LiDAR descriptor $\bar{q}$, and $\lceil * \rceil$ is a ceil operator. By summing up the features of this 2-dimensional context for all angles at each row (i.e., interval) using (3), the rotation-invariant descriptor OSM keys and LiDAR keys are built as shown in Fig. 6.
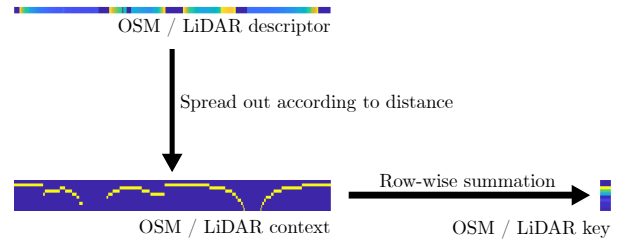


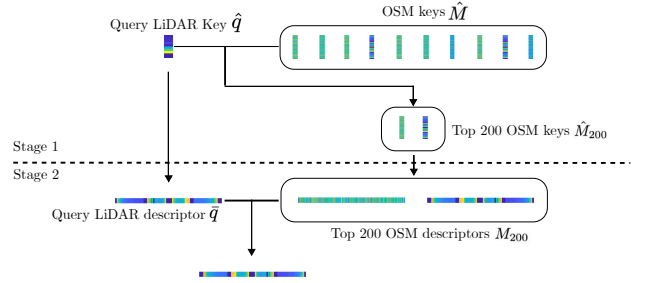Fig. 6. rotation invariant key generation



Fig. 7. Pipeline of our 2-stage method. In the first stage, we sort OSM keys that have smaller dimensions compared to the OSM descriptors for performing faster calculations. We selected 200 keys; we could only compare 200 descriptors among 14000 candidates in the KITTI dataset sequence 00.

$$\hat{q}(d) = \sum_{i=1}^{360} \hat{q}'(d,i), \quad 1 \leq d \leq R/l_b, \quad d \in Z,$$
(3)

where $R$ denotes the maximum sensor range. In this paper, we used $R = 50$. For OSM descriptors, we can apply the same process to make rotation-invariant OSM keys $\hat{M} = \{\hat{m_1}, \hat{m_2}, ..., \hat{m_N}\}$.

### E. Finding the Most Similar Descriptor

Finally, we were able to localize a vehicle by comparing the query LiDAR descriptor $\bar{q}$ and the set of OSM decriptors $M$ in two stages, as shown in Fig. 7. In the first stage, we compared the query LiDAR key $\hat{q}$ of the query LiDAR descriptor $\bar{q}$ with the set of OSM keys $\hat{M}$ to rank and extract the top 200 OSM keys $\hat{M}_{200}$ with the highest similarity using the L1-norm. The 1-stage method provides sorted candidates for the second stage.

Subsequently, in the second stage, the top 200 OSM descriptors $M_{200}$ corresponding to the $\hat{M}_{200}$ are compared once again to the query LiDAR descriptor $\bar{q}$, to re-rank the descriptors. At the second stage, $\bar{q}$ and $M$ are no longer rotation invariant. However, we compared them by rotating $\bar{q}$ for every angle. Although the comparison itself is a time-consuming process, it can be done within $0.0624s$ per frame because we have done this only with the top 200 candidates. This two-step method exhibits not only a high speed, achieved with initial screening using rotation-invariant descriptors, but also a high accuracy, achieved in comparison with the OSM descriptors. If the location of the highest-ranked OSM descriptors and that

TABLE I
INFORMATION OF DATASETS USED IN THIS PAPER

| Dataset | | Sequence | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | 00 | 02 | 05 | 06 | 07 | 08 | 09 | 10 |
| KITTI | Name | | | | | | | | |
| | # LiDAR scan | 4541 | 4661 | 2761 | 1101 | 1101 | 4071 | 1591 | 1201 |
| | # OSM contexts | 7010 | 19523 | 11074 | 11074 | 7010 | 8250 | 11036 | 11036 |
| KITTI-360 | Name | | 00 | | 05 | | 06 | | 09 |
| | # LiDAR scan | | 11518 | | 6743 | | 9699 | | 14056 |
| | # OSM contexts | | 10528 | | 13512 | | 9363 | | 7664 |

of the query LiDAR descriptor is within 5 m, the result is considered as a success.

