# Image Patch-Matching with Graph-Based Learning in Street Scenes

Rui She, Qiyu Kang, Sijie Wang, Wee Peng Tay, *Senior Member, IEEE,*
Yong Liang Guan, *Senior Member, IEEE,* Diego Navarro Navarro, and Andreas Hartmannsgruber

*Abstract*—Matching landmark patches from a real-time image captured by an on-vehicle camera with landmark patches in an image database plays an important role in various computer perception tasks for autonomous driving. Current methods focus on local matching for regions of interest and do not take into account spatial neighborhood relationships among the image patches, which typically correspond to objects in the environment. In this paper, we construct a spatial graph with the graph vertices corresponding to patches and edges capturing the spatial neighborhood information. We propose a joint feature and metric learning model with graph-based learning. We provide a theoretical basis for the graph-based loss by showing that the information distance between the distributions conditioned on matched and unmatched pairs is maximized under our framework. We evaluate our model using several street-scene datasets and demonstrate that our approach achieves state-of-the-art matching results.

*Index Terms*—Image patch-matching, graph neural network, Kullback-Leibler divergence, information distance maximization, visual place recognition

Fig. 1. Landmark patch-matching using spatial graphs in street scenes and its potential applications.

## I. INTRODUCTION

**A**S a critical and fundamental technique in visual perception, image matching is widely used in many applications, such as image retrieval [1] and vehicle re-identification [2]. Conceptually, the target of a matching task is to solve the similarity correspondence problem for contents from an image pair [3]–[5]. In landmark-based street-scene applications, semantic objects such as traffic signs, traffic lights and road-side poles [6]–[8] often serve as landmarks. The correspondence between the landmark patches captured at different locations may be further utilized as cornerstones to solve other problems, including loop-closure detection in simultaneous localization and mapping (SLAM) [7], [9], place recognition [10], [11], multi-view camera relocalization [12], landmark-LiDAR vehicle relocalization [6], [13], and landmark-based odometry estimation [8].

In traditional image patch-matching methods, handcrafted local features using pixel statistics or gradient information, such as SIFT [14], SURF [15], HOG [16] and ORB [17],

are used. The similarity of a feature pair is commonly computed using different predefined metrics, like the $\mathcal{L}_2$ distance and cosine distance. Moreover, a circular pattern with an adjustable radius is exploited in BRISK [18] and FREAK [19], which provides more efficient neighborhood information for computing relevant pixel statistics. However, these handcrafted features are not robust to viewpoint changes, varying illuminations and transformations. Consequently, the matching performance for methods based on such handcrafted local features is often unstable [20].

With the rapid development of artificial intelligence techniques, deep learning methods, such as convolutional neural networks (CNNs), are widely used in image matching [21]–[23]. In this case, high-dimensional features are exploited to replace handcrafted features in image representations. In the joint feature and metric learning method [24]–[26], the representations and similarity metrics are combined in an end-to-end learning framework, in which high-level features of the images are extracted, and their similarities are learned simultaneously. The feature descriptor learning method [20], [27]–[30] focuses on high-level feature learning and tries to keep matched samples close and unmatched samples far from each other in the corresponding feature space. The similarity is computed using a predefined similarity metric. In these approaches, matching is based on learning feature representations of each image separately and does not exploit the relationships between objects in the images. Recent keypoint-

based learning methods such as D2-Net [31], ASLFeat [32] and SuperGlue [33] perform the point-level correspondence based on the detected keypoints and their descriptors for the input images, which can also be used for the image matching task [34].
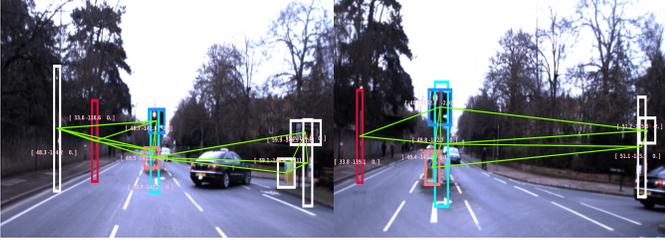


Fig. 2. Landmark patches matching in two full-sized images sampled from the Oxford Radar RobotCar dataset. The matched landmark patches are labeled with the same colored bounding boxes, while the white bounding box indicates that the landmark patch in one image has no matched pair in the other image. Green lines indicate the constructed graph edges in our model.

Unlike other image patch-matching tasks, rich spatial information for landmark patches is often available. For example, lamp posts along a road are usually spaced at equal intervals, and their relative locations with respect to (w.r.t.) each other in the environment provide additional information for the matching task. In special street scenes like the downtown or central business district (CBD), landmark patch-matching has advantages over conventional pixel-/point-level matching due to the presence of dynamic objects, such as vehicles and pedestrians. These dynamic objects captured by the vehicular cameras may have more matched pixels across different frames than static landmarks. However, matching these objects is useless or even harmful for tasks such as place recognition. To mitigate this issue, in this work, we perform the patch-matching task based on static landmarks such as traffic lights, traffic signs, poles, and windows.

Inspired by graph-level representation learning [35], [36], we propose to construct a graph for the neighborhood of an image patch and use graph-level representations to enrich the landmark patch embedding. We identify each landmark patch as a vertex of a graph and find the $K$-nearest neighbors based on estimated spatial information. In the literature, there exist various spatial information estimation techniques like structure-from-motion (SfM) [37], monocular or stereo depth estimation [38] and optical attenuation masks [39]. In this paper, for the sake of illustration, we choose an off-the-shelf monocular depth estimation method from [38] to estimate the landmark spatial relations. However, any other spatial estimation or augmented ranging sensors like LiDAR or depth camera can also be utilized in our framework. We form a clique whose vertex embeddings are learnable via a graph neural network (GNN) [40], [41]. This graph is utilized in our proposed patch-matching framework for object information characterization. The final matching score is an average of the graph and vertex embedding similarity.

We also introduce two landmark patch-matching datasets derived from the street-scene KITTI dataset [42] and the Oxford Radar RobotCar dataset [43]. Our paper focuses on matching image patches of specific *static* roadside objects from two full-sized images taken by cameras onboard vehicles. See

Fig. 1 for an illustration. More specifically, we focus on static roadside objects including traffic lights, signs, lamp posts, and even windows on a building facade. This is because in most landmark-based applications, other transient static objects like parked cars, are inappropriate landmarks or do not have sufficient distinctive features. Due to complex environmental conditions like dynamic element occlusion, e.g., due to pedestrians, vehicles, or the scene viewpoint changing (especially when turning at sharp corners or traversing a stretch in opposite directions), the landmark patches may have dramatic differences in appearance. We refer the readers to the supplementary material for more details on the landmark patch-matching datasets' preparation. For a concrete illustration, some examples of matched or unmatched landmark patches are presented in Fig. 2.

Our contributions are summarized as follows:

- We propose a landmark patch-matching method with graph-based learning for vehicles in street scenes, which extends the feature representation approach used in traditional image patch-matching tasks and incorporates spatial relationship information.
- We analyze the fundamental principle and properties of the proposed graph-based loss function from an information-theoretic perspective.
- We introduce two landmark patch-matching datasets, which contain challenging street-view landmark patches captured in an autonomous driving environment.
- We empirically demonstrate that our method achieves state-of-the-art performance on the landmark patch-matching task when compared to various other benchmarks.

The rest of this paper is organized as follows. In Section II, related works are discussed. Our model and framework are introduced in Section III, where we also provide a theoretical analysis of our graph-based loss. We present experimental results in Section IV and conclude the paper in Section V. The proofs for all lemmas or propositions proposed in this paper are provided in the appendices.

## II. RELATED WORKS

Since deep learning-based methods play dominant roles in the image-matching problem, we only discuss deep learning-based works here. Deep learning-based methods include feature descriptor learning, joint feature and metric learning, as well as keypoint-based correspondence learning.

**Feature descriptor learning.** High-level features of an image are first extracted using a neural network like a CNN so that matched samples are close while unmatched samples are distant under a similarity metric, which is chosen to be a feature distance function. In many models [22], [23], pairwise or triplet loss is used to train the neural networks. To improve performance, in [44], a regularization is designed by maximizing the spread of local feature descriptors over the descriptor space, from which a better embedding for image-level features is obtained. To ensure many samples are accessible to the descriptor network within a few epochs, L2-Net [20] uses a progressive sampling strategy. Furthermore, HardNet [27] is designed to fully utilize the hard negative

samples by making the closest positive sample far away from the closest negative sample in a batch. The reference [28] overcomes the hard sample learning issue by use of exponential Siamese and triplet losses, which naturally pay more attention to hard samples and less attention to easy ones. SOSNet is studied in [45] to learn better local descriptors, where the second-order similarity (SOS) is introduced into the loss function as a regularization. Moreover, [29] designs two second-order components, i.e., the second-order spatial information and the second-order descriptor space similarity, to achieve feature map re-weighting and global descriptors learning, respectively. The paper [46] proposes topology consistent descriptors (TCDesc) based on neighborhood information of descriptors, which can be combined with other methods via the triplet loss.

**Joint feature and metric learning.** In joint feature and metric learning, the similarity metric is not predefined and is instead set as a trainable network together with the feature extraction network. In this case, the matching task is regarded as a binary classification task by resorting to the similarity metric network with a classification loss function. As a classical method, MatchNet proposed by [24] extracts high-level features by using deep CNNs and measures the feature similarity using fully connected (FC) layers. To compare the different network architectures for the matching task, several networks, including SiameseNet, Pseudo-SiameseNet and 2-channel network, are investigated in [21], [47]. The 2-channel network merges the two images into a 2-channel image to achieve faster convergence. The SiameseNet and Pseudo-SiameseNet both use two branches based on the same structure to extract high-dimensional features, with and without the shared weights respectively. Using the normalized cross-correlation (NCC) as a metric, [25] proposes NCC-Net, which utilizes robust matching layers to measure the similarity of feature pairs. To tackle cross-spectral image matching, AFD-Net is proposed by [26] to aggregate multi-level feature differences, which can strengthen the discrimination of the network.

**Keypoint-based correspondence learning.** In keypoint-based correspondence learning, the main procedure is to construct neural networks to perform keypoint detection and description and to measure or learn the keypoints' similarity for matching inference. For instance, LIFT [48] is designed based on a united deep network architecture where keypoints are detected in the first network, the orientation for cropped regions is estimated in the second network, and the feature description is performed in the third network. Here, the Euclidean distance is used to measure the similarity of features. The SuperPoint approach [49] introduces a self-supervised domain adaptation framework named Homographic Adaptation into interest point detection and description. The D2-Net [31] makes use of a single CNN to perform dense feature description and detection simultaneously, where the detection, instead of being based on low-level image structures, is postponed to the high-level structures, which are also used for image descriptions. Based on the D2-Net backbone architecture, ASLFeat [32] is equipped with three lightweight effective modifications, which have better local shape estimation and more accurate keypoint localization. The above methods all measure the point-level correspondence based on Euclidean

distances. On the other hand, SuperGlue [33] is designed using attention GNNs and the Sinkhorn algorithm for keypoint-based feature matching. LoFTR [50] achieves accurate semi-dense matches with Transformers including self and cross-attention layers. Generally speaking, all the above keypoint-based correspondence learning methods can be used to perform the image matching task with further operations on the keypoint matching scores [33].

To improve image matching performance, spatial information is used in [33], [50], [51] through spatial verification, graph learning, and cross attention. In spatial verification, spatial information is usually used for the transformation calibration w.r.t. the key points or objects, as well as a correspondence auxiliary for direction or location w.r.t. the objects of interest [51]. This can introduce global information to improve local correspondence. In particular, transformation optimization methods like RANSAC [52], fast spatial measure (FSM) [53], hough pyramid matching (HPM) [54] and pairwise geometric matching (PGM) [55], can filter out weak correspondences for keypoints or local features obtained by key feature detection and descriptors such as SIFT [14], SURF [15], and ORB [17]. The region-based or object-based verification methods such as Objects in Scene to Objects in Scene (OS2OS) [51] and block-based image matching [56], make use of the relative positions of local patches to refine the whole image matching. Different from the above approaches, our method uses distance-based spatial information for the neighborhood graph construction, rather than for transformation correction or weak correspondence filtering.

Graph learning methods such as SuperGlue [33], GLMNet [57], and joint graph learning and matching network (GLAM) [58], are exploited to represent local features based on the neighborhood graphs for keypoints. The graphs are constructed based on the detected keypoints or the corresponding features, and GNNs are used to learn graph representations. These methods achieve more robust and stable representations for the corresponding features based on spatial information.

Different from the above methods, our approach focuses on the *neighborhood information based on landmark distances*, which is used for patch-level, rather than point-level, representation and not used to filter weak or invalid correspondences. Moreover, we also adopt GNNs to represent the patch-level neighborhood graphs, which is demonstrated to be beneficial for the landmark patch-matching task.

