# Pix2Map: Cross-modal Retrieval for Inferring Street Maps from Images

Xindi Wu [1*]    KwunFung Lau[1†]    Francesco Ferroni[2‡]    Aljoša Ošep[1]    Deva Ramanan[1,2]

[1]Carnegie Mellon University    [2]Argo AI

xindiw@princeton.edu, kwun.fung.lau@intel.com, fferroni@nvidia.com, {aosep, deva}@andrew.cmu.edu

pix2map.github.io

## Abstract

*Self-driving vehicles rely on urban street maps for autonomous navigation. In this paper, we introduce Pix2Map, a method for inferring urban street map topology directly from ego-view images, as needed to continually update and expand existing maps. This is a challenging task, as we need to infer a complex urban road topology directly from raw image data. The main insight of this paper is that this problem can be posed as cross-modal retrieval by learning a joint, cross-modal embedding space for images and existing maps, represented as discrete graphs that encode the topological layout of the visual surroundings. We conduct our experimental evaluation using the Argoverse dataset and show that it is indeed possible to accurately retrieve street maps corresponding to both seen and unseen roads solely from image data. Moreover, we show that our retrieved maps can be used to update or expand existing maps and even show proof-of-concept results for visual localization and image retrieval from spatial graphs.*

## 1. Introduction

We propose *Pix2Map*, a method for inferring road maps directly from images. More precisely, given the camera images, *Pix2Map* generates a topological map of the visible surroundings, represented as a spatial graph. Such maps encode both geometric and semantic scene information such as lane-level boundaries and locations of signs [53] and serve as powerful *priors* in virtually all autonomous vehicle stacks. In conjunction with on-the-fly sensory measurements from lidar or camera, such maps can be used for localization [3] and path planning [40]. As map maintenance and expansion to novel areas are challenging and expensive, often requiring manual effort [39, 51], automated map maintenance and expansion has been gaining interest in the community [10, 11, 28, 33, 34, 36, 39, 41].
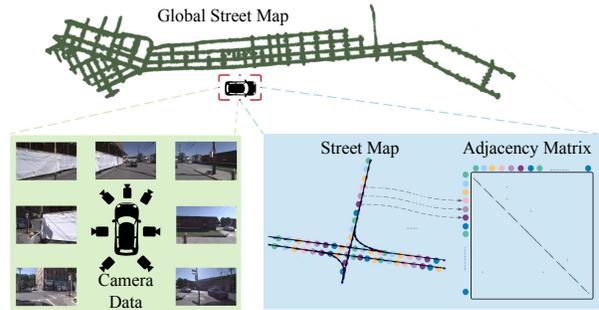


Figure 1. **Illustration of our proposed Pix2Map for cross-modal retrieval.** Given unseen 360° ego-view images collected from seven ring cameras (*left*), *Pix2Map* predicts the local street map by retrieving from the existing street map library (*right*), represented as an adjacency matrix. The local street maps can be further used for global high-definition map maintenance (*top*).

**Why is it hard?** To estimate urban street maps, we need to learn to map continuous images from ring cameras to discrete graphs with varying numbers of nodes and topology in bird's eye view (BEV). Prior works that estimate road topology from monocular images first process images using Convolutional Neural Networks or Transformers to extract road lanes and markings [10] or road centerlines [11] from images. These are used in conjunction with recurrent neural networks for the generation of polygonal structures [13] or heuristic post-processing [34] to estimate a spatial graph in BEV. This is a very difficult learning problem: such methods need to jointly learn to estimate a non-linear mapping from image pixels to BEV, as well as to estimate the road layout and learn to generate a discrete spatial graph.

**Pix2Map.** Instead, our core insight is to simply sidestep the problem of graph generation and 3D localization from monocular images by recasting *Pix2Map* as a *cross-modal retrieval* task: given a set of test-time ego-view images, we (i) compute their visual embedding and then (ii) retrieve a graph with the closest graph embedding in terms of cosine similarity. Given recent multi-city autonomous vehicle datasets [14], it is straightforward to construct pairs of ego-view images and street maps, both for training and testing.

---

[*]Now at Princeton, work done while at Carnegie Mellon University.
[†]Now at Intel.
[‡]Now at Nvidia.

We train image and graph encoders to operate in the same embedding space, making use of recent techniques for cross-modal contrastive learning [45]. Our key technical contribution is a novel but simple graph *encoder*, based on sequential transformers from the language community (*i.e.*, BERT [19]) that extract fixed-dimensional embeddings from street maps of arbitrary size and topology.

In fact, we find that even naive nearest-neighbor retrieval performs comparably to leading techniques for map generation [10, 11], *i.e.*, returning the graph paired with the best-matching image in the training set. Moreover, we demonstrate that cross-modal retrieval via *Pix2Map* performs even better due to its ability to learn graph embeddings that regularize the output space of graphs. In addition, cross-modal retrieval has the added benefit of allowing one to expand the retrieval graph library with *un*paired graphs that lack camera data, leveraging the insight that retrieval need not be limited to the same (image, graph) training pairs used for learning the encoders. This suggests that *Pix2Map* can be further improved with augmented road graph topologies that capture *potential* road graph updates (for which paired visual data might not yet exist). Beyond mapping, we show pilot experiments for visual localization, and the *inverse* method, *Map2Pix*, which retrieves a close-matching image from an image library given a graph. While not our primary focus, such approaches may be useful for generating photorealistic simulated worlds [41].

We summarize our **main contributions** as follows: We (i) show that dynamic street map construction from cameras can be posed as a cross-modal retrieval task and propose an contrastive image-graph model based on this framing. Building on recent advances in multimodal representation learning, we train a graph encoder and an image encoder with a shared latent space. We (ii) demonstrate empirically that this approach is effective and perform ablation studies to highlight the impacts of architectural decisions. Our approach outperforms existing graph generation methods from image cues by a large margin. We (iii) further show that it is possible to retrieve similar graphs to those in previously unseen areas without access to the ground truth graphs for those areas, and demonstrate the generalization ability to novel observations.

## 2. Related Work

Maps are ubiquitous in robotics: given a pre-built map of the environment, autonomous agents can localize themselves via live sensory data and plan their future trajectories [18]. Since the dawn of robotics, mapping and localization have been vibrant fields of research [54], tackled using different types of sensors, ranging from line-laser RGB-D sensors for indoor mapping [12, 24, 48], to lidar [4] and/or cameras [20, 21, 42], commonly used outdoors [9, 14, 22, 52]. In the following, we focus on map construction and main-

tenance. For localization, we refer to prior work [40, 54].

**Map Representation.** Several map representations have been proposed in the community, ranging from full 3D maps, represented as meshes [48, 55], voxel grids [12, 24, 32, 56], and (semantic) point clouds [4, 16]. In visual localization [50], point clouds are often constructed using structure-from-motion methods [50] and additionally store visual descriptors that aid matching-based visual localization. The aforementioned representations can be used for highly accurate 6-DoF camera localization. However, they are heavy in storage (and, consequentially, transmission) [62], which limits their applicability in outdoor environments.

**High-Definition (HD) Maps.** Alternatively, High Definition (HD) maps store key semantic information, such as road layout and traffic light sign positions [53], together with their attributes and connectivity information. As shown in [40], such sparse and storage-efficient maps can be used as priors for centimeter-precise vehicle location in conjunction with vehicle sensors, such as cameras and lidars. While immensely useful, HD maps are difficult to create and maintain [27, 39, 51] and often require manual annotations and post-processing, rendering map construction and maintenance costly. Therefore, a problem of great importance is the automation of map construction and maintenance directly from sensory data.