## IV. EXPERIMENTS

### A. Benchmark Datasets

We used two of the publicly available datasets that included LiDAR scans as query LiDAR scans: KITTI [20] and KITTI-360 [21]. Both provide several sequences containing LiDAR scans, images, and poses of residential areas. In this paper, we utilized LiDAR scans and poses as query data. We converted all of the LiDAR scans into LiDAR descriptors using the method illustrated in Section III-C. Information of the KITTI and KITTI-360 datasets are summarized in Table I.

For the reference map, we used OpenStreetMap, which provides 'points', 'lines', 'multiline strings, 'multi polygons', and 'other relations' data. Among these data, we utilized lines and multipolygons that represent lines, such as roads and rivers, and polygons, such as buildings and districts on the map. By filtering them, we extracted the road and building information. Because road data are inconsistent and sparse, interpolation was necessary to make them consistent and dense. We used the 1 m interpolation for this experiment. For each KITTI and KITTI-360 dataset sequence, we exported OpenStreetMap data of a rectangular area that encompass the entire trajectory of the sequence. Using exported data, we generated OSM descriptors for each location with interpolated road and building data using the method illustrated in Section III-B.

### B. Comparison Methods

We compared our method with PointNetVLAD [11], a state-of-the-art LiDAR point cloud retrieval method. We used the same parameter setting as PointNetVLAD: use ground-removed point cloud, point cloud within 25 m, normalized it into with zero as the mean and range between $[-1, 1]$, and 4096 points. Using the pre-trained model provided by the author, we converted LiDAR scans into 256-dimensional descriptors.

We only compared the performances at the KITTI dataset sequence 00, because while our method can perform a Li-DAR localization at every sequence, PointNetVLAD can be performed only if one sequence fully encompasses the other sequence. In the KITTI and KITTI-360 datasets, only KITTI-360 sequence 09 encompasses KITTI sequence 00. Therefore, we compared the results of PointNetVLAD using KITTI sequence 00 as query LiDAR scans and KITTI-360 sequence 09 as the reference map.

## V. RESULTS

In this section, we present the results of analyzing the performance of our method in an outdoor large-scale LiDAR localization, especially in kidnapped situations. We compared the performance of PointNetVLAD, which also performed outdoor large-scale LiDAR localizations. We evaluated the accuracies of each method using the top one, top five, and top 10 candidates, to show that our method is robust even with a smaller number of candidates. In particular, we describe how our method is specialized for autonomous driving situations, because it has a rotation-invariant characteristic that can distinguish intersections and shows higher accuracy with a smaller number of candidates.

### A. LiDAR Localization Performance on the KITTI dataset

First, we evaluated the performances on twelve KITTI and KITTI-360 datasets using the top $N$ candidates. In this paper, we used $N = 200$. Fig. 8 and Fig. 9 illustrate the qualitative results of the KITTI and KITTI-360 dataset sequences, respectively. Fig. 10 shows the top five candidates for the query LiDAR point cloud. It shows our descriptor's rotation-invariant characteristic, because they find the right reference OSM descriptors even when the query LiDAR descriptor and reference descriptor's directions are different.

The evaluation was performed using every query LiDAR scan in the KITTI and KITTI-360 dataset sequences, even if a scan revisits the same place multiple times. This suggests robustness against rotation, because our OSM keys are rotation-invariant. Especially, the result at KITTI-360 sequence 00 shows our method's rotation-invariant performance, indicating over 95% accuracy at intersections, where the vehicle passed from at least two different directions. In practice, autonomous vehicles usually follow a lane in normal situations and make decisions at intersections, where they can only change their direction. Therefore, the ability to distinguish intersections is important for autonomous driving.

The quantitative results of our method are presented in Table II. The magnitude of the difference between the top one and top five is similar to that between the top five and top 10. This demonstrates that our method can yield as effective results even with a lower number of candidates. Several sequences show poor performance, which we will analyze in Section V-D.