## III. LANDMARK PATCH-MATCHING WITH GRAPH-BASED LEARNING

In this section, we first introduce our graph-based learning framework to find matched landmark patch pairs that are extracted from two images taken from on-vehicle cameras. The images may be taken from different perspectives and our framework can also identify those patches that are unmatched. Fig. 2 shows examples of matched and unmatched landmark patch pairs. We then discuss the theoretical basis for our graph-based learning approach.
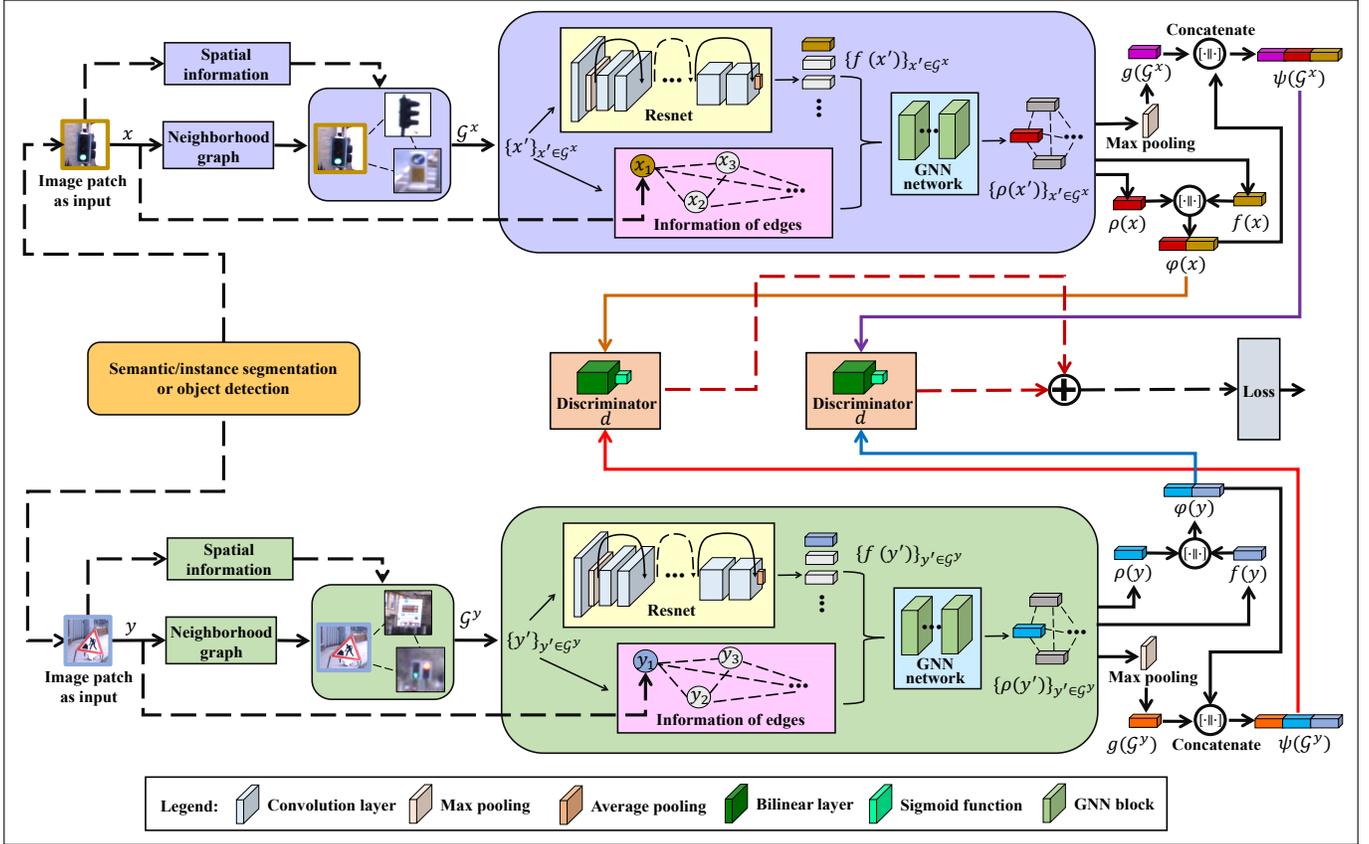
Fig. 3. VGIDM: landmark patch-matching with the graph-based learning. The Resnet $f$ shown in the framework is a shared network serving as the feature descriptor function $f$ to extract high-dimensional features from patches. Likewise, the discriminator $d$ is also shared to make a decision for the vertex-to-graph correspondence. The model takes as input a pair of image patches that correspond to street scene landmarks.

## A. Framework Overview

Similar to other patch-matching datasets like the multi-view stereo (MVS) dataset [59] and the DTU dataset [60], in our work, the landmark patches are extracted from the full-sized images and the matching ground truths are established using 3D points. More specifically, the landmark patches are extracted using well-known object detection techniques like Faster R-CNN [61]. To distinguish the full-sized images from the landmark patches, we use the term *frame* to denote the full-sized image from which the patches are extracted. We refer the readers to Section IV-A for more details on the preparation of the landmark patch-matching datasets.

We assume that the spatial information (i.e., approximate relative distances between landmark objects) of landmark patches is available. The spatial information can be obtained from range estimation methods like the monocular depth estimation networks [62]–[64] in both the training and the testing phases. To construct a graph, we let the landmark patches of a frame be vertices of the graph. For each patch or vertex $x$, we find the $K$ nearest neighbors in terms of spatial locations as indicated by the observed spatial information. An example of the constructed graph is shown in Fig. 2. For the vertex $x$, we form a complete graph or clique with its $K$ nearest neighbors found. Let $\mathcal{G}^x$ denote this neighborhood graph.

Our image patch-matching framework is illustrated in Fig. 3. In this framework, the inputs for the model are image patches obtained by semantic or instance segmentation methods, e.g.,

Mask R-CNN [65], or object detection methods, e.g., Faster R-CNN [61]. These two kinds of methods can extract objects of interest such as traffic lights and traffic signs from image frames. Two main modules, Resnet [66] and GNN, are respectively used for image feature extraction and neighborhood graph embedding, where the GNN can be the graph attention network (GAT) [40], graph convolutional network (GCN) [67], GraphSAGE [68] or any other GNN architecture. Given the vertex and graph embedding features from our model, we maximize the empirical information distance between the cases where patches are matched and unmatched. We call our image patch-matching approach *Vertex-Graph-learning-and-Information-Distance-Maximization (VGIDM)*. The details are given as follows.

## B. Model Details

Our objective is to determine if two landmark patches from different frames are matched with each other. In VGIDM, the feature extraction module $f$ first learns embeddings for the input landmark patch as well as the patches in its neighborhood graph. The model then makes use of a learnable graph embedding module $g$ to represent the neighborhood graph-level and vertex-level feature readout features. Finally, it uses a decision-making module to compute the matching classification.

**Feature extraction for patches.** We use the Resnet $f$ to extract high-dimensional features for each landmark patch $x$. The Resnet output is denoted by $f(x) \in \mathbb{R}^n$. Recall that for a

patch $x$, we form a neighborhood graph $\mathcal{G}^x$. For the graph $\mathcal{G}^x$, applying $f$ to each node in $\mathcal{G}^x$, we have $\{f(x')\}_{x' \in \mathcal{G}^x}$.

**Embedding representation for the neighborhood graph and its vertices.** The graph $\mathcal{G}^x$ is input to a GNN network $g$ to obtain a graph-level embedding representation $g(\mathcal{G}^x)$. Specifically, the vertex features $(f(x'))_{x' \in \mathcal{G}^x} \in \mathbb{R}^{|\mathcal{G}^x| \times n}$ are updated via the GNN, which consists of several layers of neighborhood aggregation and node update [40], [41], followed by some activation functions and a final pooling layer. The vertex embeddings $(\rho(x'))_{x' \in \mathcal{G}^x} \in \mathbb{R}^{|\mathcal{G}^x| \times n}$ are obtained from the last graph convolutional/attentional layer of the GNN, while the graph-level embedding representation $g(\mathcal{G}^x) \in \mathbb{R}^n$ is obtained as the output of the last pooling layer.

Compared to $f(x) \in \mathbb{R}^n$ which extracts features directly from the patch $x$, $\rho(x) \in \mathbb{R}^n$ learns a feature embedding with additional information from its neighborhood, while $g(\mathcal{G}^x) \in \mathbb{R}^n$ learns an embedding for the surrounding environment itself.

**Correspondence comparison.** Suppose that $x$ and $y$ are landmark patches from two different frames, respectively. If $x$ and $y$ are patches for the same real-world object, we say that they are *matched* and denote this event as $x \leftrightarrow y$. Otherwise, they are *unmatched* and denoted as $x \nleftrightarrow y$. For any patch pair $(x, y)$, we denote the matching ground truth label as $\mathbf{1}_{\{x \leftrightarrow y\}}$, where $\mathbf{1}_{\{\cdot\}}$ is the indicator function. In order to compare the correspondence between the patch pair $(x, y)$, we design a decision-making mechanism based on the patch features. For the two patches $x$ and $y$, we respectively obtain $f(x)$ and $f(y)$ as the features from the Resnet, $\rho(x)$ and $\rho(y)$ as the vertex-level embedding features, and $g(\mathcal{G}^x)$ and $g(\mathcal{G}^y)$ as the graph-level embedding features from the GNN network.

Let the ensemble vertex embedding for a patch $x$ be

$$\varphi(x) = \rho(x) \| f(x) \tag{1}$$

and the neighborhood graph embedding for $\mathcal{G}^x$ be

$$\psi(\mathcal{G}^x) = g(\mathcal{G}^x) \| \varphi(x) = g(\mathcal{G}^x) \| \rho(x) \| f(x), \tag{2}$$

where $\|$ is the concatenation operation.

The ensemble vertex feature for $x$ and the graph embedding for $\mathcal{G}^y$ are input to a discriminator $d$ consisting of a *bilinear* layer of the form:

$$d(a, b) = \sigma(a^\top \times M \times b), \tag{3}$$

where $M \in \mathbb{R}^{n \times m}$ is a trainable matrix and $\sigma(\cdot)$ denotes the sigmoid function. In particular, the matrix $M$ is designed as

$$M = \begin{bmatrix} \mathbf{0} & M_{12} & \mathbf{0} \\ M_{21} & M_{22} & M_{23} \end{bmatrix}, \tag{4}$$

where $M_{12}$, $M_{21}$, $M_{22}$ and $M_{23}$ serve as the matrix blocks with learnable parameters and $\mathbf{0}$ denotes the zero matrix. The specifically designed block matrix (4) is to restrict the

comparison between the features. Inputting $(\varphi(x), \psi(\mathcal{G}^y))$ to the discriminator $d$, we have

$$d(\varphi(x), \psi(\mathcal{G}^y))$$

$$= \sigma \left( \begin{bmatrix} \rho(x)^\top & f(x)^\top \end{bmatrix} \times M \times \begin{bmatrix} g(\mathcal{G}^y) \\ \rho(y) \\ f(y) \end{bmatrix} \right) \tag{5}$$

$$= \sigma \Big( \rho(x)^\top M_{12} \rho(y) + f(x)^\top M_{21} g(\mathcal{G}^y)$$

$$+ f(x)^\top M_{22} \rho(y) + f(x)^\top M_{23} f(y) \Big). \tag{6}$$

The first term $\rho(x)^\top M_{12} \rho(y)$ in (6) is used to compare the vertex embeddings of $x$ and $y$ obtained from the GNN. This emphasizes the domain part of the embedding. The second term $f(x)^\top M_{21} g(\mathcal{G}^y)$ and third term $f(x)^\top M_{22} \rho(y)$ are used for the comparison of the vertex $x$ and the neighborhood graph of $y$. This helps to constrain GNN learning. The last term $f(x)^\top M_{23} f(y)$ is used to compare the Resnet features for the two vertices, which updates the Resnet training. The same procedure is performed analogously for $\varphi(y) = \rho(y) \| f(y)$ and $\psi(\mathcal{G}^x) = g(\mathcal{G}^x) \| \varphi(x)$.

The learnable discriminator $d(\varphi(x), \psi(\mathcal{G}^y))$ from (6) utilizes the ensemble vertex embedding $\varphi(x)$ and the neighborhood graph embedding $\psi(\mathcal{G}^y)$. The vertex-level embedding $\rho(x)$ and graph-level embedding $g(\mathcal{G}^x)$ contain information from the vertex feature $f(x)$ (output of Resnet) due to the incorporation of neighborhood information from the GNN. In the case of a large number of frames, the neighborhood graphs can be quite different as they typically consist of vertices from different frames. As a result, the embeddings $\rho(x), g(\mathcal{G}^x)$ and $\rho(y), g(\mathcal{G}^y)$ can have different features to some degree even if $x \leftrightarrow y$. Therefore, it may be appropriate to use the original vertex feature $f(x)$ to constrain the graph learning for the vertex-level embedding $\rho(y)$ and the graph-level embedding $g(\mathcal{G}^y)$. The comparisons between $f(x)$ and $\rho(y)$ or $g(\mathcal{G}^y)$ can emphasize the principal component for the learned graph features. When vertices $x$ and $y$ are matched, $\rho(y)$ and $g(\mathcal{G}^y)$ essentially contain the information of $f(x)$. Therefore, comparing $f(x)$ with $\rho(y)$ and $g(\mathcal{G}^y)$ can introduce more information with neighborhood characteristics for the matching process.

**Loss function and matching score.** Let $\mathcal{M}$ be a training set consisting of patch pairs $(x, y)$. Define the graph-based learning objective function as $L_{\mathrm{empID}}$ given in (8), which depends on the discriminator $d$ in (6). We show that $L_{\mathrm{empID}}$ is the empirical version of an information distance between the distributions conditioned by matched and unmatched pairs in Proposition 1. We set our overall loss as

$$\min_{\varphi, \psi, d} \{-L_{\mathrm{empID}}\}, \tag{9}$$

to maximize the information distance.

In the testing phase, the final matching score is given by

$$S_{\mathrm{match}}(x, y) = \frac{d(\varphi(x), \psi(\mathcal{G}^y)) + d(\varphi(y), \psi(\mathcal{G}^x))}{2}, \tag{10}$$

$$L_{\text{empID}} = \frac{1}{|\mathcal{M}|} \sum_{(x,y)\in\mathcal{M}} \left\{ \mathbf{1}_{\{x\leftrightarrow y\}} \frac{1}{2}\Big( \log[d(\varphi(x),\psi(\mathcal{G}^y))] + \log[d(\varphi(y),\psi(\mathcal{G}^x))] \Big) \right.$$

$$\left. + \mathbf{1}_{\{x\not\leftrightarrow y\}} \frac{1}{2}\Big( \log[1 - d(\varphi(x),\psi(\mathcal{G}^y))] + \log[1 - d(\varphi(y),\psi(\mathcal{G}^x))] \Big) \right\} \tag{7}$$

$$= \frac{1}{2} \underbrace{\frac{1}{|\mathcal{M}|} \sum_{(x,y)\in\mathcal{M}} \left\{ \mathbf{1}_{\{x\leftrightarrow y\}} \log[d(\varphi(x),\psi(\mathcal{G}^y))] + \mathbf{1}_{\{x\not\leftrightarrow y\}} \log[1 - d(\varphi(x),\psi(\mathcal{G}^y))] \right\}}_{L_{\text{empID}-1}}$$

$$+ \frac{1}{2} \underbrace{\frac{1}{|\mathcal{M}|} \sum_{(x,y)\in\mathcal{M}} \left\{ \mathbf{1}_{\{y\leftrightarrow x\}} \log[d(\varphi(y),\psi(\mathcal{G}^x))] + \mathbf{1}_{\{y\not\leftrightarrow x\}} \log[1 - d(\varphi(y),\psi(\mathcal{G}^x))] \right\}}_{L_{\text{empID}-2}} \tag{8}$$

and the prediction function for whether there is a match is given by

$$A_S^{\text{test}}(x,y) = \begin{cases} 1, & \text{if } S_{\text{match}}(x,y) > \Gamma, \\ 0, & \text{otherwise}, \end{cases} \tag{11}$$

where $\Gamma$ is a predefined threshold. A decision "1" indicates that $x$ and $y$ are matched and "0" otherwise.