**HD Map Construction and Maintenance.** Several methods for HD map estimation rely on various sensor modalities. Li *et al*. [35] propose a method that generates a topological map (represented as a spatial graph) of a city from satellite images. Wang *et al*. [57] propose a collaborative approach that fuses several sources of information (data from airplanes, drones, and cars), such that consequent manual human post-processing can be minimized. Several methods tackle map construction directly from on-board vehicle sensory data. Often, methods tackle this challenging problem by first detecting road features in images (*e.g.*, segment lanes) [5, 15, 38, 58, 59], and then utilize the camera and lidar sensory data to estimate precise road layout in 3D space. To generate spatial graphs, the aforementioned methods employ generative recurrent neural networks [28, 36], or optimization-based approaches [37]. Unlike the aforementioned works, we estimate spatial graphs directly from images, by-passing explicit lane estimation.

**Pixel Segmentation.** State-of-the-art methods for graph generation from monocular images [10, 11] tend to first segment road lanes or centerlines in images, followed by graph generation using Polygon-RNN [13]. Instead of road lanes, HDMapNet [34] utilize methods for semantic/instance BEV maps (from cameras and/or lidar *e.g.*, [23, 29, 47, 49, 60]), followed by heuristic post-processing to obtain vectorized HD maps. HDMapGen [41] builds on recent developments in generative graph modeling [61] to construct con-
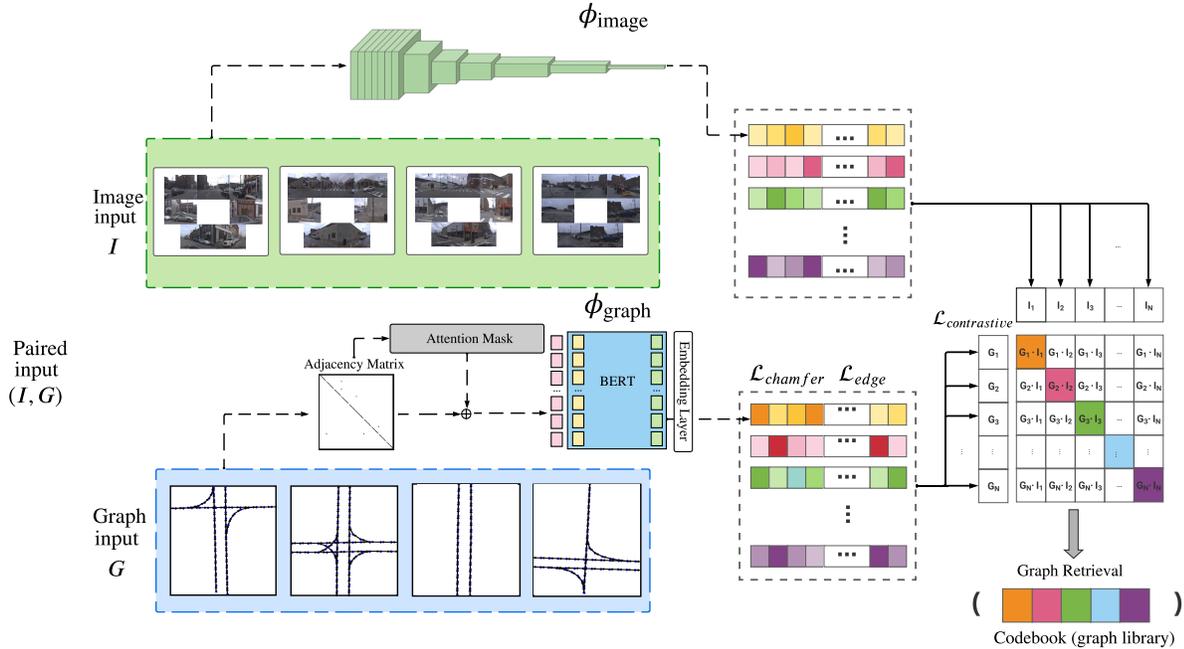
Figure 2. **Pix2Map**: The *graph encoder (bottom)* computes a graph embedding vector $\phi_{\text{graph}}$ for each street map in a batch. The *image encoder, (top)* outputs an image embedding $\phi_{\text{image}}$ for each corresponding image stack. We then build a similarity matrix for a batch, that contrasts the image and graph embeddings. We highlight that the adjacency matrix of a given graph is used as the attention mask for our transformer-based graph encoder.

trol points of central lane lines and their connectivity in a hierarchical manner. However, generating graphs *conditioned* on a particular (image) input remains an open problem. Unlike HDMapGen, *Pix2Map* sidesteps generative modeling, and instead directly retrieves a graph from a large database whose embedding vector is *most similar* to image embeddings in terms of cosine distance. We show that *Pix2Map* can also be used to keep HD maps up-to-date [6, 7, 33, 44].

## 3. Method

In this section, we formalize our *Pix2Map* approach as a cross-modal retrieval task. As shown in Fig. 2, given training pairs of images and graphs, we learn image and graph encoders that map both inputs to a common fixed-dimensional space via contrastive learning. We then use the learned encoders to retrieve a graph (from a training library) with the most similar embedding to the test image.

### 3.1. Problem Formulation

We construct a library of image-graph pairs $(I, G)$ where $I$ is a list of 7 ego-view images from a camera ring and $G$ is a street map represented as a graph $G = (V, E)$. Each vertex $v \in V$ represents a lane *node* and $E \in \{0, 1\}^{|V| \times |V|}$ that encodes the connectivity between nodes, stored in an adjacency matrix. Lane nodes have a position attribute $(x, y)$ in a local egocentric "birds-eye-view" coordinate frame, such that $(0, 0)$ is the ego-vehicle location (Fig. 3). Im-
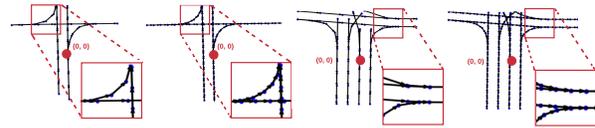


Figure 3. Two street map examples from Pittsburgh (*left*) and Miami (*right*) as *segment graphs* and our resampled *node graph*, with ego-vehicle origin at $(0, 0)$.

portantly, different graphs $G$ may have different numbers of lane nodes and connectivity information.

**Graph Representation.** The Argoverse dataset represents a street map as a *segment graph*, where each vertex represents a lane *segment*. Lane segments are represented as polylines with 10 $(x, y)$ points. We convert this *segment graph* to a *node graph* by defining each $(x, y)$ point as a graph node and adding a directed edge between successive points in a polyline (Fig. 3). We further resample the *segment graphs* by fitting degree-3 spline curves to lane segments, ensuring that connected nodes throughout the graph are approximately equidistant $(2m)$. We use this library for training the image and graph encoders, as detailed below.

### 3.2. Image Encoder

Given an image $I$, we use ResNet18 [26] as a feature extractor (without fully connected layers) as an image encoder that learns fixed-dimensional embedding vectors $\phi_{\text{image}}(I) \in R^{512}$. To process $n$ input images (we

3

use $n = 7$ images throughout our work), we stack them channel-wise. We experiment with ImageNet-pre-trained weights, as well as training "from scratch". In the pre-trained case, we replace the first convolution layer with one that stacks the original pre-trained filters $n$ times, with each weight divided by $n$. We reuse the original weights when making the new convolutional filters so the benefits of the pre-trained weights will be preserved. We ablate different training strategies and encoder architectures in Sec. 4.