### B. Comparison with PointNetVLAD

Table III and Fig. 11 show our method's performance compared with PointNetVLAD. The KITTI dataset sequence 00

TABLE II
KITTI AND KITTI-360 LOCALIZATION ACCURACY (%)

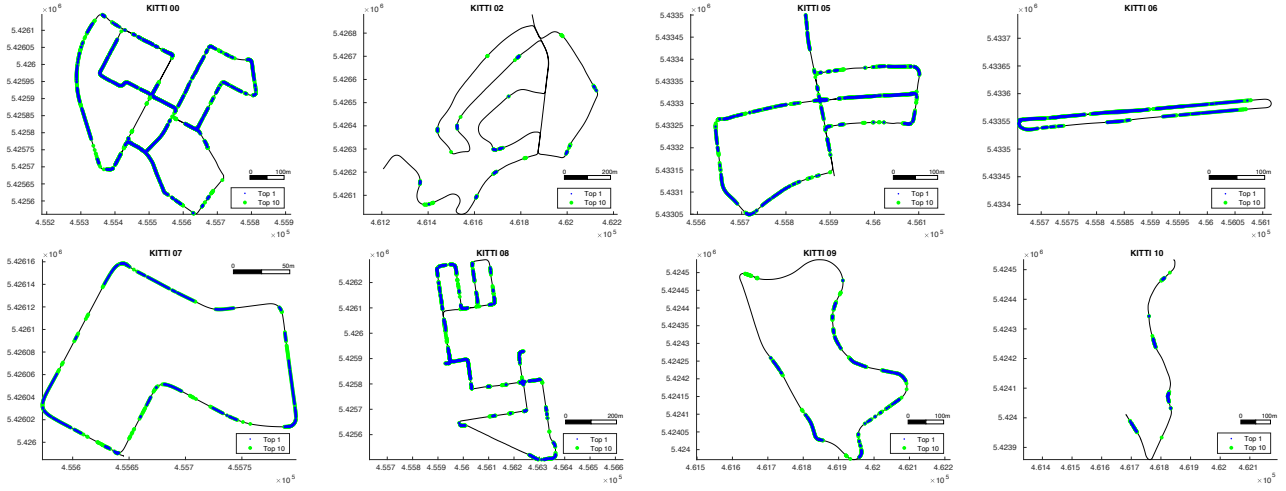| Method ($l_b$) | | KITTI | | | | | | | | KITTI-360 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 00 | 02 | 05 | 06 | 07 | 08 | 09 | 10 | 00 | 05 | 06 | 09 |
| **Ours (2m)** | top1 | 45.72 | **2.45** | 36.87 | 59.67 | 37.06 | 27.09 | 19.80 | 11.49 | **40.04** | 8.36 | **8.97** | 34.40 |
| | top5 | 54.22 | **3.15** | **45.31** | 65.40 | **51.04** | 35.37 | 24.07 | 12.66 | **47.78** | 11.51 | **11.56** | 42.72 |
| | top10 | 57.74 | **3.43** | **50.49** | 68.12 | 54.22 | 38.79 | 25.83 | **13.66** | **52.34** | 13.76 | **13.45** | 46.74 |
| **Ours (5m)** | top1 | **48.34** | 1.85 | **37.52** | **62.13** | **37.33** | **29.89** | **22.25** | **11.66** | 39.89 | **9.31** | 8.35 | **35.64** |
| | top5 | **56.86** | 2.36 | 44.62 | **67.03** | 50.77 | **36.70** | **26.21** | **12.99** | 46.67 | **12.62** | 10.81 | **44.61** |
| | top10 | **61.15** | 2.64 | 49.51 | **70.75** | **56.49** | **39.50** | **29.23** | 13.57 | 50.55 | **15.10** | 12.76 | **48.78** |
| **Ours (10m)** | top1 | 44.13 | 1.37 | 30.32 | 57.77 | 36.60 | 25.84 | 21.68 | 11.07 | 37.70 | 7.50 | 7.45 | 32.42 |
| | top5 | 52.87 | 1.72 | 35.97 | 61.58 | 50.40 | 32.60 | 26.15 | 12.32 | 43.79 | 10.08 | 9.89 | 40.59 |
| | top10 | 56.20 | 1.93 | 38.90 | 64.49 | 54.59 | 35.64 | 28.41 | 12.99 | 47.24 | 12.03 | 11.58 | 44.62 |