### C. Theoretical Basis

In this subsection, we discuss the theoretical basis for the graph-based learning objective function $L_{\text{empID}}$ defined in (8). To make the analysis tractable, we assume that patch pairs $(x,y)$ are randomly generated from a distribution $\mathbb{P}$. Let $\mathbb{E}$ be the expectation operator. We start with a simplifying assumption as follows.

**Assumption 1.** $\varphi(x)$ and $\psi(\mathcal{G}^y)$ are continuous random variables induced from $\mathbb{P}$.

In practice, due to the chosen activation functions used in Resnet $f$ and the GNN network $g$, their outputs typically satisfy the continuity requirement of Assumption 1.

In our analysis, the discriminator $d$ is assumed to be a general function without necessarily having the form (3).

Let $\mathbb{A}$ be the set of all possible $(\varphi(x),\psi(\mathcal{G}^y))$ where $\int_{\mathbb{A}} p(\varphi(x),\psi(\mathcal{G}^y))\mathrm{d}(\varphi(x),\psi(\mathcal{G}^y)) = 1$, $p : \mathbb{A} \mapsto \mathbb{R}_+$ is a probability density whose set of discontinuities has Lebesgue measure zero.

For any given landmark patches $x$ and $y$, we assume that $\mathbb{P}(x\leftrightarrow y) > 0$ and $\mathbb{P}(x\not\leftrightarrow y) > 0$. The probability densities of $(\varphi(x),\psi(\mathcal{G}^y))$ conditioned on $x\leftrightarrow y$ and $x\not\leftrightarrow y$ are denoted by $p(\varphi(x),\psi(\mathcal{G}^y) \mid x\leftrightarrow y)$ and $p(\varphi(x),\psi(\mathcal{G}^y) \mid x\not\leftrightarrow y)$, respectively.[1]

We discuss only $L_{\text{empID}-1}$ in (8) since $L_{\text{empID}-2}$ is symmetrical to it. The expectation form of $L_{\text{empID}-1}$ is given by

$$L_{\text{ID}} = L_{\text{ID}}(\varphi,\psi,d)$$
$$= \mathbb{E}\big[\mathbf{1}_{\{x\leftrightarrow y\}} \log d(\varphi(x),\psi(\mathcal{G}^y))\big]$$
$$+ \mathbb{E}\big[\mathbf{1}_{\{x\not\leftrightarrow y\}} \log(1 - d(\varphi(x),\psi(\mathcal{G}^y)))\big]. \tag{12}$$

In minimizing the loss in (9), in the asymptotic regime $|\mathcal{M}| \to \infty$, we aim at $\max_{\varphi,\psi,d} L_{\text{ID}}$. Let $D(\cdot \| \cdot)$ denote the Kullback-Leibler (KL) divergence.

**Proposition 1** (Relationship with KL divergence). *Suppose Assumption 1 holds. For a vertex embedding $\varphi$ and a neighborhood graph embedding $\psi$, let $L_{\text{ID}}^{d^*}(\varphi,\psi) = \max_d L_{\text{ID}}(\varphi,\psi,d)$, where $d^*$ is the corresponding optimal discriminator. Then*

$$D(p(\varphi(x),\psi(\mathcal{G}^y) \mid x\leftrightarrow y) \| p(\varphi(x),\psi(\mathcal{G}^y) \mid x\not\leftrightarrow y)) \tag{13}$$
$$\geq \frac{1}{\mathbb{P}(x\leftrightarrow y)}\Big( L_{\text{ID}}^{d^*}(\varphi,\psi) + H_b(\mathbb{P}(x\leftrightarrow y)) \Big). \tag{14}$$

*where $H_b(p) = -p\log p - (1-p)\log(1-p)$ is the binary entropy function.*

*Proof.* See Appendix A. $\square$

**Remark 1.** *Proposition 1 suggests that maximizing $L_{\text{ID}}$ over $(\varphi,\psi,d)$ helps to distinguish between the matched and unmatched patch pairs since their conditional distributions are forced to be very different in terms of the KL divergence.*

We next consider how the graph-based learning objective function $L_{\text{ID}}$ in (12) is influenced by perturbations in the discriminator $d$.

**Proposition 2** (Effect of discriminator perturbation). *Suppose Assumption 1 holds. Let $\varepsilon$ be a sufficiently small perturbation to the discriminator $d$. Then, $|L_{\text{ID}}(\varphi,\psi,d+\varepsilon) - L_{\text{ID}}(\varphi,\psi,d)| = O(\varepsilon)$. Furthermore, we have $|\max_d L_{\text{ID}}(\varphi,\psi,d+\varepsilon) - \max_d L_{\text{ID}}(\varphi,\psi,d)| = O(\varepsilon^2)$.*

*Proof.* See Appendix B. $\square$

In the following, we consider how the GNN embedding of the neighborhood graph $\mathcal{G}^x$ of a vertex $x$ affects the matching effectiveness under further assumptions.

For two landmark patches $x$ and $y$, if their neighborhood graphs $\mathcal{G}^x$ and $\mathcal{G}^y$ have vertices corresponding to the same set of objects, i.e., the patch and spatial information procedure identifies the same objects as the neighbors of $x$ and $y$, we write $\mathcal{G}^x \leftrightarrow \mathcal{G}^y$.

**Assumption 2.** *The ranges of $\varphi(\cdot)$ and $\psi(\cdot)$ are finite sets. The embedding $\varphi(x) = \varphi(y)$ for landmark patches $x$ and $y$ are the same if $x\leftrightarrow y$. If furthermore $\mathcal{G}^x \leftrightarrow \mathcal{G}^y$, then $\psi(\mathcal{G}^x) = \psi(\mathcal{G}^y)$.*

---

[1] Here we abuse notations $p(\varphi(x),\psi(\mathcal{G}^y)|x\leftrightarrow y)$ to denote the conditional probability density of $(\varphi(x),\psi(\mathcal{G}^y))$ given that $x$ and $y$ are matched. This is to avoid the cluttered notation $p_{(\varphi(x),\psi(\mathcal{G}^y))|x\leftrightarrow y}(\cdot,\cdot)$.

While the Resnet $f$ and GNN block $g$ are in general continuous functions of their inputs, Assumption 2 can be satisfied by restricting to a finite number of objects of interest in the environment, assuming that frames are captured from approximately the same perspectives (e.g., from an on-vehicle camera of a vehicle traveling along a fixed road) so that landmark patches of the same object are within a certain similarity distance of each other. Finally, the outputs of $f$ and $g$ can be quantized into discrete ranges, which implies $\varphi$ and $\psi$ have finite sets of ranges. For the same object $o$ in the environment but under two different frames $\mathcal{F}_1$ and $\mathcal{F}_2$, Assumption 2 says that the outputs from the embedding $\varphi$ are the same for the two frames. This implicitly assumes that $\varphi$ is robust to perturbation in its input. Furthermore, the outputs of the embedding $\psi$ are also the same if the patch and spatial information are noiseless.

**Proposition 3.** *Suppose Assumption 2 holds, and $x$ and $y$ are landmark patches of frames $\mathcal{F}_1$ and $\mathcal{F}_2$ (based on the same environment), respectively. Let $m(x \leftrightarrow y) = \mathbb{P}(\mathcal{G}^x \leftrightarrow \mathcal{G}^y \mid x \leftrightarrow y)$ and $m(x \not\leftrightarrow y) = \mathbb{P}(\mathcal{G}^x \not\leftrightarrow \mathcal{G}^y \mid x \not\leftrightarrow y)$. Then we have*

$$\|p(\varphi(x), \psi(\mathcal{G}^y) \mid x \leftrightarrow y) - p(\varphi(x), \psi(\mathcal{G}^y) \mid x \not\leftrightarrow y)\|_{\mathrm{TV}}$$
$$\geq \min_{x \leftrightarrow y} m(x \leftrightarrow y) \sum_{(\varphi(x), \psi(\mathcal{G}^y)) \in \mathbb{B}} \Big\{ p(\varphi(x), \psi(\mathcal{G}^y) \mid \mathcal{G}^x \leftrightarrow \mathcal{G}^y, x \leftrightarrow y)$$
$$- p(\varphi(x), \psi(\mathcal{G}^y) \mid \mathcal{G}^x \not\leftrightarrow \mathcal{G}^y, x \not\leftrightarrow y) \Big\}$$
$$+ \min_{x \leftrightarrow y} m(x \leftrightarrow y) - \max_{x \not\leftrightarrow y} m(x \not\leftrightarrow y) - 1, \qquad (15)$$

*where $\mathbb{B} = \big\{ (\varphi(x), \psi(\mathcal{G}^y)) : p(\varphi(x), \psi(\mathcal{G}^y) \mid x \leftrightarrow y) \geq p(\varphi(x), \psi(\mathcal{G}^y) \mid x \not\leftrightarrow y) \big\}$. Here, $\|\cdot\|_{\mathrm{TV}}$ denotes the total variation distance, and $\min_{x \leftrightarrow y}$ and $\max_{x \not\leftrightarrow y}$ denotes minimization over all matched patch pairs $(x, y)$ and maximization over all unmatched patch pairs, respectively.*

*Proof.* See Appendix C. □

In the ideal case where the patch and spatial information are noiseless, we have $\min_{x \leftrightarrow y} m(x \leftrightarrow y) = 1$ and $\max_{x \not\leftrightarrow y} m(x \not\leftrightarrow y) = 0$. Then the right-hand side of (15) simplifies to

$$\sum_{(\varphi(x), \psi(\mathcal{G}^y)) \in \mathbb{B}} \Big\{ p(\varphi(x), \psi(\mathcal{G}^y) \mid \mathcal{G}^x \leftrightarrow \mathcal{G}^y, x \leftrightarrow y)$$
$$- p(\varphi(x), \psi(\mathcal{G}^y) \mid \mathcal{G}^x \not\leftrightarrow \mathcal{G}^y, x \not\leftrightarrow y) \Big\}. \qquad (16)$$

In this case, we also have $p(\varphi(x), \psi(\mathcal{G}^y) \mid x \leftrightarrow y) = p(\varphi(x), \psi(\mathcal{G}^y) \mid \mathcal{G}^x \leftrightarrow \mathcal{G}^y, x \leftrightarrow y)$ and $p(\varphi(x), \psi(\mathcal{G}^y) \mid x \not\leftrightarrow y) = p(\varphi(x), \psi(\mathcal{G}^y) \mid \mathcal{G}^x \not\leftrightarrow \mathcal{G}^y, x \not\leftrightarrow y)$. Furthermore, from Assumption 2, any $(\varphi(x), \psi(\mathcal{G}^y))$ such that $p(\varphi(x), \psi(\mathcal{G}^y) \mid \mathcal{G}^x \leftrightarrow \mathcal{G}^y, x \leftrightarrow y) > 0$ implies that $p(\varphi(x), \psi(\mathcal{G}^y) \mid \mathcal{G}^x \not\leftrightarrow \mathcal{G}^y, x \not\leftrightarrow y) = 0$. These probability measures are thus mutually singular and have a total variation distance of 1. Therefore, in the ideal case, the model perfectly distinguishes between $p(\varphi(x), \psi(\mathcal{G}^y) \mid x \leftrightarrow y)$ and $p(\varphi(x), \psi(\mathcal{G}^y) \mid x \not\leftrightarrow y)$.

## IV. EXPERIMENTS

### A. Datasets

As there are no existing standard datasets for street-scene landmark patch-matching, we introduce in this paper two new

datasets: the Landmark KITTI dataset and the Landmark Oxford dataset,[2] which are derived from the street-scene KITTI dataset [42] and the Oxford Radar RobotCar dataset [43], respectively.

Both datasets contain image frames and LiDAR scans captured from onboard cameras and Velodyne LiDAR sensors. The landmark patches are extracted from the full-sized image frames using the object detection neural network Faster R-CNN [61]. To facilitate detection efficacy, we manually label several street-scene compact landmark objects including traffic lights, traffic signs, poles, and facade windows for the sampled frames. The labels are used to train Faster R-CNN, which is used to produce landmark object detection for the image frames. The detected landmarks in bounding boxes are then used to obtain the landmark patches for our matching experiments with some intentionally included background, shown in Fig. 4 for example.



**(a) Landmark patches from KITTI Dataset**



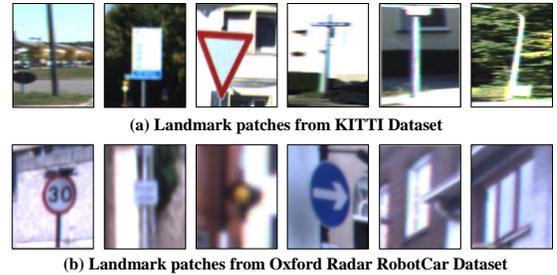**(b) Landmark patches from Oxford Radar RobotCar Dataset**

Fig. 4. (a) and (b) are landmark patch samples (displayed with intentionally included background) from the KITTI dataset and Oxford Radar RobotCar dataset respectively.

To establish the patch-matching ground truth, we use the vehicle locations and collected LiDAR scans to build the 3D LiDAR reference map similar to the operations in [69]. The 3D reference map is used to determine the landmark locations by projecting the 3D LiDAR points to the image frames. The LiDAR points reflected from the landmark patch are read out to get the global locations of the corresponding landmark objects. We then compute the $\mathcal{L}_2$ distance of each landmark patch pair from two frames to determine the patch-matching ground truth. Some details of the two landmark patch-matching datasets are introduced as follows. More dataset preparation details are given in the supplementary material.