### 3.3. Graph Encoder

Given a graph $G = (V, E)$, we would like to produce a fixed dimensional embedding $\phi_{\text{graph}}(G) \in R^{512}$, which is invariant to the orderings of the graph nodes. Unlike pixels in an image or words in a sentence, nodes in graphs do not have an inherent order. We construct such a graph encoder using a Transformer architecture inspired by sequence-to-sequence architectures from the language [19] defined on sequential tokens. Our encoder treats lane nodes as a collection of tokens and edges as masks for attention processing:

$$v_{l+1} = \sum_{\{w:E(v,w)=1\}} \text{Value}(v_l)\text{Softmax}_w[\text{Query}(v_l)\text{Key}(w)], \quad (1)$$

where $v_l$ is the embedding for vertex $v$ at layer $l$ and $v_0$ is initialized to its position $(x, y)$. We omit multiple attention heads and layer norm operations for brevity.

Embeddings are fed into a Transformer that computes new embeddings by taking an attention-weighted average of embeddings from nodes $w$ adjacent to $v$ (as encoded in the adjacency matrix $E$). We apply $M = 7$ transformer layers (similar to BERT [19]). Finally, we average (or mean pool) all output embeddings to produce a final fixed-dimensional embedding for graph $G$, regardless of the number of nodes or their connectivity: $\phi_{\text{graph}}(G) = \frac{1}{|V|} \sum_{v \in V} v_M \in R^{512}$. In Sec. 4 we ablate various design choices for our graph encoder, including the usage of periodic positional embeddings and encoding the edge connectivity information.

### 3.4. Image-Graph Contrastive Learning

To learn a joint embedding space, we follow the cross-modal contrastive formalism of [45], and briefly describe it here for completeness. Given N image-graph pairs $(I, G)$ within a batch, our model jointly learns the encoders $\phi_{\text{image}}(\cdot)$ and $\phi_{\text{graph}}(\cdot)$ such that the cosine similarity of the N correct image-graph pairs will be high and the $N^2 - N$ incorrect pairs will be low. We define cosine similarity between image $i$ and graph $j$ as:

$$\alpha_{ij} = \frac{\langle \phi_{\text{image}}(I_i), \phi_{\text{graph}}(G_j) \rangle}{||\phi_{\text{image}}(I_i)|| ||\phi_{\text{graph}}(G_j)||}. \quad (2)$$

We then compute bidirectional contrastive losses composed of an image-to-graph loss $\ell^{(I \to G)}$ and a graph-to-image loss

$\ell^{(G \to I)}$, following the form of the InfoNCE loss [43]:

$$\ell_i^{(I \to G)} = -\log \frac{\exp \alpha_{ii}}{\sum_j \exp \alpha_{ij}}, \quad (3)$$

$$\ell_i^{(G \to I)} = -\log \frac{\exp \alpha_{ii}}{\sum_j \exp \alpha_{ji}}. \quad (4)$$

The contrastive loss is then computed as a weighted combination of the two, averaged over all positive image-graph pairs in each minibatch:

$$\ell_{contrastive} = \frac{1}{2N} \sum_{i=1}^{N} \left( \ell_i^{(I \to G)} + \ell_i^{(G \to I)} \right). \quad (5)$$

The above penalizes all incorrect image-graph pairs equally. We found it beneficial to penalize false matches between pairs with similar graphs (measured using graph metrics, Sec. 4) less severely, as similar graphs should intuitively have similar embeddings. We measure this similarity of the graphs after aligning the vertices. Formally, given a ground truth graph, $G_0$ and candidate match $G_i$, we first establish a correspondence between each vertex $v \in V_0$ and its closest match $\pi_i(v) = v_i \in V_i$ (in terms of Euclidean distance between vertices). Given such corresponding vertices, we compute both a Chamfer Distance [2] (CD) and a binary cross-entropy (BCE) loss between the ground-truth binary adjacency matrix $E_0$ and the permuted matrix $E_i$:

$$\ell_{chamfer} = \sum_{v \in V_0} \sum_i \alpha_i \text{Distance}(v, \pi_i(v)), \quad (6)$$

$$\ell_{edge} = \sum_{v,w \in V_0} \text{BCE}(\sum_i \alpha_i E_i(\pi_i(v), \pi_i(w)) + \epsilon, E_0(v, w)), (7)$$

where $\alpha_i = \text{softmax}_i \, \alpha_{i0}$. The final loss is then:

$$\ell = \omega_1 \ell_{contrastive} + \omega_2 \ell_{chamfer} + \omega_3 \ell_{edge}, \quad (8)$$

where $\omega_1 = 1, \omega_2 = 1, \omega_3 = 1/10$. To ensure that the BCE loss remains finite, we add a small non-zero $\epsilon$ to ensure that the edge probabilities are strictly positive. To speed up the loss computation, we ignore edges $v, w \in V_0$ that are missing for all graphs in the batch $E_i(\pi_i(v), \pi_i(w)) = 0, \forall i$. A key difficulty in evaluating graph edge losses such as our own or the Rand Loss (described in Sec. 4.2) is that they assume that a vertex-wise correspondence is already known between the predicted and target graphs. A more theoretically optimal framework may *search* over one-to-one vertex correspondences that jointly minimize the Chamfer and edge loss, *e.g.*, by solving a bipartite matching problem [30].

### 3.5. Pix2Map via Cross-Modal Retrieval

Given the learned encodings above, we now use them for regressing maps from pixel image input via retrieval.

4

Denoting a graph library as $\mathbb{G}$, we retrieve image $I \Rightarrow G^*$, where $G^* = \mathrm{argmax}_{G \in \mathbb{G}_{\text{retrieval}}} \langle \phi_{\text{image}}(I), \phi_{\text{graph}}(G) \rangle$. Note that the graph library used for retrieval $\mathbb{G}_{\text{retrieval}}$ need not be the same as the one used to train the image-graph encoders. Formally, let encoders be trained on a collection of image-graph pairs, written as $\mathbb{D}_{\text{train}} := \{(I, G)|I \in \mathbb{I}_{\text{train}}, G \in \mathbb{G}_{\text{train}}\}$. $\mathbb{G}_{\text{retrieval}}$ need not be equivalent to $\mathbb{G}_{\text{train}}$, and furthermore, the set of corresponding images $\mathbb{I}_{\text{retrieval}}$ is not needed.

This has several important properties. *Firstly*, we can populate the graph library $\mathbb{G}_{\text{retrieval}}$ with additional graphs, not available during training, or cull the library to a particular subset of training graphs corresponding to a given city neighborhood. *Secondly*, we can enlarge the graph library with graphs that have no corresponding images, including augmented variants of real street maps that capture potential map updates (such as the potential addition of a lane at a particular intersection, for which no real-world imagery would be available). *Thirdly*, the above algorithm returns a ranked list of graphs, including near-ties. This can be used to generate multiple graphs that could correspond to an image input. *Finally*, similar observations hold for retrieving street map images using a graph (i.e., *Map2Pix*). *Map2Pix* is more likely to be a one-to-many task, as the same street geometry can be associated with different visual pixels depending on the time of day or weather conditions.

## 4. Experiments

In this section, we first discuss our evaluation test-bed (Sec. 4.1) that we use to conduct the experimental evaluation. Then, we perform ablation studies to highlight each component's contribution (Sec. 4.4). We compare our method to a recent state-of-the-art in Sec. 4.3 and, finally, highlight several use cases of our *Pix2Map* to automated map maintenance and expansion, and vehicle localization.