Fig. 8. Top 10 accuracies of the KITTI dataset sequences. Blue: successful localization with the top one candidate. Green: successful localization with the top 10 candidates.
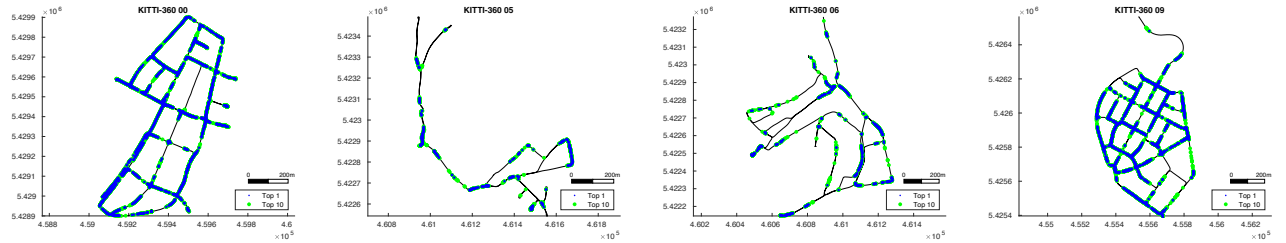


Fig. 9. Top 10 accuracies of the KITTI-360 dataset sequences. Blue: successful localization with the top one candidate. Green: successful localization with the top 10 candidates.

was used for comparison, as mentioned in Section IV-B. Our method shows a comparable result to PointNetVLAD, even though we did not utilize LiDAR maps and deep-learning networks. In particular, our method outperforms PointNetVLAD at the top one accuracy.

### C. Runtime Analysis

The OSM descriptor is also a lightweight descriptor compared to the PointNetVLAD. Our method consists of two stages; therefore, we compare the performance and runtime for 1-stage and 2-stage results. Table IV shows the runtimes for each method executed in the same hardware setting except for GPU usage in PointNetVLAD. All experiments were carried out on a PC with an Intel i7-6700 CPU at 3.40 GHz and 32 GB memory. For PointNetVLAD experiments, we used an additional GeForce GTX 1080 Ti GPU. The localization time is proportional to the size of the descriptor, because we used the same comparison algorithm for all the descriptors. The sizes of the descriptors are 1-by-256, 1-by-10, and 1-by-360 for PointNetVLAD, OSM keys, and OSM descriptors, respectively. However, the descriptor generation time is comparable with PointNetVLAD, even though PointNetVLAD was processed on a GPU.

Additionally, Table IV shows that there is a 0.0534 s difference between the 1-stage method and 2-stage method. This means that it took 0.0534 s to compare the query de-
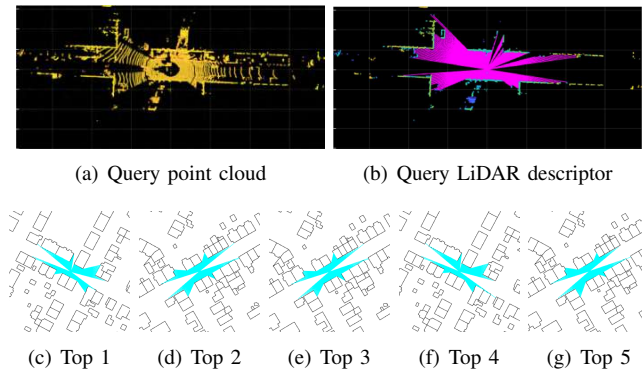
(a) Query point cloud      (b) Query LiDAR descriptor



(c) Top 1    (d) Top 2    (e) Top 3    (f) Top 4    (g) Top 5

Fig. 10. Top 5 reference OSM descriptors against query LiDAR point cloud.