**Landmark KITTI Dataset.** The KITTI dataset[3] contains street-scene image frames and their corresponding LiDAR point clouds collected in Karlsruhe, Germany. We use the object labels provided by [70] to detect landmark patches for all frames including traffic lights, traffic signs and poles. An example is shown in Fig. 5. We do not include windows as landmarks in this dataset due to the lack of labels. Furthermore, to avoid "trivial matchings" between consecutive images, a minimum difference of 2m between the image frames is also set. The aforementioned operations are performed to obtain the landmark patch-matching ground truth by projecting the 3D LiDAR scans to the image frames. Finally, 1500 frames are selected for landmark patch-matching experiments. The dataset is randomly split into training and testing sets, with a

---

ratio around 2 : 1. In both training and testing, we select frame pairs that are captured at locations with relative distances not more than 25m to ensure the presence of common landmarks.
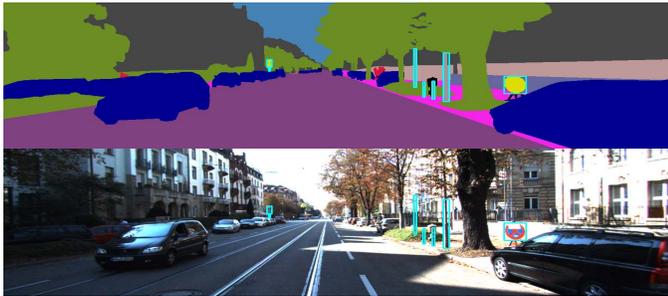


Fig. 5. A semantic segmentation image and its corresponding real image, both with bounding box labels, from the KITTI dataset.

**Landmark Oxford Dataset.** The Oxford Radar RobotCar dataset[4] contains image frames and LiDAR scans captured on the streets in Oxford, UK. We manually label landmarks including traffic lights, traffic signs, poles, and facade windows for 500 sampled frames. An example is shown in Fig. 6. We then train Faster R-CNN to obtain the landmarks for all $29,687$ frames. To avoid "trivial matchings" between consecutive images, a minimum difference of 2m between the image frames is also set. Finally, 3000 frames are selected for landmark patch-matching experiments. The remaining steps are similar to that for the Landmark KITTI dataset.
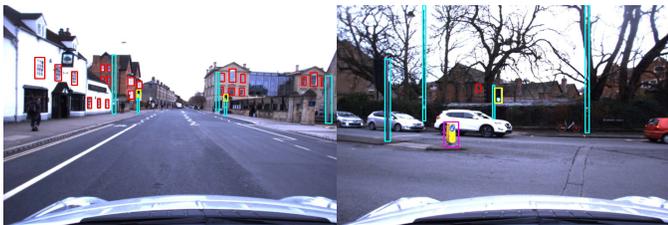


Fig. 6. Examples of the ground truth landmark bounding box labels for the Oxford Radar RobotCar dataset.

*B. Experimental Details*

**Baseline Methods.** We compare VGIDM with several baseline methods, including MatchNet [24], SiameseNet [47], HardNet [27], SOSNet [45], D2-Net [31], ASLFeat [32], SuperGlue [33] and LoFTR [50]. The MatchNet and SiameseNet are regarded as joint feature and metric learning methods, combining deep CNNs and an FC layer to learn features and their metrics. The decision-making process for the matching task is based on the output of the FC layer. HardNet and SOSNet focus on similarity measures to distinguish the learned high-dimensional features, where the feature descriptors are almost all based on deep CNNs consisting of several convolution layers with batch normalization (BN) or rectified linear units (ReLUs). In testing, the Euclidean distance between the output patch features is used for the decision-making. D2-Net, ASLFeat, SuperGlue and LoFTR are based on keypoint correspondence and perform the matching task according to the ratio of the matched keypoints among the whole set of

[4]http://ori.ox.ac.uk/datasets/radar-robotcar-dataset

keypoints. In this regard, a patch pair with a large enough ratio of matched keypoints is regarded as a match.

**Model Setting.** We use Resnet18 in [66] for the feature descriptor $f$, with output feature dimension 512 after 17 convolution layers. In VGIDM, we choose several GNNs for the neighborhood graph embedding, including GAT [40], GCN [67] and GraphSAGE [68]. When using GAT, the network contains 2 GAT blocks with the exponential linear unit (ELU). For each GAT block, we use 4 attention heads, which compute 512 hidden features in total. As for GCN and GraphSAGE, they both contain 2 corresponding blocks with ReLU, where there are 512 hidden features in each block. Further details of our model architecture are provided in the supplementary material. The Adam optimizer is selected with a learning rate of 0.0001 to train the model by minimizing its corresponding loss in (9). The number of training epochs is 150 for all datasets.

**VGIDM with Image Depth Estimation.** To test VGIDM in the case where precise depth information like that provided by LiDAR is unavailable, we construct neighborhood graphs using estimated image depth and with different GNNs in the backbone. Specifically, we include an image depth estimation method called Monocular Depth Prediction Module proposed in [38]. Based on the image depth estimation, we can obtain the rough relative locations of the landmarks in the street scenes and use them to construct a neighborhood graph for each landmark in the test procedure. The depth estimation performance is provided in the supplementary material. The estimated image pixel depths are transformed to 3D locations w.r.t. the camera using its intrinsic matrix. We then use the estimated locations to test VGIDM. In this depth estimation method, the pre-trained ResNeXt101 model from [38] is utilized in our experiments, and the images are from the two landmark datasets. We extract the predicted depth points from the static roadside landmarks, including traffic lights, traffic signs, and poles, to compute the locations of the objects. Therefore, we can construct the neighborhood graphs and test the VGIDM.

**Implementation.** For a given sequence of street scene frames captured by a vehicular camera, we perform the following training steps: i) We use object detection methods like faster R-CNN [61] to extract landmark patches for each frame. The landmarks include traffic lights, traffic signs, poles, and windows. ii) We manually label matching landmark patches. To determine the global locations of these landmarks, we combine vehicular Global Positioning System (GPS) information with data from LiDAR or stereo cameras. With this information, we are able to establish the ground truth for the matching landmark patches between two frames captured at the same location. iii) We take the global locations of landmarks to construct the neighborhood graph for each landmark patch based on $K$-NN. iv) We train the VGIDM using landmark patch pairs with ground-truth labels. The details of VGIDM with training loss and test score are given in Section III-B.

During testing, we perform steps i and iii as above but in step iii, we create neighborhood graphs by estimating the relative locations of landmarks using a stereo camera or a depth estimation method, which replaces the need for GPS and LiDAR information. The ground truth for computing the testing performance is found based on GPS and LiDAR information.

TABLE I

Matching performance on the Landmark KITTI dataset. The best and the second-best result for each criterion are highlighted in RED and BLUE respectively.

| Methods | Precision | Recall | $F_1$-Score | AUC |
|---------|-----------|--------|-------------|-----|
| MatchNet [24] | $0.9039 \pm 0.0027$ | $0.9483 \pm 0.0105$ | $0.9255 \pm 0.0050$ | $0.8229 \pm 0.0055$ |
| SiameseNet [47] | $0.7953 \pm 0.0124$ | $0.8960 \pm 0.0208$ | $0.8426 \pm 0.0159$ | $0.8328 \pm 0.0162$ |
| HardNet [27] | $0.9041 \pm 0.0016$ | $0.9562 \pm 0.0177$ | $0.9294 \pm 0.0093$ | $0.8261 \pm 0.0088$ |
| SOSNet [45] | $0.9042 \pm 0.0015$ | $0.9563 \pm 0.0160$ | $0.9294 \pm 0.0083$ | $0.8261 \pm 0.0080$ |
| D2-Net [31] | $0.9115 \pm 0.0031$ | $0.8789 \pm 0.0131$ | $0.8949 \pm 0.0076$ | $0.8115 \pm 0.0082$ |
| ASLFeat [32] | $0.9189 \pm 0.0022$ | $0.9008 \pm 0.0082$ | $0.9098 \pm 0.0048$ | $0.8312 \pm 0.0057$ |
| SuperGlue [33] | $0.9067 \pm 0.0039$ | $0.9125 \pm 0.0123$ | $0.9096 \pm 0.0072$ | $0.8155 \pm 0.0093$ |
| LoFTR [50] | $0.9069 \pm 0.0025$ | $0.9243 \pm 0.0110$ | $0.9154 \pm 0.0059$ | $0.8197 \pm 0.0064$ |
| VGIDM (GAT) [ours] | $0.9425 \pm 0.0020$ | $0.9733 \pm 0.0050$ | $0.9577 \pm 0.0026$ | $0.8977 \pm 0.0038$ |
| VGIDM (GCN) [ours] | $0.9340 \pm 0.0027$ | $0.9753 \pm 0.0076$ | $0.9543 \pm 0.0044$ | $0.8847 \pm 0.0063$ |
| VGIDM (GraphSAGE) [ours] | $0.9464 \pm 0.0042$ | $0.9653 \pm 0.0129$ | $0.9557 \pm 0.0073$ | $0.9007 \pm 0.0098$ |



Fig. 7. Examples of matched and mismatched pairs from the Landmark KITTI dataset. A green or red box indicates a correct or incorrect prediction result, respectively. "GT" stands for ground truth.

The remaining steps are the same as those used during training.

## C. Performance Evaluation

**Performance on Landmark KITTI Dataset.** Table I summarizes the test performance of models trained with 150 training epochs on the Landmark KITTI dataset. The evaluation uses statistics including mean value and standard deviation from 5 experiments. From Table I, we observe that VGIDM (with GAT, GCN or GraphSAGE) outperforms the other baseline methods under all four criteria, with a slight performance difference among these VGIDM models. This implies that graph-based learning makes a positive difference in matching efficiency. Moreover, we observe that VGIDM with GAT has a more stable performance than the other methods. Several examples of the matching prediction are shown in Fig. 7.

**Performance on Landmark Oxford Dataset.** From Table II, we observe that the VGIDM variants with different GNNs outperform the other benchmark methods on almost all measures. Since the Oxford Radar RobotCar and KITTI datasets have different image qualities and are collected in

different street scenes, the performances on both datasets are different. From Tables I and II, we also observe that nearly all the methods have better performance on the Landmark KITTI dataset compared with the Landmark Oxford dataset. This may be caused by more similarity among the window patches in the Landmark Oxford dataset, which makes distinguishing them more difficult. A few matching prediction examples from the Landmark Oxford dataset are shown in Fig. 8.

**Performance Analysis.** The VGIDM variants (with GAT, GCN or GraphSAGE) not only make use of landmark patch information but also the neighborhood information in the decision-making process for the matching task. Other feature descriptor learning as well as joint feature and metric learning methods such as MatchNet, SiameseNet, HardNet and SOSNet, depend only on the individual image patch rather than the neighborhood relationships. An erroneous match can happen between patches from two similar but distinct objects. VGIDM mitigates this error by using the neighborhood information. However, VGIDM requires more computing resources for neighborhood graph processing.

On the other hand, keypoint-based learning methods such as D2-Net, ASLFeat, SuperGlue and LoFTR, suffer from low

TABLE II

Matching performance on the Landmark Oxford dataset. The best and the second-best result for each criterion are highlighted in RED and BLUE respectively.

| Methods | Precision | Recall | $F_1$-Score | AUC |
|---|---|---|---|---|
| MatchNet [24] | $0.8742 \pm 0.0068$ | $0.9589 \pm 0.0047$ | $0.9146 \pm 0.0025$ | $0.7723 \pm 0.0119$ |
| SiameseNet [47] | $0.7210 \pm 0.0086$ | $0.8968 \pm 0.0109$ | $0.7992 \pm 0.0076$ | $0.7748 \pm 0.0088$ |
| HardNet [27] | $0.8762 \pm 0.0011$ | $0.9533 \pm 0.0094$ | $0.9131 \pm 0.0048$ | $0.7747 \pm 0.0047$ |
| SOSNet [45] | $0.8763 \pm 0.0010$ | $0.9544 \pm 0.0086$ | $0.9136 \pm 0.0044$ | $0.7752 \pm 0.0043$ |
| D2-Net [31] | $0.8032 \pm 0.0033$ | $0.9005 \pm 0.0084$ | $0.8491 \pm 0.0052$ | $0.6194 \pm 0.0078$ |
| ASLFeat [32] | $0.8729 \pm 0.0036$ | $0.9048 \pm 0.0073$ | $0.8886 \pm 0.0035$ | $0.7548 \pm 0.0062$ |
| SuperGlue [33] | $0.8639 \pm 0.0054$ | $0.8747 \pm 0.0077$ | $0.8692 \pm 0.0049$ | $0.7305 \pm 0.0099$ |
| LoFTR [50] | $0.8515 \pm 0.0020$ | $0.9837 \pm 0.0060$ | $0.9129 \pm 0.0029$ | $0.7346 \pm 0.0045$ |
| VGIDM (GAT) [ours] | $0.9052 \pm 0.0047$ | $0.9517 \pm 0.0044$ | $0.9278 \pm 0.0040$ | $0.8263 \pm 0.0092$ |
| VGIDM (GCN) [ours] | $0.8918 \pm 0.0051$ | $0.9515 \pm 0.0077$ | $0.9206 \pm 0.0037$ | $0.8025 \pm 0.0087$ |
| VGIDM (GraphSAGE) [ours] | $0.8938 \pm 0.0046$ | $0.9648 \pm 0.0052$ | $0.9279 \pm 0.0045$ | $0.8104 \pm 0.0096$ |



Fig. 8. Examples of matched and mismatched pairs from the Landmark Oxford dataset. A green or red box indicates a correct or incorrect prediction result, respectively. "GT" stands for ground truth.

pixel resolution of the image patches. As a landmark can be far away from the camera on the vehicle, its corresponding image patch can be small. As a result, it is more likely for these models to make mistakes in the matching decision.