### 4.1. Evaluation Test-Bed

**Dataset.** For evaluation we use Argoverse dataset [14], which provides seven ring camera images ($1920 \times 1200$) recorded at 30 Hz with overlapping fields of view, providing $360°$ coverage. Crucially, Argoverse contains street maps that capture the geometry and connectivity of road lanes. Such map annotations are not available in other autonomous vehicle datasets such as nuScenes [9]. We perform the experiments across two cities in the United States, including Pittsburgh ($86km$) and Miami ($204km$).

**Splits.** Argoverse provides train, validation, and test splits. Note that validation and test splits may include regions that spatially overlap with the regions included in the training set. However, recordings of these regions were collected at different data collection runs at different times. To evaluate realistic applications of *map-updating* (where one trains on, *e.g.*, Pittsburgh up to 2021 and tests on Pittsburgh 2022+)

and *map-expansion* (where one trains on the neighborhood of Squirrel Hill and tests on Shadyside), we split up the union of (test+val) into those regions that spatially overlap the trainset and those that do not. We refer to these as *MapUpdate* and *MapExpand* test sets (Fig. 5). We present results for both settings, but default to *MapUpdate* for diagnostics unless otherwise specified.

**Map Preprocessing.** The key component of HD maps is the central line of drivable lanes. We extract the subgraphs corresponding to $40m \times 40m$ spatial windows. We use the adjacency matrix to represent the node connectivity. An edge connects two nodes if they are immediately reachable following the traffic flow *i.e.*, a subgraph of nodes in a given lane corresponds to a directed path. Moreover, an edge exists between two lanes if the first node of the second lane follows directly from the last node of the first, either because one lane continues to another or because one can turn from one lane to the other. We make sure to rotate the node positions and lanes to align with the driving direction.

**Implementation.** We train on a single NVIDIA A100 GPU, and the training dataset contains up to 512 samples in one batch. The model is trained for a total of 40 epochs, where a single epoch takes 40 minutes of wall-clock time. We use the Adam optimizer with a learning rate of 2e-4. We use a pretrained ResNet18 for the image encoder. In order to support an input containing several images, we duplicate and stack the filters of the input conv layer corresponding to the number of images. We then divide the parameters by the number of images per input example, giving us a model that initially returns an identical output to the original model. We extract the feature representation immediately before the fully connected layer. For the graph encoder, we use a Bert model with mean pooling and no positional embeddings to model the pairwise intersections between each of the nodes. For each node, we pass its adjacencies and its coordinates. We apply the attention mask which indicates to the model which tokens should be attended to and which should not.

### 4.2. Metrics

To quantitatively evaluate the quality of the retrieved graphs, we design three types of metrics to capture the difference between the retrieved graph $G_1 = (V_1, E_1)$ and the ground truth $G_2 = (V_2, E_2)$.

**Spatial Point Discrepancy.** We first introduce metrics that represent lane graph nodes $v \in V$ as $(x, y)$ points corresponding to the lane centroid, ignoring edge connectivity. We can then use metrics to measure differences between set of points. *Chamfer Distance* computes the closest point in $v_2 \in V_2$ for every $v_1 \in V_1$ (and vice versa, to ensure symmetry). *Maximum Mean Discrepancy (MMD)* [25, 41] measures the squared distance between point centroids in

| Methods | Chamfer $10^1$ | RandLoss $10^{-2}$ | MMD $10^{-1}$ | U. density $10^{-1}$ | U. reach $10^{-1}$ | U. conn. $10^{-1}$ |
|---|---|---|---|---|---|---|
| PINET [31] | 4.9244 | 10.8935 | 4.2983 | 2.8194 | 7.4194 | 2.9231 |
| TOPO-PRNN [11] | 7.4811 | 9.2813 | 5.7726 | 3.9371 | 6.8297 | 1.3934 |
| TOPO-TR [11] | 3.0140 | **7.1603** | 4.6431 | 2.2467 | 3.3091 | 1.1530 |
| *Pix2Map*-Unimodal | 4.3967 | 9.0764 | 4.1873 | 1.8391 | 3.2746 | 1.7734 |
| *Pix2Map*-Single | 2.6819 | 7.5204 | 4.0848 | 2.5339 | **3.0134** | **1.0291** |
| *Pix2Map* (ours) | **2.0882** | 7.7562 | **3.9621** | **1.4354** | 3.2893 | 1.5532 |

Table 1. **Baseline comparisons.** For fair comparisons with the prior art [11], in this experiment, we (i) train *Pix2Map* using frontal $50m \times 50m$ road-graphs (as opposed to our default setting of predicting the surrounding $40m \times 40m$ area). Moreover, we (ii) train *Pix2Map* with a single frontal view (*Pix2Map*-Single) to ensure consistent comparisons to baselines. Importantly, even in this setting, **our method still outperforms baselines by a large margin**: 2.6819 in terms of Chamfer distance, as compared to 3.0140 obtained by the closest competitor, TOPO-TR [11].

a Hilbert space using Gaussian kernels $\langle \varphi(v_1), \varphi(v_2) \rangle_H = k(x_1 - x_2, y_1 - y_2)$:

$$\text{MMD}(G_1, G_2) = \left\| \frac{1}{|V_1|} \sum_{v_1 \in V_1} \varphi(v_1) - \frac{1}{|V_2|} \sum_{v_2 \in V_2} \varphi(v_2) \right\|_H^2 .$$

**Edge Connectivity.** The above metrics evaluate the quality of only the retrieved graph nodes, but not their edge connectivity. We define a RandLoss similar to (6), as:

$$\text{RandLoss} = \sum_{v, w \in V_1} \mathbb{1}_{[E_2(\pi(v), \pi(w)) \neq E_1(v, w)]},$$

where $\mathbb{1}$ is an indicator function for mismatching edge labels between a pair of nodes in the graph $G_1$ and their corresponding pair in graph $G_2$. This metric is also known as the Rand index for clustering evaluation [46].

**Urban Planning.** We also report a set of metrics motivated by the urban planning literature [1, 17, 41], evaluating the degree to which we are able to reconstruct the following key properties of urban HD maps. **Connectivity** is the number of edges relative to the number of lane nodes. **Density** is the number of edges relative to the maximum possible number of edges. **Reach** is designed to capture urban development and is defined as the total distance covered by the lanes:

$$\text{Connectivity} = \frac{\|E\|_0}{|V|}, \quad \text{Density} = \frac{\|E\|_0}{|V|(|V| - 1)},$$
$$\text{Reach} = \sum_{(v, w): E(v, w) = 1} \text{len}(v, w).$$

We report the absolute relative error [41] for these metrics.