TABLE III
PERFORMANCE AGAINST POINTNETVLAD (%)

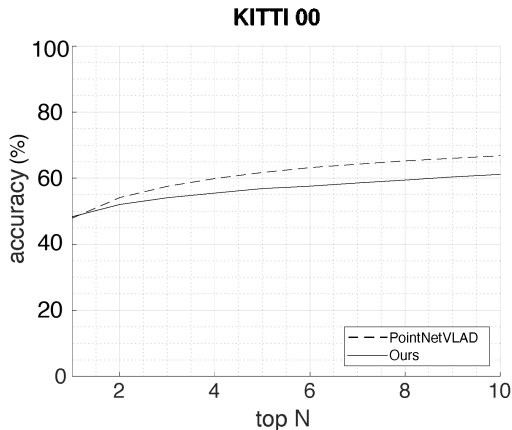| method | Top 1 | Top 5 | Top 10 |
|---|---|---|---|
| **PointNetVLAD** | 47.90 | **61.73** | **66.79** |
| **Ours** | **48.34** | 56.86 | 61.15 |



Fig. 11. Top 10 accuracies for our method and PointNetVLAD. Our method performs comparable global localization without a prior LiDAR map and ourperforms with the top one candidate.

scriptor and the top 200 OSM descriptors. If we compare all descriptors with out-of-sorting in stage 1, it might take $1.872 = 0.0534 * (7010/200)$ more seconds for every query LiDAR scan, which takes 8,499 more seconds for the entire KITTI dataset sequence 00. Table V and Fig. 12 compares the accuracy of the 1-stage method and 2-stage method.

### D. Failure Cases

As mentioned above, several sequences showed poor performances. In this subsection, we analyze the failure cases in three scenarios. In the first scenario, as shown in failure example 1 in Fig. 13, the building is blocked by fences and trees or the building is on a hill. In this case, the vehicle perceives that there is no building around the vehicle, but the OSM descriptor is surrounded by buildings. Dynamic point removal algorithms such as ERASOR [22] can help resolving this problem. In the second case, a similar building pattern was repeated over a long interval. In failure example 2, as shown in

TABLE IV
RUNTIME COMPARISON PER FRAME

| | PointNetVLAD | Ours (1-Stage) | Ours (2-Stage) |
|---|---|---|---|
| Descriptor Generation | **0.0690s** | 0.1714s | 0.1714s |
| Localization | 0.0167s | **0.0090s** | 0.0624s |
| Total | **0.0857s** | 0.1804s | 0.2338s |

TABLE V
ACCURACY OF 2-STAGE METHOD (%)

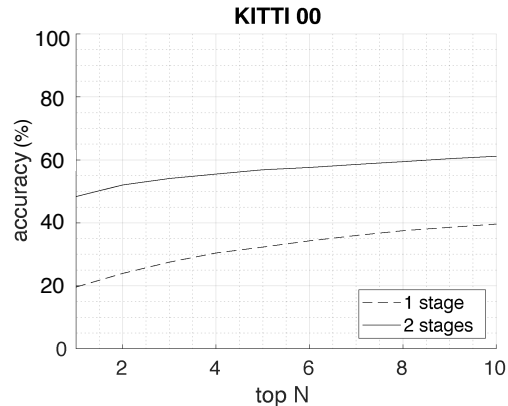| method | Top 1 | Top 5 | Top 10 |
|---|---|---|---|
| **1 stage** | 19.60 | 32.31 | 39.60 |
| **2 stages** | **48.34** | **56.86** | **61.15** |



Fig. 12. Global localization accuracy of the 1-stage method and 2-stage method.

Fig. 13, both sides of the vehicle are full of buildings without a gap between the buildings. Therefore, our descriptor loses its discrimination power in these areas. In the last case, if there are no buildings around the vehicle, such as a highway or off-road, the building-based descriptor cannot be made, and our method does not work.