**Cross-Validation.** To evaluate the generalization capability of VGIDM, we train VGIDM on the Landmark Oxford dataset and test it on the Landmark KITTI dataset. From Table III and Table I, we observe that the VGIDM variants still outperform the other baselines in all metrics. From Table III and Table II, when we train on the Landmark KITTI dataset and test on the Landmark Oxford dataset, VGIDM outperforms the other baselines in precision and AUC. Since the Landmark Oxford dataset contains windows as landmarks while not the Landmark KITTI dataset, the test performance on the Landmark Oxford dataset deteriorates more significantly. In Tables I and II, the baselines D2-Net, ASLFeat, SuperGlue and LoFTR do not perform training and test on the same dataset. Since there is no point-level ground-truth for our landmark patches, we adopt pre-trained models for these baselines from the literature [31]–[33], [50]. The other baselines are trained and tested on the same datasets.

### D. Ablation Study

**Feature Pair Discrimination.** We perform ablation studies on different feature pairs for the matching task shown in Table IV. The feature pair comparison include $d(f(x), f(y))$ for Resnet features, $d(\rho(x), \rho(y))$ for vertex embeddings, $d(\varphi(x), \varphi(y))$ for ensemble vertex embeddings, as well as $d(\psi(\mathcal{G}^x), \psi(\mathcal{G}^y))$ for neighborhood graph embeddings, where the learnable layer $d$ given by (3) as the metric. For each feature comparison, we train the corresponding models for $150$ epochs and select the optimal test result for the matching task based on the Landmark Oxford dataset. From Table IV, we observe that our proposed vertex-to-graph comparison outperforms the other feature pairs in most metrics. The feature pairs containing graph information, like $\psi(\mathcal{G}^x)$ and $\psi(\mathcal{G}^y)$, have an advantage over those based only on vertex information, like $f(x)$ and $f(y)$. This demonstrates the benefit of utilizing graph information.

**Discriminator Function.** We evaluate the effectiveness of the learnable discriminator $d$ by comparing it with other discriminator functions. Specifically, we replace the learnable discriminator $d$ with either cosine similarity or $\mathcal{L}_2$ distance in (6). We evaluate the patch-matching task on the Landmark

TABLE III
CROSS-VALIDATION ON THE LANDMARK KITTI DATASET OR LANDMARK OXFORD DATASET USING THE TRAINED MODEL BASED ON THE LANDMARK OXFORD DATASET OR LANDMARK KITTI DATASET, RESPECTIVELY. THE "BEST IN TABLE I" OR "BEST IN TABLE II" METHOD REFERS TO THE BEST-PERFORMING BASELINE OUT OF D2-NET, ASLFEAT, SUPERGLUE AND LOFTR FROM TABLE I OR TABLE II.

| Cross-Validation | Methods | Precision | Recall | $F_1$-Score | AUC |
|---|---|---|---|---|---|
| Oxford dataset (Train) & KITTI dataset (Test) | VGIDM (GAT) | 0.9238 ± 0.0020 | 0.9047 ± 0.0039 | 0.9141 ± 0.0027 | 0.8403 ± 0.0041 |
| | VGIDM (GCN) | 0.9084 ± 0.0019 | **0.9579** ± 0.0074 | 0.9325 ± 0.0041 | 0.8341 ± 0.0050 |
| | VGIDM (GraphSAGE) | **0.9255** ± 0.0029 | 0.9483 ± 0.0102 | **0.9368** ± 0.0061 | **0.8597** ± 0.0080 |
| Baselines tested on KITTI dataset | Best in Table I | 0.9189 ± 0.0022 | 0.9243 ± 0.0110 | 0.9154 ± 0.0059 | 0.8312 ± 0.0057 |
| KITTI dataset (Train) & Oxford dataset (Test) | VGIDM (GAT) | **0.9022** ± 0.0031 | 0.8930 ± 0.0040 | 0.8976 ± 0.0024 | **0.8013** ± 0.0051 |
| | VGIDM (GCN) | 0.8909 ± 0.0043 | 0.8845 ± 0.0124 | 0.8877 ± 0.0074 | 0.7799 ± 0.0096 |
| | VGIDM (GraphSAGE) | 0.8918 ± 0.0055 | 0.9098 ± 0.0064 | 0.9007 ± 0.0043 | 0.7893 ± 0.0099 |
| Baselines tested on Oxford dataset | Best in Table II | 0.8729 ± 0.0036 | **0.9837** ± 0.0060 | **0.9129** ± 0.0029 | 0.7548 ± 0.0062 |

TABLE IV
ABLATION STUDY FOR FEATURE PAIR DISCRIMINATION.

| Feature Comparison | Precision | Recall | $F_1$-Score | AUC |
|---|---|---|---|---|
| $d(f(x), f(y))$ | 0.8817 | 0.9146 | 0.8979 | 0.7733 |
| $d(\rho(x), \rho(y))$ | 0.7772 | **0.9720** | 0.8637 | 0.5680 |
| $d(\varphi(x), \varphi(y))$ | 0.8819 | 0.9560 | 0.9175 | 0.7860 |
| $d(\psi(\mathcal{G}^x), \psi(\mathcal{G}^y))$ | 0.9029 | 0.9427 | 0.9224 | 0.8193 |
| $d(\varphi(x), \psi(\mathcal{G}^y))$ | **0.9097** | 0.9533 | **0.9310** | **0.8347** |

Oxford dataset. From Table V, we observe that the proposed learnable discriminator $d$ outperforms the other discriminators, which is likely due to the adaptability of neural networks to different feature dimensions.

TABLE V
ABLATION STUDY FOR DIFFERENT DISCRIMINATOR FUNCTIONS.

| Feature Discriminator | Precision | Recall | $F_1$-Score | AUC |
|---|---|---|---|---|
| Cosine similarity | 0.9007 | 0.8707 | 0.8854 | 0.7913 |
| $\mathcal{L}_2$ distance | 0.7277 | 0.8800 | 0.7966 | 0.5540 |
| Learnable $d$ | **0.9097** | **0.9533** | **0.9310** | **0.8347** |

TABLE VI
MATCHING PERFORMANCE OF DIFFERENT METHODS BASED ON SPATIAL NEIGHBORHOOD INFORMATION.

| Methods | Precision | Recall | $F_1$-Score | AUC |
|---|---|---|---|---|
| SIFT [71]+RANSAC [52] (+neighbors) | 0.8965 | 0.9467 | 0.9209 | 0.8093 |
| MatchNet [24] (+neighbors) | 0.9023 | 0.9600 | 0.9302 | 0.8240 |
| SuperGlue [33] (+neighbors) | 0.9225 | 0.9200 | 0.9212 | 0.8440 |
| LoFTR [50] (+neighbors) | 0.9192 | 0.9413 | 0.9302 | 0.8467 |
| VGIDM [ours] | **0.9280** | **0.9627** | **0.9450** | **0.8693** |

**Spatial Neighborhood Information.** We investigate whether the performance improvement of VGIDM is mainly due to the spatial neighborhood information. To do this, we introduce the neighborhood graphs used in VGIDM to other baseline methods. Specifically, for a given vertex, we sort its neighbors according to increasing distances from it. We then use each baseline method to compare not only the vertex pair but also the pairs of their corresponding neighbors with the same sort order. Then, we calculate the average of the predicted scores for the vertex pair and its neighbor pairs. Finally, we decide whether there is a match based on a threshold, which is a hyperparameter tuned separately to achieve the best performance for each baseline. Table VI shows results on the Landmark KITTI dataset, where the best test performance for each baseline model is selected.

We compare with VGIDM (GraphSAGE), which is trained on the Landmark Oxford dataset. We observe that including neighborhood information generally improves the performance of every baseline, but VGIDM still outperforms them. This indicates that the neighborhood graph feature representation in VGDIM has advantages in the patch-matching task. As shown in Table VII, which presents the inference runtime for one pair of frames (with around twenty patch pairs for comparison), the computational complexity of VGIDM is lower than most baselines, except MatchNet.

TABLE VII
INFERENCE RUNTIME COMPARISON FOR DIFFERENT METHODS ON THE LANDMARK KITTI DATASET.

| Methods | SIFT+ RANSAC (+neighbors) | MatchNet (+neighbors) | SuperGlue (+neighbors) | LoFTR (+neighbors) | VGIDM (GraphSAGE) |
|---|---|---|---|---|---|
| Inference runtime | 0.7035s | 0.1591s | 3.8072s | 4.3740s | 0.2013s |

*E. Computational Complexity*

To evaluate the runtime performance, we test VGIDM on an NVIDIA RTX A5000 GPU. Table VIII shows the inference runtime (mean time for one pair of frames in the testing phase) for the VGIDM variants with different GNNs. Specifically, given a pair of frames (i.e., full-size images), an average of around twenty patch pairs are compared, which takes less than 0.25 seconds. The time taken is acceptable for practical applications, such as place recognition and autonomous driving. Moreover, the amount of the parameters for these VGIDM networks with the GAT, GCN and GraphSAGE is 12.16M, 12.16M and 12.66M, respectively.

TABLE VIII
INFERENCE RUNTIME OF VGIDM ON LANDMARK OXFORD DATASET.

| Methods | VGIDM (GAT) | VGIDM (GCN) | VGIDM (GraphSAGE) |
|---|---|---|---|
| Inference Runtime | 0.2330s | 0.1953s | 0.2092s |

*F. Further Possible Applications*

*1) Application of VGIDM in Visual Place Recognition:* A possible application of VGIDM is visual place recognition. We apply our local patch-matching to obtain global frame matching to determine if two frames show the same place. In visual place recognition, we construct a bipartite graph with edges
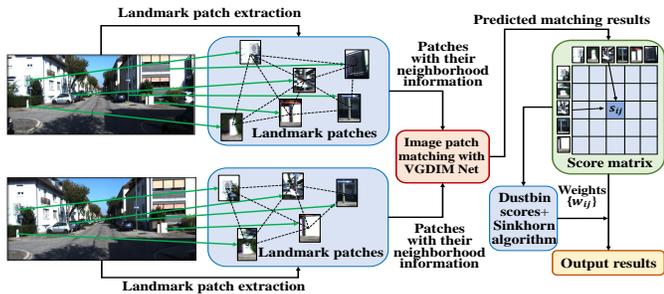
Fig. 9. The diagram of visual place recognition with VGIDM.

being the scores output by our network for all landmark patch pairs, which is used to construct a matching score matrix. Then, similar to the Optimal Matching Layer in [33], by appending learnable dustbin scores for the score matrix, the Sinkhorn algorithm is used to output the partial assignment. Finally, we obtain frame-matching results by summing up the elements in the matching score matrix with the weights from the partial assignment. The details are shown in Fig. 9, where GAT is chosen for the GNN part in VGIDM.

TABLE IX
VISUAL PLACE RECOGNITION PERFORMANCE IN KITTI DATA.

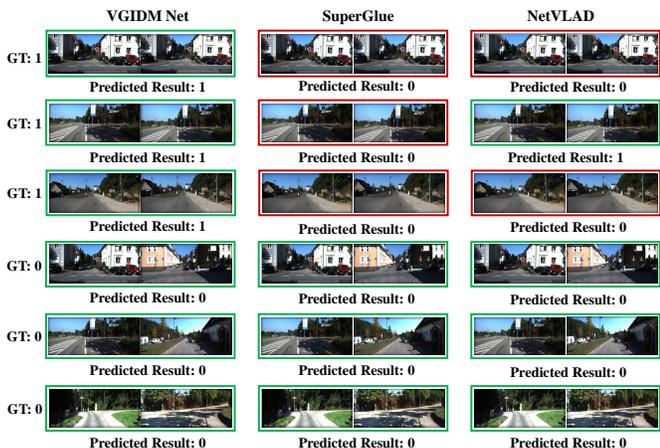| Methods | MatchNet | NetVLAD | SuperGlue | LoFTR | VGIDM |
|---|---|---|---|---|---|
| $F_1$-Score | 0.9668 | 0.9702 | 0.9694 | 0.9711 | **0.9719** |
| Accuracy | 0.9360 | 0.9424 | 0.9408 | 0.9440 | **0.9452** |



Fig. 10. Several examples of place recognition on the KITTI dataset. The prediction "1" indicates the frames are from the same place, while "0" indicates they are from different places. A green box indicates a correct prediction result while a red box an incorrect one.

TABLE X
CROSS-VALIDATION PERFORMANCE FOR VISUAL PLACE RECOGNITION ON OXFORD DATASET.

| Methods | MatchNet | NetVLAD | SuperGlue | LoFTR | VGIDM |
|---|---|---|---|---|---|
| $F_1$-Score | 0.9010 | 0.9069 | 0.9190 | 0.9207 | **0.9266** |
| Accuracy | 0.8273 | 0.8303 | 0.8563 | 0.8607 | **0.8680** |

We compare VGIDM with MatchNet [24], NetVLAD [72], SuperGlue [33], and LoFTR [50] under the place recognition task with around 600 pairs of place images from the KITTI dataset. The two places contained in each image frame pair are regarded as the same if the distance between the camera
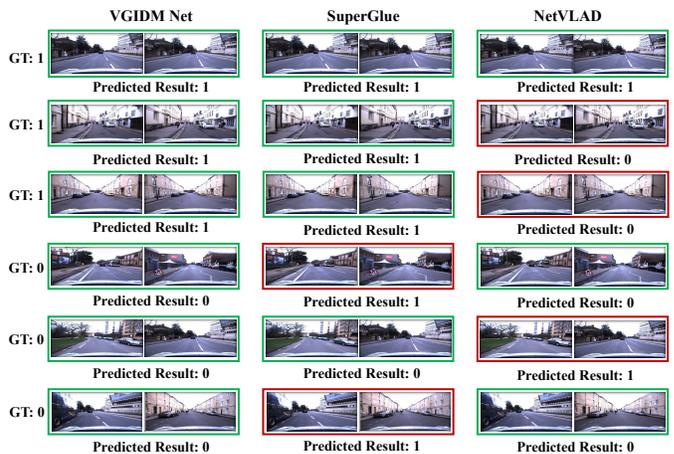


Fig. 11. Several examples of place recognition on the Oxford dataset where the model is trained on the KITTI dataset. The prediction "1" or "0" indicates the frames are from the same or different places. A green or red box indicates a correct or incorrect prediction result.

locations is less than 10 meters. To recognize the same places more accurately, the thresholds of the matching results for place image pairs are set to higher recall levels. The results and example outputs are shown in Table IX and Fig. 10, respectively. From Table IX, it is observed that all the methods perform well, with VGIDM having a slight advantage. The reason may be that there exist obvious differences among the image pairs that are not from the same place. However, unlike MatchNet, NetVLAD, SuperGlue, and LoFTR, which require the full image, VGIDM can perform place recognition by using only landmark patches and their spatial relationships. Moreover, the landmark patches based on static objects are more stable than the keypoints based on edges or corners and are not affected by noisy image pixels from non-persistent objects or surroundings.