### 4.3. Baselines

We show a visual comparison with top performers in Fig. 4 and quantitative results in Tab. 1. We first report the performance of a naive nearest-neighbor baseline (Unimodal), which returns the graph associated with the closest training image example. This unimodal approach already performs on par with the state-of-the-art Transformer
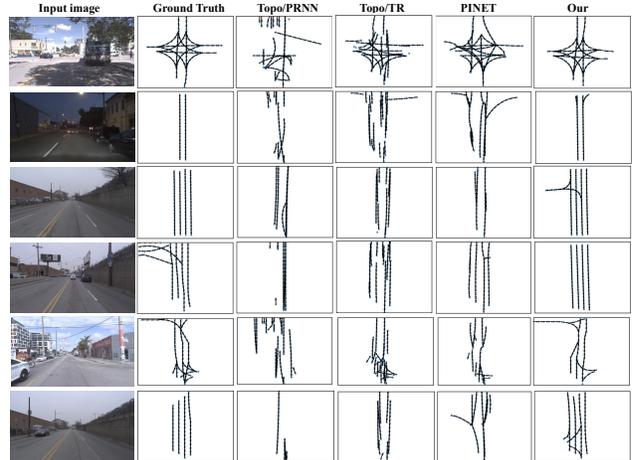


Figure 4. **Qualitative results**. From left to right: input image, ground-truth maps, maps generated by state-of-the-art methods, and, in the last column, our method. As can be seen, the retrieved maps with our method have the highest visual fidelity.

and Polygon-RNN [27] based methods TOPO [10], and PINET [31]. Our diagnostics further explore the improvement from unimodal to cross-modal retrieval. As can be seen, *Pix2Map* improves greatly over several state-of-the-art baselines in terms of the Chamfer distance, MMD, urban density error, and urban connectivity error, while performing comparably in terms of RandLoss and urban reach error. Moreover, our method is especially strong in terms of preserving the spatial point discrepancy, outperforming baselines by a large margin. We note that *Pix2Map* was designed to fully utilize the image data available in the camera ring, whereas baselines use only frontal view. For apples-to-apples comparison, we retrain *Pix2Map* with a single frontal view (*i.e.*, single camera, *Pix2Map*-Single. As can be seen, even the single-view variant of *Pix2Map* outperforms the closest competitor (TOPO-TR) across almost all metrics except RandLoss and Urban density. Our results suggest that baselines may also benefit from multi-view processing.

### 4.4. Model Analysis

In this section, we experiment with two graph representations, evaluate different design decisions on the graph and image encoders, and multimodal contrastive learning.

**Ablations.** Tab. 2 ablates several design decisions on image (Sec. 3.2) and graph (Sec. 3.3) encoders. For the analysis, we focus on Chamfer distance as an illustrative metric, as it appears to be most consistent with graph qualitative estimation.

*Image Encoder.* As can be seen in Tab. 2, row3-vs-row6, *early fusion* (1×RN18) is significantly better compared to a *late fusion* (7×RN18) variant, where we separately encode images with a distinct ResNet18 model for each camera, followed by fusion via average-pooling.

| Row | $\mathcal{E}_{img}$ | Attention Mask | Adjacency Matrix | Positional Encoding | Resampling | Chamfer $10^1$ | RandLoss $10^{-2}$ | MMD $10^{-1}$ | U. density $10^{-1}$ | U. reach $10^{-1}$ | U. conn. $10^{-1}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | $1 \times$ RN18 | ✓ | | | ✓ | 1.9241 | 9.0446 | 4.1804 | 1.2058 | 3.6333 | 1.8398 |
| 2 | $1 \times$ RN18 | | ✓ | | ✓ | 1.9309 | 8.3834 | 3.8383 | 1.1934 | 3.7428 | 1.8629 |
| 3 | $1 \times$ RN18 | ✓ | ✓ | | ✓ | **1.5908** | 7.3283 | **3.0888** | **0.7593** | **3.2997** | **0.8397** |
| 4 | $1 \times$ RN18 | ✓ | ✓ | | | 3.2663 | **6.9943** | 6.3704 | 3.6883 | 5.3219 | 4.1658 |
| 5 | $1 \times$ RN18 | ✓ | ✓ | ✓ | ✓ | 2.1564 | 9.1200 | 8.8328 | 0.8813 | 3.4290 | 1.5481 |
| 6 | $7 \times$ RN18 | ✓ | ✓ | | ✓ | 4.8129 | 11.2118 | 9.7538 | 3.9169 | 6.7794 | 2.5285 |

Table 2. **Image and graph encoder ablations.** From left to right, we ablate a) encoding each of the seven ego-images separately or using an early-fusion multiview image encoder, b) restricting the transformer attention mask to the graph-adjacency matrix or using the default fully-connected attention, c) the inclusion of the corresponding row of the attention matrix as a node input feature for the model, d) adding a positional encoding to each graph vertex, and e) resampling graph vertices to be equidistant. We find results are dramatically improved by early fusion for image encoding (row3-vs-row6) and graph vertex resampling (row3-vs-row4). Results are marginally improved by restricting the attention mask (row2-vs-row3) and adding the adjacency matrix as an input feature (row1-vs-row3 and row2-vs-row3). Perhaps surprisingly, adding in positional vertex encodings slightly decreases performance (row3-vs-row5).
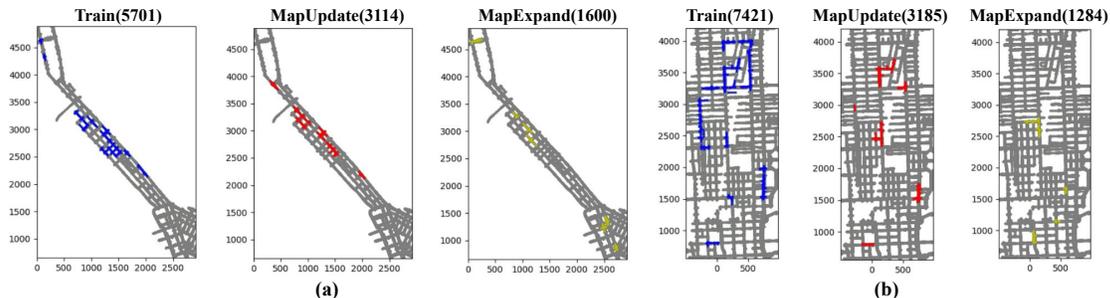


Figure 5. **Pittsburgh (a) and Miami (b) datasets.** Including training data (blue), MapUpdate test data (red) that overlap the blue but are collected at different times, and MapExpand test data (yellow) that do not non-overlap. We denote the size of each dataset (in terms of the number of image-graph pairs) in parentheses.

| $\mathcal{L}_{con.}$ | $\mathcal{L}_{edge}$ | $\mathcal{L}_{cham f.}$ | Chamfer $10^1$ | RandLoss $10^{-2}$ | MMD $10^{-1}$ | U. density $10^{-1}$ | U. reach $10^{-1}$ | U. conn. $10^{-1}$ |
|---|---|---|---|---|---|---|---|---|
| ✓ | | | 3.2249 | 9.1727 | 7.1070 | 2.7387 | 14.1033 | 3.6998 |
| ✓ | ✓ | | 2.7967 | 8.9951 | 8.8717 | 0.9808 | 10.0251 | 1.0512 |
| ✓ | | ✓ | 1.7440 | 8.5186 | **2.2480** | 0.9841 | 7.4534 | 1.5664 |
| ✓ | ✓ | ✓ | **1.5908** | **7.3283** | 3.0888 | **0.7593** | **3.2997** | **0.8397** |

Table 3. **Ablation on the training loss.** Adding a partial credit for matching to graphs with low Chamfer distance ($\mathcal{L}_{cham.}$) to the ground truth improves results considerably compared to vanilla contrastive loss ($\mathcal{L}_{con.}$). By adding in an additional edge loss ($\mathcal{L}_{edge}$), we further improve the performance.

| Methods | Chamfer $10^1$ | RandLoss $10^{-2}$ | MMD $10^{-1}$ | U. density $10^{-1}$ | U. reach $10^{-1}$ | U. conn. $10^{-1}$ |
|---|---|---|---|---|---|---|
| Unimodal | 3.2168 | 9.7596 | 7.7671 | 0.7365 | 3.9452 | 1.3661 |
| Ours | 1.5908 | 7.3283 | 3.0888 | 0.7593 | 3.2997 | 0.8397 |
| Ours++ | 1.5208 | 6.1504 | 3.0944 | 0.7407 | 3.2610 | 0.8089 |

Table 4. Cross-model retrieval (**ours**) significantly outperforms classical **unimodal** retrieval (*i.e.*, the nearest neighbor on image encoder features). Cross-modal retrieval can exploit the graph embedding space, which appears to regularize retrieval results, while the unimodal approach does not utilize any graph embedding. Moreover, our cross-modal retrieval can take advantage of larger unpaired graph libraries, which further improve performance (**ours++**). Unimodal retrieval requires paired data.