## VI. CONCLUSION

In this study, we proposed a localization method using the OpenStreetMap information which can be executed without any prior LiDAR maps. We showed that our method performs comparably with PointNetVLAD, which uses LiDAR maps and deep learning techniques, in terms of accuracy and efficiency. The newly proposed novel descriptor, context, and key for OSM and LiDAR data enables a fast and highly accurate localization performance invariant to rotation. We verified its performance on 12 sequences of the KITTI dataset and KITTI-360 dataset. Our method outperforms PointNetVLAD in terms of the top one accuracy at KITTI dataset sequence 00, and showed comparable performance in top five and ten accuracy. However, in KITTI sequence 02, buildings are blocked by other environments; therefore, our method does not perform well because a LiDAR scan measures walls around the vehicle, while the OSM descriptor recognizes that the same place is surrounded by buildings. In future work, we aim to improve our method to function in challenging environments with
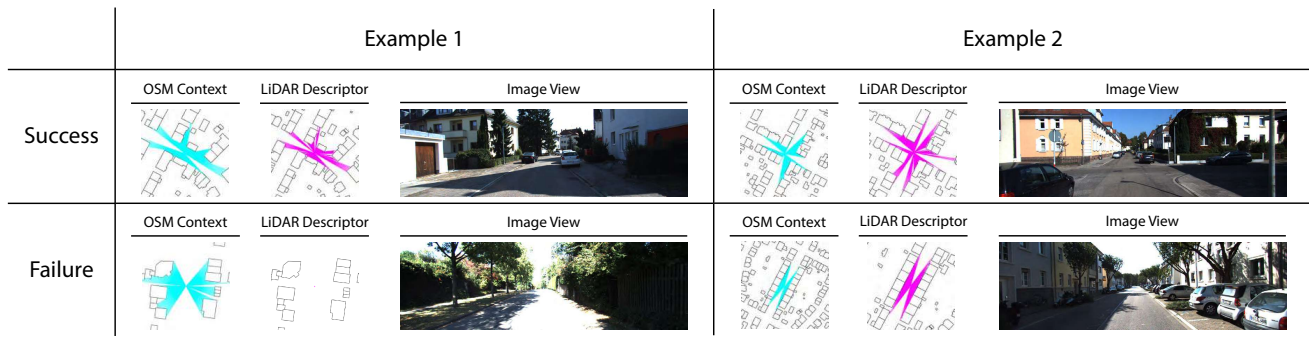
Fig. 13. Some sequences showed poor localization accuracy, with our method being weak in three situations: first, as shown in Example 1, if the buildings are blocked by fences, trees, or hills, there is no building information in LiDAR scans, which generates a blank descriptor; second, as shown in Example 2, if there are repeated building patterns, there are multiple candidates that have similar descriptors, which decreases the discrimination ability of our descriptors; finally, if there is no building around a vehicle, such as a highway or forest, a building-based descriptor cannot be made.

sparse buildings, such as the KITTI sequence 02 and MulRan sequence 'Sejong city' [23].

## REFERENCES

[1] M. A. Brubaker, A. Geiger, and R. Urtasun, "Lost! leveraging the crowd for probabilistic visual self-localization," in *Proceedings of the IEEE Conf. on computer vision and pattern recognition*, 2013, pp. 3057–3064.

[2] P. Panphattarasap and A. Calway, "Automated map reading: Image based localisation in 2-d maps using binary semantic descriptors," in *2018 IEEE/RSJ Int. Conf. on Intel. Robots and Syst. (IROS)*, 2018, pp. 6341–6348.

[3] T. Y. Tang, D. De Martini, and P. Newman, "Get to the Point: Learning Lidar Place Recognition and Metric Localisation Using Overhead Imagery."

[4] F. Yan, O. Vysotska, and C. Stachniss, "Global localization on openstreetmap using 4-bit semantic descriptors," in *2019 European Conf. on Mobile Robots (ECMR)*, 2019, pp. 1–7.

[5] M. Himstedt, J. Frost, S. Hellbach, H.-J. Böhme, and E. Maehle, "Large scale place recognition in 2D LIDAR scans using geometrical landmark relations," in *2014 IEEE/RSJ Int. Conf. on Intel. Robots and Syst.*, 2014, pp. 5030–5035.