We conduct cross-validation experiments for the visual place recognition task. Specifically, we use VGIDM (GAT) trained on the KITTI dataset, and compare it with baselines for inference on the Oxford dataset. We use 3000 frames in the experiments. From Table X, we observe that VGIDM is superior to the baselines. This suggests that VGIDM has good generalization ability. We include a few examples in Fig. 11.

*2) Application of VGIDM in Stereo Depth Estimation of Landmarks:* Another application is *stereo* depth estimation of landmarks. Similar to the steps described in Section IV-A, we obtain the landmarks from the full-sized frames captured from both the *left* and *right* stereo cameras.

Different from the experiment settings in Section IV-A where the matching is performed for landmark patches in image frames captured at different locations, here we use VGIDM to perform the matching between landmark patches captured at the *same* location but from *different* cameras. We split 3000 pairs of stereo frames into training and testing sets with a ratio of around 2 : 1. During testing, we set a high similarity threshold of 0.9 to prevent false positive matching. Since the landmark objects we have chosen have regular shapes, the original narrow landmark object detection boxes (without the intentionally added background to form the landmark patches) are sufficiently accurate to locate the landmarks in the frames.

For each of the matched landmark objects in the left and

TABLE XI
ACCURACY OF DEPTH ESTIMATION.

| Method | Monocular Depth [38] | DIFFNet [62] | VGIDM [ours] |
|---|---|---|---|
| **RMSE** (m) | 14.22 | 4.45 | **0.86** |

right frames, we compute the pixel disparity on the center line (average of left and right sides) of the original bounding box. The depth of the pixels on the landmark bounding box middle lines can be calculated using the camera focal length and distance. Following the diagram in Fig. 12, the coarse monocular depth is improved to a more accurate stereo depth as shown in Table XI. The vanilla Monocular Depth [38] only achieves $14.22$m RMSE. After applying VGIDM, we can improve the depth estimation accuracy to $0.86$m RMSE. In contrast, the current state-of-the-art DIFFNET [62] only achieves $4.45$m RMSE performance in stereo depth estimation. However, a direct application of VGIDM can only output the estimated depth for sparse pixels (only for the chosen landmarks). To improve general stereo depth estimation, it is possible to incorporate VGIDM into existing methods, e.g., using VGIDM's accurately estimated stereo landmark depths as calibration for other general stereo depth estimation algorithms like DIFFNet [62], HRDepth [63], and CADepth [64].
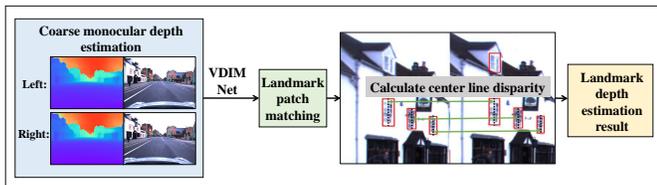


Fig. 12. Depth estimation: from coarse monocular depth to fine stereo depth.

VGIDM can serve as a module in various learning-based localization techniques, such as Detect-SLAM [73], EAO-SLAM [74], and other semantic SLAM with object-level data association [75], [76]. These techniques underscore the importance of landmark patch matching, which is central to VGIDM in real-world applications. Consequently, our approach has crucial implications for diverse applications, demonstrating its versatility and effectiveness.

## V. CONCLUSION

We have developed an image patch-matching model VGIDM that incorporates spatial information of the landmark patches through graph-based learning. We provided a theoretical basis for our approach. Extensive experiments demonstrate that our approach outperforms the current state-of-the-art baselines, which do not take into account the spatial relationships between patches. Our framework indicates that such spatial information can be beneficial to landmark patch-matching in street scenes.

In future work, it is of interest to incorporate a greater variety of objects as landmarks to adapt VGIDM to more diverse street scenes and generalize to more datasets. To achieve this, we can train VGIDM on landmarks from a wider range of classes. As our method is better suited for matching tasks in scenarios with sufficient landmarks or static objects, pixel-level matching methods can serve a complementary role in scenarios

with fewer landmarks. Additionally, combining our patch-level matching method with point-level methods shows promise in achieving more accurate pixel-level or sub-pixel-level matching. To this end, we can use our method as a post-processing step to emphasize the keypoints with more attention or to filter out weak correspondences. We can further generalize VGIDM to multi-view camera relocalization [12] to estimate camera poses by determining the matched landmarks in multiple image frames. Furthermore, the matched landmarks can serve as anchor points or interest regions to aid other applications. In LiDAR relocalization [6], [13] or LiDAR odometry estimation, by restricting the LiDAR points from matched landmarks using VGIDM, many outlier points can be removed, leading to better estimation accuracy.

## APPENDIX A
### PROOF OF PROPOSITION 1

From (12), we have

$$L_{\text{ID}}$$
$$= \mathbb{P}(x \leftrightarrow y)\mathbb{E}[\log d(\varphi(x), \psi(\mathcal{G}^y)) \,|\, x \leftrightarrow y]$$
$$\quad + \mathbb{P}(x \not\leftrightarrow y)\mathbb{E}[\log(1 - d(\varphi(x), \psi(\mathcal{G}^y))) \,|\, x \not\leftrightarrow y] \quad (17)$$
$$= \int_{\mathbb{A}} \Big\{ \mathbb{P}(x \leftrightarrow y)p(\varphi(x), \psi(\mathcal{G}^y) \,|\, x \leftrightarrow y) \log d(\varphi(x), \psi(\mathcal{G}^y))$$
$$\quad + \mathbb{P}(x \not\leftrightarrow y)p(\varphi(x), \psi(\mathcal{G}^y) \,|\, x \not\leftrightarrow y)$$
$$\quad \cdot \log(1 - d(\varphi(x), \psi(\mathcal{G}^y))) \Big\} \mathrm{d}(\varphi(x), \psi(\mathcal{G}^y)), \quad (18)$$

where $\mathbb{P}(x \leftrightarrow y)$ and $\mathbb{P}(x \not\leftrightarrow y)$ are the matched and unmatched probabilities. From (18), it is clear that $L_{\text{ID}}$ is concave in $d(\varphi(x), \psi(\mathcal{G}^y))$ for every $(x, y)$. Taking the first derivatives and setting them to zero, we obtain

$$d^*(\varphi(x), \psi(\mathcal{G}^y)) = \frac{\mathbb{P}(x \leftrightarrow y)p(\varphi(x), \psi(\mathcal{G}^y) \,|\, x \leftrightarrow y)}{p(\varphi(x), \psi(\mathcal{G}^y))}. \quad (19)$$

Substituting (19) into (18), we have

$$L_{\text{ID}}^{d^*}$$
$$= \mathbb{P}(x \leftrightarrow y)\mathbb{E}\left[\log \frac{p(\varphi(x), \psi(\mathcal{G}^y) \,|\, x \leftrightarrow y)}{p(\varphi(x), \psi(\mathcal{G}^y) \,|\, x \not\leftrightarrow y)} \,\middle|\, x \leftrightarrow y \right]$$
$$\quad + \mathbb{E}\left[\log \frac{p(\varphi(x), \psi(\mathcal{G}^y) \,|\, x \not\leftrightarrow y)}{p(\varphi(x), \psi(\mathcal{G}^y))}\right] - H_b(\mathbb{P}(x \leftrightarrow y)). \quad (20)$$

Applying Jensen's inequality to the second term in the right-hand side of (20), we have

$$\mathbb{E}\left[\log \frac{p(\varphi(x), \psi(\mathcal{G}^y) \,|\, x \not\leftrightarrow y)}{p(\varphi(x), \psi(\mathcal{G}^y))}\right]$$
$$\leq \log \mathbb{E}\left[\frac{p(\varphi(x), \psi(\mathcal{G}^y) \,|\, x \not\leftrightarrow y)}{p(\varphi(x), \psi(\mathcal{G}^y))}\right]$$
$$= \log \int_{\mathbb{A}} p(\varphi(x), \psi(\mathcal{G}^y) \,|\, x \not\leftrightarrow y)\mathrm{d}(\varphi(x), \psi(\mathcal{G}^y))$$
$$= 0. \quad (21)$$

Therefore, from (20), we have

$$L_{\text{ID}}^{d^*} \leq \mathbb{P}(x \leftrightarrow y)$$
$$\quad \cdot D(p(\varphi(x), \psi(\mathcal{G}^y) \,|\, x \leftrightarrow y) \,\|\, p(\varphi(x), \psi(\mathcal{G}^y) \,|\, x \not\leftrightarrow y))$$
$$\quad - H_b(\mathbb{P}(x \leftrightarrow y)). \quad (22)$$

$$L_{\text{ID}}(\varphi, \psi, \tilde{d}^*) - L_{\text{ID}}(\varphi, \psi, d^*)$$
$$= -\frac{\varepsilon^2}{2} \int_{\mathbb{A}} \frac{(\mathbb{P}(x \leftrightarrow y)p(\varphi(x), \psi(\mathcal{G}^y) \mid x \leftrightarrow y) + \mathbb{P}(x \nleftrightarrow y)p(\varphi(x), \psi(\mathcal{G}^y) \mid x \nleftrightarrow y))^3}{\mathbb{P}(x \leftrightarrow y)\mathbb{P}(x \nleftrightarrow y)p(\varphi(x), \psi(\mathcal{G}^y) \mid x \leftrightarrow y)p(\varphi(x), \psi(\mathcal{G}^y) \mid x \nleftrightarrow y)} \mathrm{d}(\varphi(x), \psi(\mathcal{G}^y)) + o(\varepsilon^2) \quad (25)$$

Rearranging the inequality completes the proof.

## APPENDIX B
## PROOF OF PROPOSITION 2

Let $\tilde{d} = d + \epsilon$. Similar to (17) in the proof of Proposition 1, it is easy to see

$$L_{\text{ID}}(\varphi, \psi, \tilde{d})$$
$$= \mathbb{P}(x \leftrightarrow y)\mathbb{E}[\log(d(\varphi(x), \psi(\mathcal{G}^y)) + \varepsilon) \mid x \leftrightarrow y]$$
$$+ \mathbb{P}(x \nleftrightarrow y)\mathbb{E}[\log(1 - d(\varphi(x), \psi(\mathcal{G}^y)) - \varepsilon) \mid x \nleftrightarrow y], \quad (23)$$

where the notations are the same as those in (17).

According to Taylor's series expansion theorem [77], we have

$$L_{\text{ID}}(\varphi, \psi, \tilde{d}) - L_{\text{ID}}(\varphi, \psi, d)$$
$$= \varepsilon \left\{ \mathbb{P}(x \leftrightarrow y)\mathbb{E}\left[ \frac{1}{d(\varphi(x), \psi(\mathcal{G}^y))} \,\middle|\, x \leftrightarrow y \right] \right.$$
$$\left. - \mathbb{P}(x \nleftrightarrow y)\mathbb{E}\left[ \frac{1}{1 - d(\varphi(x), \psi(\mathcal{G}^y))} \,\middle|\, x \nleftrightarrow y \right] \right\}$$
$$- \frac{\varepsilon^2}{2} \left\{ \mathbb{P}(x \leftrightarrow y)\mathbb{E}\left[ \frac{1}{(d(\varphi(x), \psi(\mathcal{G}^y)))^2} \,\middle|\, x \leftrightarrow y \right] \right.$$
$$\left. + \mathbb{P}(x \nleftrightarrow y)\mathbb{E}\left[ \frac{1}{(1 - d(\varphi(x), \psi(\mathcal{G}^y)))^2} \,\middle|\, x \nleftrightarrow y \right] \right\}$$
$$+ o(\varepsilon^2). \quad (24)$$

Furthermore, by substituting $d(\varphi(x), \psi(\mathcal{G}^y)) = d^*(\varphi(x), \psi(\mathcal{G}^y))$ given in (19) into (24), we have (25) and the proof is complete.

## APPENDIX C
## PROOF OF PROPOSITION 3

We have

$$\|p(\varphi(x), \psi(\mathcal{G}^y) \mid x \leftrightarrow y) - p(\varphi(x), \psi(\mathcal{G}^y) \mid x \nleftrightarrow y)\|_{\text{TV}}$$
$$= \sum_{(\varphi(x), \psi(\mathcal{G}^y)) \in \mathbb{B}} \left\{ m(x \leftrightarrow y)p(\varphi(x), \psi(\mathcal{G}^y) \mid \mathcal{G}^x \leftrightarrow \mathcal{G}^y, x \leftrightarrow y) \right.$$
$$+ (1 - m(x \leftrightarrow y))p(\varphi(x), \psi(\mathcal{G}^y) \mid \mathcal{G}^x \nleftrightarrow \mathcal{G}^y, x \leftrightarrow y)$$
$$- m(x \nleftrightarrow y)p(\varphi(x), \psi(\mathcal{G}^y) \mid \mathcal{G}^x \leftrightarrow \mathcal{G}^y, x \nleftrightarrow y)$$
$$\left. - (1 - m(x \nleftrightarrow y))p(\varphi(x), \psi(\mathcal{G}^y) \mid \mathcal{G}^x \nleftrightarrow \mathcal{G}^y, x \nleftrightarrow y) \right\}$$
$$\geq \sum_{(\varphi(x), \psi(\mathcal{G}^y)) \in \mathbb{B}} \left\{ \min_{x \leftrightarrow y} m(x \leftrightarrow y) \right.$$
$$\cdot p(\varphi(x), \psi(\mathcal{G}^y) \mid \mathcal{G}^x \leftrightarrow \mathcal{G}^y, x \leftrightarrow y)$$
$$- \max_{x \nleftrightarrow y} m(x \nleftrightarrow y)p(\varphi(x), \psi(\mathcal{G}^y) \mid \mathcal{G}^x \leftrightarrow \mathcal{G}^y, x \nleftrightarrow y)$$

$$\left. - p(\varphi(x), \psi(\mathcal{G}^y) \mid \mathcal{G}^x \nleftrightarrow \mathcal{G}^y, x \nleftrightarrow y) \right\} \quad (26)$$
$$\geq \min_{x \leftrightarrow y} m(x \leftrightarrow y)$$
$$\cdot \sum_{(\varphi(x), \psi(\mathcal{G}^y)) \in \mathbb{B}} \left\{ p(\varphi(x), \psi(\mathcal{G}^y) \mid \mathcal{G}^x \leftrightarrow \mathcal{G}^y, x \leftrightarrow y) \right.$$
$$\left. - p(\varphi(x), \psi(\mathcal{G}^y) \mid \mathcal{G}^x \nleftrightarrow \mathcal{G}^y, x \leftrightarrow y) \right\} + \min_{x \leftrightarrow y} m(x \leftrightarrow y)$$
$$- \max_{x \nleftrightarrow y} m(x \nleftrightarrow y) - 1, \quad (27)$$

where the inequality (26) follows from $0 \leq m(x \leftrightarrow y) \leq 1$ and $0 \leq m(x \nleftrightarrow y) \leq 1$. The proof is now complete.