*Graph Encoder.* We consider three primary variations on the graph encoder architecture. The results suggest (Tab. 2, row2-vs-row3) that restricting the attention mask

| City | Library Size $10^1$ | Chamfer $10^{-2}$ | RandLoss $10^{-1}$ | MMD $10^{-1}$ | U. density $10^{-1}$ | U. reach $10^{-1}$ | U. conn. |
|---|---|---|---|---|---|---|---|
| PIT | 5.7k | 1.5908 | 7.3283 | 3.0888 | 0.7593 | 3.2997 | 0.8397 |
| | 10k | 1.6457 | 7.6247 | 3.2848 | **0.7264** | 4.5891 | 1.6364 |
| | 20k | 1.5369 | 6.5373 | 3.1883 | 0.7581 | 3.2902 | 1.0602 |
| | 30k | 1.5239 | 6.6553 | **3.0253** | 0.8586 | 4.0642 | 0.9615 |
| | 40k | **1.5208** | **6.1504** | 3.0944 | 0.7407 | **3.2610** | **0.8089** |
| MIA | 7.4k | 1.4747 | 6.8693 | 3.4033 | 1.0948 | 4.6253 | 1.1910 |
| | 10k | 1.4991 | 6.2315 | 3.3118 | 1.2784 | 5.5209 | 1.3679 |
| | 20k | 1.3878 | 8.0234 | 3.3910 | 1.1290 | 4.1237 | 1.3249 |
| | 30k | 1.4012 | 7.1898 | 3.2773 | 1.2444 | 4.2471 | 1.2298 |
| | 40k | 1.3878 | 7.6305 | 3.3351 | 1.2523 | 5.3894 | 1.3385 |
| | 60k | 1.3080 | 6.3369 | 3.18879 | 1.1972 | 4.7578 | 1.1977 |
| | 80k | 1.2711 | 6.2852 | 3.19506 | 1.0123 | 4.6827 | 1.1651 |
| | 100k | **1.2462** | **6.2740** | 3.1277 | **0.9884** | **3.8521** | **1.1397** |

Table 5. **Ablation with larger map-graph libraries.** As we grow the graph retrieval library (including maps without corresponding image views), we observe performance grows consistently with the size of the retrieval library.

marginally improves the results. Furthermore, encoding the edge connectivity information (adding the adjacency matrix as an input feature, row1-vs-row3) slightly improves the performance. Perhaps surprisingly, positional vertex encodings slightly decrease performance (row3-vs-row5).

*Graph Representation.* Tab. 2 (row3-vs-row4) shows that the proposed graph vertex resampling dramatically improves the results. As we mentioned in Sec. 3.1, we resample the *segment graphs* and the connected nodes throughout the graph to be approximately equidistant.
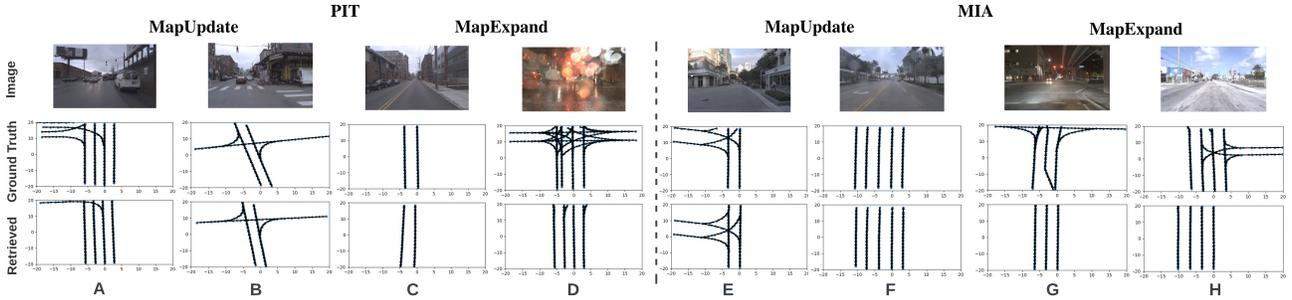
Figure 6. Given a set of ego-view images (front camera, *top*), we plot the ground-truth graph (*middle*), followed by the *Pix2Map* predictions (*bottom*) for both the **MapUpdate** (*left*) and **MapExpand** (*right*) tasks for two cities, Pittsburgh (**PIT**) and Miami (**MIA**). In general, retrieval results are quite accurate, particularly for **MapUpdate**, where train and test samples are drawn from the same geographic regions. Inclement weather such as heavy rain is challenging (row **PIT MapExpand**, column **D**) due to the degraded visual signal.

*Loss.* Tab. 3 ablates our loss function (Sec. 3.4). Compared to the naïve contrastive loss $\mathcal{L}_{contrastive}$, which weighs all incorrect image-graph matches equally, we find adding in partial credit for matching to graphs with low Chamfer distance $\mathcal{L}_{chamfer}$, to the ground-truth improves results considerably while adding in an additional edge loss $\mathcal{L}_{edge}$ further improves performance. *We use the most performant combination (row3) for further experiments.*

**Unimodal v.s. Cross-modal.** To quantify the benefits of the cross-modal training scheme, we compare *Pix2Map* to its image encoder alone, evaluating it as a unimodal image-encoding-based retriever. Specifically, in *Pix2Map*, we directly retrieve graphs by finding the graph embedding that is closest to the input image embedding in the multimodal embedding space. However, in this ablation, we instead find the image embedding in the training set which is closest to the input image embedding and return its corresponding graph. We find that using both modalities improves performance on almost all metrics, as shown in Tab. 4. The improvements range from a 37.0% decrease in RandLoss to a 60.2% decrease in MMD, with one increase of less than a percentage point in urban density error. However, note that while unimodal ablation performs broadly worse than *Pix2Map*, it still performs better than any baseline shown in Tab. 1 in terms of Chamfer, MMD, and Urban Density.

**Augmenting the Graph Library.** One of the benefits of our cross-modal retrieval approach is that we can match (or retrieve from) arbitrary collections of graphs that are different from (or larger than) the training graphs used to learn cross-modal encoders. This allows us to make use of graphs that do not have corresponding images. Interestingly, the Argoverse dataset provides such data, as maps include many locations for which no imagery is provided. By sampling random ego-vehicle positions in Miami and Pittsburgh, we grow both the Pittsburgh library and the Miami library to $40k$, thus significantly expanding our retrieval graph library. We summarize these results in Tab. 5. We observe that performance improves across the suite of metrics

as the graph retrieval library grows larger. This suggests there is a significant potential to further improve our results by simply growing our graph retrieval dataset using existing maps, without access to corresponding image pairs.