[6] L. He, X. Wang, and H. Zhang, "M2DP: A novel 3D point cloud descriptor and its application in loop closure detection," in *2016 IEEE/RSJ Int. Conf. on Intel. Robots and Syst. (IROS)*, 2016, pp. 231–237.

[7] G. Kim, B. Park, and A. Kim, "1-day learning, 1-year localization: Long-term lidar localization using scan context image," *IEEE Robotics and Automat. Letters*, vol. 4, no. 2, pp. 1948–1955, 2019.

[8] G. Kim and A. Kim, "Scan context: Egocentric spatial descriptor for place recognition within 3d point cloud map," in *2018 IEEE/RSJ Int. Conf. on Intel. Robots and Syst. (IROS)*, 2018, pp. 4802–4809.

[9] G. Kim, S. Choi, and A. Kim, "Scan Context++: Structural Place Recognition Robust to Rotation and Lateral Variations in Urban Environments," *IEEE Transactions on Robotics*, 2021.

[10] W. Lu, Y. Zhou, G. Wan, S. Hou, and S. Song, "L3-net: Towards learning based lidar localization for autonomous driving," in *Proceedings of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition*, 2019, pp. 6389–6398.

[11] M. A. Uy and G. H. Lee, "Pointnetvlad: Deep point cloud based retrieval for large-scale place recognition," in *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition*, 2018, pp. 4470–4479.

[12] W. Lu, G. Wan, Y. Zhou, X. Fu, P. Yuan, and S. Song, "Deepvcp: An end-to-end deep neural network for point cloud registration," in *Proceedings of the IEEE/CVF Int. Conf. on Computer Vision*, 2019, pp. 12–21.

[13] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3d classification and segmentation," in *Proceedings of the IEEE Conf. on computer vision and pattern recognition*, 2017, pp. 652–660.

[14] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, "NetVLAD: CNN architecture for weakly supervised place recognition," in *Proceedings of the IEEE Conf. on computer vision and pattern recognition*, 2016, pp. 5297–5307.

[15] G. Floros, B. Van Der Zander, and B. Leibe, "Openstreetslam: Global vehicle localization using openstreetmaps," in *2013 IEEE Int. Conf. on Robotics and Automat.*, 2013, pp. 1054–1059.

[16] P. Ruchti, B. Steder, M. Ruhnke, and W. Burgard, "Localization on openstreetmap data using a 3d laser scanner," in *2015 IEEE Int. Conf. on Robotics and Automat. (ICRA)*, 2015, pp. 5260–5265.

[17] O. Vysotska and C. Stachniss, "Exploiting building information from publicly available maps in graph-based SLAM," in *2016 IEEE/RSJ Int. Conf. on Intel. Robots and Syst. (IROS)*, 2016, pp. 4511–4516.

[18] B. Suger and W. Burgard, "Global outer-urban navigation with openstreetmap," in *2017 IEEE Int. Conf. on Robotics and Automat. (ICRA)*, 2017, pp. 1417–1422.

[19] A. Milioto, I. Vizzo, J. Behley, and C. Stachniss, "RangeNet++: Fast and Accurate LiDAR Semantic Segmentation," in *IEEE/RSJ Intl. Conf. on Intel. Robots and Syst. (IROS)*, 2019.

[20] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets Robotics: The KITTI Dataset," *Int. Journal of Robotics Research (IJRR)*, 2013.

[21] J. Xie, M. Kiefel, M.-T. Sun, and A. Geiger, "Semantic Instance Annotation of Street Scenes by 3D to 2D Label Transfer," in *Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[22] H. Lim, S. Hwang, and H. Myung, "ERASOR: Egocentric ratio of pseudo occupancy-based dynamic object removal for static 3D point cloud map building," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 2272–2279, 2021.

[23] G. Kim, Y. S. Park, Y. Cho, J. Jeong, and A. Kim, "Mulran: Multimodal range dataset for urban place recognition," in *2020 IEEE Int. Conf. on Robotics and Automat. (ICRA)*, 2020, pp. 6246–6253.