## REFERENCES

[1] X. Zhang, S. Wang, Z. Li, and S. Ma, "Landmark image retrieval by jointing feature refinement and multimodal classifier learning," *IEEE Trans. Cybern.*, vol. 48, no. 6, pp. 1682–1695, Jun. 2017.

[2] J. Zhu, H. Zeng, J. Huang, S. Liao, Z. Lei, C. Cai, and L. Zheng, "Vehicle re-identification using quadruple directional deep learning features," *IEEE Trans. Intell. Transp. Syst.*, vol. 21, no. 1, pp. 410–420, Mar. 2019.

[3] S. Wang, X. Guo, Y. Tie, L. Qi, and L. Guan, "Deep local feature descriptor learning with dual hard batch construction," *IEEE Trans. Image Process.*, vol. 29, pp. 9572–9583, Oct. 2020.

[4] D. Quan, S. Wang, Y. Li, B. Yang, H. Ning, J. Chanussot, H. Biao, and L. Jiao, "Multi-relation attention network for image patch matching," *IEEE Trans. Image Process.*, vol. 30, pp. 7127–7142, Aug. 2021.

[5] S. Liao and A. C. Chung, "Nonrigid brain MR image registration using uniform spherical region descriptor," *IEEE Trans. Image Process.*, vol. 21, no. 1, pp. 157–169, Jun. 2011.

[6] N. Engel and etc., "Deeplocalization: Landmark-based self-localization with deep neural networks," in *Proc. IEEE Intell. Transp. Syst. Conf.*, 2019, pp. 926–933.

[7] Y. Wang, Y. Qiu, P. Cheng, and X. Duan, "Robust loop closure detection integrating visual-spatial-semantic information via topological graphs and CNN features," *J. Remote Sensing*, vol. 12, no. 23, p. 3890, Oct. 2020.

[8] Z. Zhu, T. Oskiper, S. Samarasekera, R. Kumar, and H. S. Sawhney, "Ten-fold improvement in visual odometry using landmark matching," in *Proc. IEEE Int. Conf. Comput. Vision*, 2007, pp. 1–8.

[9] J. Vincent, M. Labbé, J. S. Lauzon, F. Grondin, P. M. Comtois-Rivet, and F. Michaud, "Dynamic object tracking and masking for visual SLAM," in *Proc. IEEE Int. Conf. Intell. Robots Syst.*, 2020, pp. 4974–4979.

[10] N. Sünderhauf, S. Shirazi, A. Jacobson, F. Dayoub, E. Pepperell, B. Upcroft, and M. Milford, "Place recognition with convnet landmarks: Viewpoint-robust, condition-robust, training-free," in *Proc. Robot.: Sci. Syst.*, 2017, pp. 5702–5708.

[11] S. Hausler, S. Garg, M. Xu, M. Milford, and T. Fischer, "Patch-NetVLAD: Multi-scale fusion of locally-global descriptors for place recognition," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2021, pp. 14 141–14 152.

[12] F. Xue, X. Wu, S. Cai, and J. Wang, "Learning multi-view camera relocalization with graph neural networks," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2020, pp. 11 372–11 381.

[13] J. Zhang and S. Singh, "Visual-lidar odometry and mapping: Low-drift, robust, and fast," in *Proc. IEEE Int. Conf. Robot. Automat.*, 2015, pp. 2174–2181.

[14] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *Int. J. Comput. Vision*, vol. 60, no. 2, pp. 91–110, Jan. 2004.

[15] H. Bay, A. Ess, T. Tuytelaars, and L. Van Gool, "Speeded-up robust features (SURF)," *J. Comput. Vision Image Understanding*, vol. 110, no. 3, pp. 346–359, Jun. 2008.

[16] P. Felzenszwalb, D. McAllester, and D. Ramanan, "A discriminatively trained, multiscale, deformable part model," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2008, pp. 1–8.

[17] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, "ORB: An efficient alternative to SIFT or SURF," in *Proc. IEEE Int. Conf. Comput. Vision*, 2011, pp. 2564–2571.

[18] S. Leutenegger, M. Chli, and R. Y. Siegwart, "BRISK: Binary robust invariant scalable keypoints," in *Proc. IEEE Int. Conf. Comput. Vision*, 2011, pp. 2548–2555.

[19] A. Alahi, R. Ortiz, and P. Vandergheynst, "Freak: Fast retina keypoint," in *Proc. IEEE Int. Conf. Comput. Vision*, 2012, pp. 510–517.

[20] Y. Tian, B. Fan, and F. Wu, "L2-net: Deep learning of discriminative patch descriptor in euclidean space," in *Proc. IEEE Int. Conf. Comput. Vision*, 2017, pp. 661–669.

[21] S. Zagoruyko and N. Komodakis, "Learning to compare image patches via convolutional neural networks," in *Proc. IEEE Int. Conf. Comput. Vision*, 2015, pp. 4353–4361.

[22] V. Kumar BG, G. Carneiro, and I. Reid, "Learning local image descriptors with deep siamese and triplet convolutional networks by minimising global loss functions," in *Proc. IEEE Int. Conf. Comput. Vision*, 2016, pp. 5385–5394.

[23] V. Balntas, E. Riba, D. Ponsa, and K. Mikolajczyk, "Learning local feature descriptors with triplets and shallow convolutional neural networks," in *Proc. British Mach. Vision Conf.*, 2016, pp. 1–11.

[24] X. Han, T. Leung, Y. Jia, R. Sukthankar, and A. C. Berg, "MatchNet: Unifying feature and metric learning for patch-based matching," in *Proc. IEEE Int. Conf. Comput. Vision*, 2015, pp. 3279–3286.

[25] A. Subramaniam, P. Balasubramanian, and A. Mittal, "NCC-net: Normalized cross correlation based deep matcher with robustness to illumination variations," in *Proc. IEEE Winter Conf. Appl. Comput. Vision*, 2018, pp. 1944–1953.

[26] D. Quan, X. Liang, S. Wang, S. Wei, Y. Li, H. Ning, and L. Jiao, "AFD-Net: Aggregated feature difference learning for cross-spectral image patch matching," in *Proc. IEEE Int. Conf. Comput. Vision*, 2019, pp. 3017–3026.

[27] A. Mishchuk, D. Mishkin, F. Radenovic, and J. Matas, "Working hard to know your neighbor's margins: Local descriptor learning loss," in *Proc. Conf. Neural Inform. Process. Syst.*, 2017, pp. 1–10.

[28] S. Wang, Y. Li, X. Liang, D. Quan, B. Yang, S. Wei, and L. Jiao, "Better and faster: Exponential loss for image patch matching," in *Proc. IEEE Int. Conf. Comput. Vision*, 2019, pp. 4812–4821.

[29] T. Ng, V. Balntas, Y. Tian, and K. Mikolajczyk, "SOLAR: Second-order loss and attention for image retrieval," in *Proc. Eur. Conf. Comput. Vision*, 2020, pp. 253–270.

[30] Y. Miao, Z. Lin, X. Ma, G. Ding, and J. Han, "Learning transformation-invariant local descriptors with low-coupling binary codes," *IEEE Trans. Image Process.*, vol. 30, pp. 7554–7566, Aug. 2021.

[31] M. Dusmanu, I. Rocco, T. Pajdla, M. Pollefeys, J. Sivic, A. Torii, and T. Sattler, "D2-net: A trainable CNN for joint description and detection of local features," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2019, pp. 8092–8101.

[32] Z. Luo, L. Zhou, X. Bai, H. Chen, J. Zhang, Y. Yao, S. Li, T. Fang, and L. Quan, "ASLFeat: Learning local features of accurate shape and localization," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2020, pp. 6589–6598.

[33] P. E. Sarlin, D. DeTone, T. Malisiewicz, and A. Rabinovich, "SuperGlue: Learning feature matching with graph neural networks," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2020, pp. 4938–4947.

[34] M. Amiri and H. R. Rabiee, "RASIM: A novel rotation and scale invariant matching of local image interest points," *IEEE Trans. Image Process.*, vol. 20, no. 12, pp. 3580–3591, May 2011.

[35] P. Veličković, W. Fedus, W. L. Hamilton, P. Liò, Y. Bengio, and R. D. Hjelm, "Deep graph infomax," in *Proc. Int. Conf. Learn. Representations*, 2018, pp. 1–13.

[36] F. Y. Sun, J. Hoffman, V. Verma, and J. Tang, "InfoGraph: Unsupervised and semi-supervised graph-level representation learning via mutual information maximization," in *Proc. Int. Conf. Learn. Representations*, 2020, pp. 1–13.

[37] J. L. Schonberger and J. M. Frahm, "Structure-from-motion revisited," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2016, pp. 4104–4113.

[38] W. Yin, J. Zhang, O. Wang, S. Niklaus, L. Mai, S. Chen, and C. Shen, "Learning to recover 3D scene shape from a single image," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2021, pp. 204–213.

[39] H. Farid and E. P. Simoncelli, "A differential optical range camera," in *Proc. Annu. Meeting Optical Soc. Amer.*, 1996, pp. 1–10.

[40] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph attention networks," in *Proc. Int. Conf. Learn. Representations*, 2018, pp. 1–12.

[41] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and S. Y. Philip, "A comprehensive survey on graph neural networks," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 32, no. 1, pp. 4–24, Mar. 2020.

[42] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? The KITTI vision benchmark suite," in *Proc. IEEE Int. Conf. Comput. Vision*, 2012, pp. 3354–3361.

[43] D. Barnes, M. Gadd, P. Murcutt, P. Newman, and I. Posner, "The Oxford Radar RobotCar Dataset: A radar extension to the Oxford RobotCar Dataset," in *Proc. IEEE Int. Conf. Robot. Automat.*, 2020, pp. 6433–6438.

[44] X. Zhang, F. X. Yu, S. Kumar, and S. F. Chang, "Learning spread-out local feature descriptors," in *Proc. IEEE Int. Conf. Comput. Vision*, 2017, pp. 4595–4603.

[45] Y. Tian, X. Yu, B. Fan, F. Wu, H. Heijnen, and V. Balntas, "SOSNet: Second order similarity regularization for local descriptor learning," in *Proc. IEEE Int. Conf. Comput. Vision*, 2019, pp. 11 016–11 025.

[46] H. Pan, Y. Chen, Z. He, F. Meng, and N. Fan, "TCDesc: Learning topology consistent descriptors for image matching," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 521, pp. 436–444, Aug. 2021.

[47] I. Melekhov, J. Kannala, and E. Rahtu, "Siamese network features for image matching," in *Proc. IEEE Int. Conf. Pattern Recognit.*, 2016, pp. 378–383.

[48] K. M. Yi, E. Trulls, V. Lepetit, and P. Fua, "LIFT: Learned invariant feature transform," in *Proc. Eur. Conf. Comput. Vision*, 2016, pp. 467–483.

[49] D. DeTone, T. Malisiewicz, and A. Rabinovich, "Superpoint: Self-supervised interest point detection and description," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit. Workshops*, 2018, pp. 224–236.

[50] J. Sun, Z. Shen, Y. Wang, H. Bao, and X. Zhou, "LoFTR: Detector-free local feature matching with transformers," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2021, pp. 8922–8931.

[51] J. Brogan, A. Bharati, D. Moreira, A. Rocha, K. W. Bowyer, P. J. Flynn, and W. J. Scheirer, "Fast local spatial verification for feature-agnostic large-scale image retrieval," *IEEE Trans. Image Process.*, vol. 30, pp. 6892–6905, 2021.

[52] M. A. Fischler and R. C. Bolles, "Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography," *Commun. ACM*, vol. 24, no. 6, pp. 381–395, 1981.

[53] J. Philbin, O. Chum, M. Isard, J. Sivic, and A. Zisserman, "Object retrieval with large vocabularies and fast spatial matching," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2007, pp. 1–8.

[54] Y. Avrithis and G. Tolias, "Hough pyramid matching: Speeded-up geometry re-ranking for large scale image retrieval," *Int. J. Comput. Vision*, vol. 107, pp. 1–19, 2014.

[55] X. Li, M. Larson, and A. Hanjalic, "Pairwise geometric matching for large-scale object retrieval," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2015, pp. 5153–5161.

[56] Y. Wang, R. Zhao, L. Liang, X. Zheng, Y. Cen, and S. Kan, "Block-based image matching for image retrieval," *J. Vis. Commun. Image Representation*, vol. 74, p. 102998, 2021.

[57] B. Jiang, P. Sun, and B. Luo, "GLMNet: Graph learning-matching convolutional networks for feature matching," *Pattern Recognit.*, vol. 121, p. 108167, 2022.