**Applications.** In the Appendix, we present several experiments that demonstrate the practical applications of *Pix2Map*. In particular, we present experiments on expanding (*MapExpand*) and updating (*MapUpdate*) existing maps using *Pix2Map*, see Fig. 6. Furthermore, we show that *Pix2Map* can be used for visual localization by generating a heatmap of possible ego-vehicle locations in a city-level map based on input images. Finally, we visually demonstrate the inverse *Map2Pix* task for retrieving ego-camera images given a street map. This experiment shows that it is possible to retrieve egocentric camera data using street maps, which can be used to synthesize virtual worlds consistent with the query road geometry.

## 5. Conclusion

In this work, we propose a significantly different approach to inferring high-definition maps from cameras. Rather than learning a nonlinear mapping from image pixels to BEV and generating a discrete spatial graph, we suggest a retrieval-based approach. Our experiments indicate that learning a multimodal embedding space for camera data and map data holds promise, and we hope our work can serve as an essential building block for map expansion and updating in the autonomous driving field. Beyond map maintenance, we also demonstrate that our approach can be used as a novel form of visual localization. While these results are encouraging, there are numerous potentially impactful future directions to explore. For instance, instead of performing graph retrieval, one could utilize the latent space to generate new unseen graphs using a graph-based decoder architecture.

# References

[1] Sawsan AlHalawani, Yong-Liang Yang, Peter Wonka, and Niloy J Mitra. What makes london work like london? *Computer Graphics Forum*, 2014. 6

[2] Harry G Barrow, Jay M Tenenbaum, Robert C Bolles, and Helen C Wolf. Parametric correspondence and chamfer matching: Two new techniques for image matching. Technical report, SRI International Menlo Park, CA AI Center, 1977. 4

[3] Ioan Andrei Barsan, Shenlong Wang, Andrei Pokrovsky, and Raquel Urtasun. Learning to localize using a lidar intensity map. *arXiv preprint arXiv:2012.10902*, 2020. 1

[4] Jens Behley and Cyrill Stachniss. Efficient Surfel-Based SLAM using 3D Laser Range Data in Urban Environments. In *RSS*, 2018. 2

[5] Karsten Behrendt and Ryan Soussan. Unsupervised labeled lane markers using maps. In *ICCV Workshops*, 2019. 2

[6] Julie Stephany Berrio, Stewart Worrall, Mao Shan, and Eduardo Nebot. Long-term map maintenance pipeline for autonomous vehicles. *IEEE TPAMI*, 2021. 3

[7] Sagar Ravi Bhavsar, Andrei Vatavu, Timo Rehfeld, and Gunther Krehl. Sensor fusion-based online map validation for autonomous driving. In *IVS*, 2020. 3

[8] Marcus A Brubaker, Andreas Geiger, and Raquel Urtasun. Map-based probabilistic visual self-localization. *IEEE TPAMI*, 2015. 1

[9] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *CVPR*, 2020. 2, 5

[10] Yigit Baran Can, Alexander Liniger, Danda Pani Paudel, and Luc Van Gool. Structured bird's-eye-view traffic scene understanding from onboard images. In *ICCV*, 2021. 1, 2, 6

[11] Yigit Baran Can, Alexander Liniger, Danda Pani Paudel, and Luc Van Gool. Topology preserving local road network estimation from single onboard camera image. In *CVPR*, 2022. 1, 2, 6

[12] Vincent Cartillier, Zhile Ren, Neha Jain, Stefan Lee, Irfan Essa, and Dhruv Batra. Semantic mapnet: Building allocentric semanticmaps and representations from egocentric views. *arXiv preprint arXiv:2010.01191*, 2020. 2

[13] Lluis Castrejon, Kaustav Kundu, Raquel Urtasun, and Sanja Fidler. Annotating object instances with a polygon-rnn. In *CVPR*, 2017. 1, 2

[14] Ming-Fang Chang, John Lambert, Patsorn Sangkloy, Jagjeet Singh, Slawomir Bak, Andrew Hartnett, De Wang, Peter Carr, Simon Lucey, Deva Ramanan, and James Hays. Argoverse: 3d tracking and forecasting with rich maps. In *CVPR*, 2019. 1, 2, 5

[15] Wensheng Cheng, Hao Luo, Wen Yang, Lei Yu, Shoushun Chen, and Wei Li. Det: A high-resolution dvs dataset for lane extraction. In *CVPR Workshops*, 2019. 2

[16] Sungjin Cho, Chansoo Kim, Jaehyun Park, Myoungho Sunwoo, and Kichun Jo. Semantic point cloud mapping of lidar based on probabilistic uncertainty modeling for autonomous driving. *Sensors*, 2020. 2

[17] Hang Chu, Daiqing Li, David Acuna, Amlan Kar, Maria Shugrina, Xinkai Wei, Ming-Yu Liu, Antonio Torralba, and Sanja Fidler. Neural turtle graphics for modeling city road layouts. In *ICCV*, 2019. 6

[18] Nachiket Deo, Eric Wolff, and Oscar Beijbom. Multimodal trajectory prediction conditioned on lane-graph traversals. In *CoRL*, 2022. 2

[19] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 2, 4

[20] Jakob Engel, Thomas Schöps, and Daniel Cremers. Lsd-slam: Large-scale direct monocular slam. In *ECCV*, 2014. 2

[21] Jakob Engel, Jörg Stückler, and Daniel Cremers. Large-scale direct slam with stereo cameras. In *IROS*, 2015. 2

[22] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the KITTI vision benchmark suite. In *CVPR*, 2012. 2

[23] Thomas Gilles, Stefano Sabatini, Dzmitry Tsishkou, Bogdan Stanciulescu, and Fabien Moutarde. Gohome: Graph-oriented heatmap output for future motion estimation. In *ICRA*, 2022. 2

[24] Margarita Grinvald, Fadri Furrer, Tonci Novkovic, Jen Jen Chung, Cesar Cadena, Roland Siegwart, and Juan Nieto. Volumetric instance-aware semantic mapping and 3d object discovery. *IEEE RAL*, 2019. 2

[25] Ehsan Hajiramezanali, Arman Hasanzadeh, Krishna Narayanan, Nick Duffield, Mingyuan Zhou, and Xiaoning Qian. Variational graph recurrent neural networks. In *NeurIPS*, 2019. 5

[26] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 3

[27] Namdar Homayounfar, Wei-Chiu Ma, Shrinidhi Kowshika Lakshmikanth, and Raquel Urtasun. Hierarchical recurrent attention networks for structured online maps. In *CVPR*, 2018. 2, 6

[28] Namdar Homayounfar, Wei-Chiu Ma, Justin Liang, Xinyu Wu, Jack Fan, and Raquel Urtasun. Dagmapper: Learning to map by discovering lane topology. In *ICCV*, 2019. 1, 2

[29] Anthony Hu, Zak Murez, Nikhil Mohan, Sofía Dudas, Jeffrey Hawke, Vijay Badrinarayanan, Roberto Cipolla, and Alex Kendall. Fiery: Future instance prediction in bird's-eye view from surround monocular cameras. In *ICCV*, 2021. 2

[30] Richard M Karp, Umesh V Vazirani, and Vijay V Vazirani. An optimal algorithm for on-line bipartite matching. In *ACM STOC*, 1990. 4

[31] Yeongmin Ko, Younkwan Lee, Shoaib Azam, Farzeen Munir, Moongu Jeon, and Witold Pedrycz. Key points estimation and point instance segmentation approach for lane detection. *IEEE Trans. ITS*, 2021. 6