[58] H. Liu, T. Wang, Y. Li, C. Lang, Y. Jin, and H. Ling, "Joint graph learning and matching for semantic feature correspondence," *Pattern Recognit.*, vol. 134, p. 109059, 2023.

[59] S. Winder, G. Hua, and M. Brown, "Picking the best daisy," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2009, pp. 178–185.

[60] H. Aanæs, A. L. Dahl, and K. Steenstrup Pedersen, "Interesting interest points," *Int. J. Comput. Vision*, vol. 97, no. 1, pp. 18–35, Jun. 2012.

[61] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 6, pp. 1137–1149, Jun. 2016.

[62] H. Zhou, D. Greenwood, and S. Taylor, "Self-supervised monocular depth estimation with internal feature fusion," *arXiv preprint arXiv:2110.09482*, 2021.

[63] X. Lyu, L. Liu, M. Wang, X. Kong, L. Liu, Y. Liu, X. Chen, and Y. Yuan, "HR-Depth: High resolution self-supervised monocular depth estimation," *arXiv preprint arXiv:2012.07356*, 2020.

[64] J. Yan, H. Zhao, P. Bu, and Y. Jin, "Channel-wise attention-based network for self-supervised monocular depth estimation," in *Proc. Int. Conf. 3D Vision*, 2021, pp. 464–473.

[65] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask R-CNN," in *Proc. IEEE Int. Conf. Comput. Vision*, 2017, pp. 2961–2969.

[66] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Int. Conf. Comput. Vision*, 2016, pp. 770–778.

[67] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *Proc. Int. Conf. Learn. Representations*, 2017, pp. 1–14.

[68] W. L. Hamilton, R. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in *Proc. Conf. Neural Inform. Process. Syst.*, 2017, pp. 1025–1035.

[69] M. A. Uy and G. H. Lee, "PointNetVLAD: Deep point cloud based retrieval for large-scale place recognition," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2018, pp. 4470–4479.

[70] H. Alhaija, S. Mustikovela, L. Mescheder, A. Geiger, and C. Rother, "Augmented reality meets computer vision: Efficient data generation for urban driving scenes," *Int. J. Compu. Vision*, vol. 126, no. 9, pp. 961–972, Mar. 2018.

[71] D. G. Lowe, "Object recognition from local scale-invariant features." *Int. J. Comput. Vision*, vol. 60, no. 2, pp. 91–110, 2004.

[72] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, "NetVLAD: CNN architecture for weakly supervised place recognition," in *Proc. IEEE Int. Conf. Comput. Vision Pattern Recognit.*, 2016, pp. 5297–5307.

[73] F. Zhong, S. Wang, Z. Zhang, and Y. Wang, "Detect-SLAM: Making object detection and slam mutually beneficial," in *Proc. IEEE Winter Conf. Appl. Comput. Vision*, 2018, pp. 1001–1010.

[74] Y. Wu, Y. Zhang, D. Zhu, Y. Feng, S. Coleman, and D. Kerr, "EAO-SLAM: Monocular semi-dense object slam based on ensemble data association," in *Proc. IEEE Int. Conf. Intell. Robots Syst.*, 2020, pp. 4966–4973.

[75] Z. Qian, K. Patath, J. Fu, and J. Xiao, "Semantic SLAM with autonomous object-level data association," in *Proc. IEEE Int. Conf. Robot. Automat.*, 2021, pp. 11 203–11 209.

[76] X. Lin, J. Ruan, Y. Yang, L. He, Y. Guan, and H. Zhang, "Robust data association against detection deficiency for semantic slam," *IEEE Trans. Autom. Sci. Eng.*, 2023.

[77] R. H. Crowell and W. E. Slesnick, *Calculus with Analytic Geometry*. WW Norton, 1968.

## A. Landmark Datasets for Image Patch Matching

In this section, we present two landmark patch matching datasets,[5] named the *Landmark KITTI Dataset* and the *Landmark Oxford Dataset*, derived from the street scene KITTI dataset and the Oxford Radar RobotCar Dataset respectively.

We first briefly introduce the two original public datasets both of which contain image frames and LiDAR scans captured from onboard cameras and Velodyne LiDAR sensors. The KITTI dataset is a public dataset[6] with multi-sensor data for autonomous driving. It contains street scene image frames and their corresponding LiDAR point clouds, which are captured in Karlsruhe, Germany, using the Point Grey Flea 2 (FL2-14S3C-C) Camera and Velodyne HDL-64E Laserscanner, respectively. The frame resolution is $1241 \times 376$ pixels. The Oxford Radar RobotCar dataset[7] contains image frames and LiDAR scans captured on the streets in Oxford, UK, by the Point Grey Grasshopper2 (GS2-FW-14S5C-C) Camera and Velodyne HDL-32E Laserscanner, respectively. The resolution of each frame in this dataset is $1280 \times 960$ pixels.

We extract the landmark object patches from the full-sized image frames of the two original street scene datasets using an object detection neural network. In the literature on landmark-based applications, Edge Boxes are used to detect a bounding box around a patch that contains a large number of internal contours compared to the number of contours exiting from the box, which indicates the presence of an object in the enclosed patch. DeepLabV3+ is used to extract significant landmark regions. However, all of the aforementioned patch extraction or landmark detection approaches are not stable when removing dynamic objects and many noisy regions are presented. By contrast, in our datasets, we use Faster R-CNN as the stable landmark object detector to locate the region of interest for static roadside objects including traffic lights, traffic signs, poles, and facade windows. To facilitate the detection efficacy, we manually labeled those objects using the frames from the street scene KITTI dataset and the Oxford Radar RobotCar dataset. In Faster R-CNN, we choose Resnet50 with Feature Pyramid Network (FPN) as the backbone, which is already pretrained on the Imagenet dataset. During training, we use Adam optimizer with learning rate 0.0002 and weight decay 0.0001 to train the detector for 50 epochs. The training batch size is set as 2 and random horizontal flipping is used for data augmentation.

We next introduce our landmark patch matching datasets. For both the Landmark KITTI dataset and the Landmark Oxford dataset, the full-sized image frames are captured by stereo cameras, and we *only use the left frames* to extract landmark patches. The details like the landmark object bounding box labels and the patch matching ground truth are described separately for each dataset as follows.

**Landmark KITTI Dataset.** The segmentation labels are semantic segmentation masks. To perform landmark object detection, we need to first convert the semantic segmentation

[5]https://github.com/AI-IT-AVs/Landmark_patch_datasets

[6]http://www.cvlibs.net/datasets/kitti/

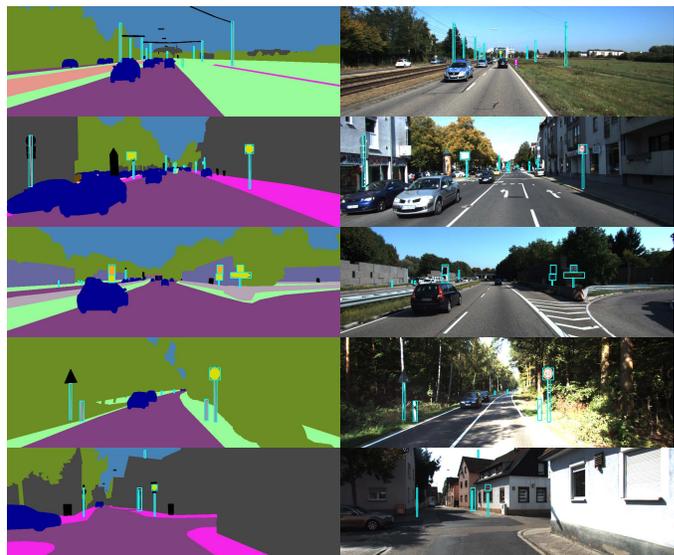[7]http://ori.ox.ac.uk/datasets/radar-robotcar-dataset



Fig. 13. Several examples of ground truth landmark bounding box labels based on semantic segmentation masks in the KITTI dataset. Left: semantic segmentation images with bounding box labels. Right: real images with bounding box labels.



Fig. 14. Several examples of the ground truth landmark bounding box labels for the Oxford Radar RobotCar dataset.

labels to object bounding box labels. We use the skimage.measure.label to label connected regions for pixel classes including traffic lights, traffic signs and poles. See Fig. 13 for an example. In some rare cases, multiple poles may overlap and the connected region algorithm outputs an inaccurate bounding box. We manually exclude these overlapped objects in the generated bounding box labels. As mentioned above, Faster R-CNN trained using the labels is used to produce the object detection results for all the other unlabeled frames contained in the dataset.

We project the surrounding LiDAR points onto the image frame plane using the intrinsic camera matrix and extrinsic

TABLE XII
NEURAL NETWORK MODELS AND THE PARAMETERS IN THE IMAGE PATCH MATCHING FRAMEWORK.

| Mapping | Models | Layers (model parameters) | Dimension of Outputs |
|---|---|---|---|
| $f$ | Resnet | Resnet18 (without the last FC layer, with 17 convolution layers) | 512 |
| $g$ | GAT/ GCN/ GraphSAGE | GAT block 1 (4 attention heads, $4 \times 128$ hidden features & ELU) | 512 |
| | | GAT block 2 (4 attention heads, $4 \times 128$ hidden features & ELU)/ | 512 |
| | | GCN block 1 (512 hidden features & ReLU) | |
| | | GCN block 2 (512 hidden features & ReLU)/ | |
| | | GraphSAGE block 1 ([512, 512] hidden features, ReLU & BatchNorm) | |
| | | GraphSAGE block 2 (512 hidden features) | |
| $d$ | Discrimiator | Bilinear layer (four $512 \times 512$ hidden partitioned matrices & Sigmoid function) | 1 |



(a) Landmark KITTI Dataset
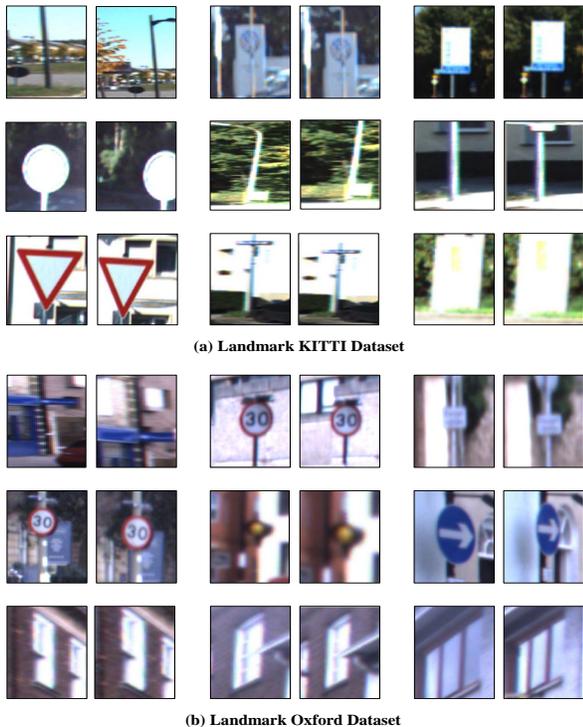


(b) Landmark Oxford Dataset

Fig. 15. (a) and (b) are examples of landmark patch pairs from the Landmark KITTI dataset and the Landmark Oxford Dataset respectively.

camera matrix. Here, we have used the sensors' information (i.e., vehicle global ground truth locations) to accumulate collected LiDAR scans to build the 3D LiDAR reference map. Due to the limited LiDAR field of view, a single LiDAR scan may not have any LiDAR point corresponding to some landmarks. To avoid this, we build a unified 3D LiDAR reference map similar to that in PointNetVLAD. Based on the 3D reference map, the LiDAR points reflected from the landmark patch are read out to obtain the global locations of the corresponding landmark objects. We apply DBSCAN to filter out some outlier points and obtain compact landmark objects. We then compute the $\mathcal{L}_2$ distance of each landmark patch pair from two frames to determine the patch matching ground truth. We have also gone through all the frames manually to remove or correct a few noisy landmark objects. Finally, for each detected landmark object, we intentionally expand its bounding box by 15 pixels on each side to include some background information. See Fig. 15 for an example.

**Landmark Oxford Dataset.** To build the Landmark Oxford

dataset, we manually labeled landmarks including traffic lights, traffic signs, poles, and facade windows for 500 frames. See Fig. 14 for examples. Compared with the Landmark KITTI Dataset, we additionally include the window class in this dataset. (Window labels are not available for the Landmark KITTI Dataset yet. We will enrich the Landmark KITTI Dataset with window labels in future work.) We then train Faster R-CNN to obtain the landmarks for all 29,687 frames. Similar operations are performed to obtain the final landmark patches with matching ground truths. See Fig. 15 for some landmark patch examples.

### B. Detailed Model Parameters

The details of the model setting mentioned in Section IV-B of the paper are provided in the following Table XII.

### C. Monocular Depth Estimation for VGIDM

In our work, we assume the spatial information of the segmentation is available to construct the neighborhood graphs in VGIDM. In Section IV of the paper, we perform evaluations on the two landmark datasets where Monocular Depth Prediction Module is used to obtain the spacial relationships among landmark patches contained in full-size images. The reported AbsRel of this depth estimation method is around 14 meters. Fig. 16 shows a few examples of the predicted depth. We observe that many objects in the predicted depth visualization are well distinguished from their surroundings.
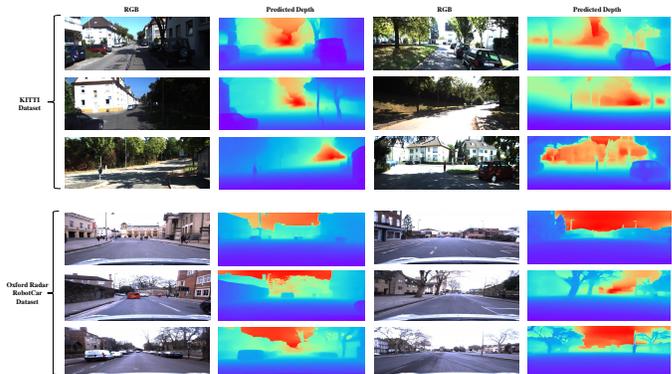


Fig. 16. Image depth estimation results from the Monocular Depth Prediction Module for the KITTI dataset and the Oxford Radar RobotCar Dataset.