[32] Deyvid Kochanov, Aljoša Ošep, Jörg Stückler, and Bastian Leibe. Scene flow propagation for semantic mapping and object discovery in dynamic street scenes. In *IROS*, 2016. 2

[33] John Lambert and James Hays. Trust, but verify: Cross-modality fusion for hd map change detection. In *NeurIPS*, 2021. 1, 3

[34] Qi Li, Yue Wang, Yilun Wang, and Hang Zhao. Hdmapnet: An online hd map construction and evaluation framework. *arXiv preprint arXiv:2107.06307*, 2021. 1, 2

[35] Zuoyue Li, Jan Dirk Wegner, and Aurélien Lucchi. Topological map extraction from overhead images. In *ICCV*, 2019. 2

[36] Justin Liang, Namdar Homayounfar, Wei-Chiu Ma, Shenlong Wang, and Raquel Urtasun. Convolutional recurrent network for road boundary extraction. In *CVPR*, 2019. 1, 2

[37] Martin Liebner, Dominik Jain, Julian Schauseil, David Pannen, and Andreas Hackelöer. Crowdsourced hd map patches based on road model inference and graph-based slam. In *IVS*, 2019. 2

[38] Lizhe Liu, Xiaohao Chen, Siyu Zhu, and Ping Tan. Condlanenet: a top-to-down lane detection framework based on conditional convolution. In *ICCV*, 2021. 2

[39] Rong Liu, Jinling Wang, and Bingqi Zhang. High definition map for automated driving: Overview and analysis. *The Journal of Navigation*, 2020. 1, 2

[40] Wei-Chiu Ma, Ignacio Tartavull, Ioan Andrei Bârsan, Shenlong Wang, Min Bai, Gellert Mattyus, Namdar Homayounfar, Shrinidhi Kowshika Lakshmikanth, Andrei Pokrovsky, and Raquel Urtasun. Exploiting sparse semantic hd maps for self-driving vehicle localization. In *IROS*, 2019. 1, 2

[41] Lu Mi, Hang Zhao, Charlie Nash, Xiaohan Jin, Jiyang Gao, Chen Sun, Cordelia Schmid, Nir Shavit, Yuning Chai, and Dragomir Anguelov. Hdmapgen: A hierarchical graph generative model of high definition maps. In *CVPR*, 2021. 1, 2, 5, 6

[42] Raul Mur-Artal, Jose Maria Martinez Montiel, and Juan D Tardos. Orb-slam: a versatile and accurate monocular slam system. *TRO*, 2015. 2

[43] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 4

[44] David Pannen, Martin Liebner, and Wolfram Burgard. Hd map change detection with a boosted particle filter. In *ICRA*, 2019. 3

[45] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021. 2, 4

[46] William M Rand. Objective criteria for the evaluation of clustering methods. *Journal of the American Statistical association*, 1971. 6

[47] Thomas Roddick and Roberto Cipolla. Predicting semantic map representations from images using pyramid occupancy networks. In *CVPR*, 2020. 2

[48] Radu Alexandru Rosu, Jan Quenzel, and Sven Behnke. Semi-supervised semantic mapping through label propagation with semantic texture meshes. *IJCV*, 2020. 2

[49] Avishkar Saha, Oscar Mendez, Chris Russell, and Richard Bowden. Translating images into maps. In *ICRA*, 2022. 2

[50] Torsten Sattler, Bastian Leibe, and Leif Kobbelt. Fast image-based localization using direct 2d-to-3d matching. In *ICCV*, 2011. 2

[51] Heiko G Seif and Xiaolong Hu. Autonomous driving in the icity—hd maps as a key challenge of the automotive industry. *Engineering*, 2016. 1, 2

[52] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving: Waymo open dataset. In *CVPR*, 2020. 2

[53] Sebastian Thrun. Robotic mapping: A survey. *Exploring artificial intelligence in the new millennium*, 2002. 1, 2

[54] Sebastian Thrun, Wolfram Burgard, and Dieter Fox. *Probabilistic Robotics (Intelligent Robotics and Autonomous Agents)*. The MIT Press, 2005. 2

[55] Julien PC Valentin, Sunando Sengupta, Jonathan Warrell, Ali Shahrokni, and Philip HS Torr. Mesh based semantic modelling for indoor and outdoor scenes. In *CVPR*, 2013. 2

[56] Vibhav Vineet, Ondrej Miksik, Morten Lidegaard, Matthias Nießner, Stuart Golodetz, Victor A Prisacariu, Olaf Kähler, David W Murray, Shahram Izadi, Patrick Pérez, et al. Incremental dense semantic stereo fusion for large-scale semantic scene reconstruction. In *ICRA*, 2015. 2

[57] Shenlong Wang, Min Bai, Gellert Mattyus, Hang Chu, Wenjie Luo, Bin Yang, Justin Liang, Joel Cheverie, Sanja Fidler, and Raquel Urtasun. Torontocity: Seeing the world with a million eyes. *arXiv preprint arXiv:1612.00423*, 2016. 2

[58] Dong Wu, Manwen Liao, Weitian Zhang, and Xinggang Wang. Yolop: You only look once for panoptic driving perception. *arXiv preprint arXiv:2108.11250*, 2021. 2

[59] Pan Xingang, Shi Jianping, Luo Ping, Wang Xiaogang, and Tang Xiaoou. Spatial as deep: Spatial cnn for traffic scene understanding. In *AAAI*, 2018. 2

[60] Weixiang Yang, Qi Li, Wenxi Liu, Yuanlong Yu, Yuexin Ma, Shengfeng He, and Jia Pan. Projecting your view attentively: Monocular road scene layout estimation via cross-view transformation. In *CVPR*, 2021. 2

[61] Jiaxuan You, Rex Ying, Xiang Ren, William Hamilton, and Jure Leskovec. Graphrnn: Generating realistic graphs with deep auto-regressive models. In *ICML*, 2018. 2

[62] Qunjie Zhou, Sergio Agostinho, Aljosa Osep, and Laura Leal-Taixe. Is geometry enough for matching in visual localization? In *ECCV*, 2022. 2

# Supplementary material for Pix2map:
# Cross-modal Retrieval for Inferring Street Maps from Images

## Abstract

*In this supplement, we provide various experiments to illustrate the practical uses of Pix2Map. These experiments include:*

- *Map Expansion and Update, in which we present experiments on expanding and updating existing maps,*

- *Visual Localization, by generating a heatmap of possible locations for the ego-vehicle on a city-level map,*

- *Map2Pix, which is visually demonstrated by retrieving ego-camera images using street maps.*

## A. Applications

In this section, we discuss how our method can be used for practical purposes, and show that graph library retrieval can greatly improve various downstream applications such as expansion (*MapExpand*) and update (*MapUpdate*) given existing maps, visual image-to-HD map localization and *Map2Pix*.

### A.1. Map Expansion and Update

We use our graph retrieval method to mimic map expansion (*MapExpand*) and map update (*MapUpdate*) using data splits. For map expansion, we retrieve local graphs corresponding to recordings obtained in a "new traversal" to expand the existing map. For map updates, we similarly retrieve local maps to update the global map.

| City | Task type | Chamfer $10^1$ | RandLoss $10^{-2}$ | MMD $10^{-1}$ | U. density $10^{-1}$ | U. reach $10^{-1}$ | U. conn. $10^{-1}$ |
|------|-----------|---------|----------|-----|-----------|---------|---------|
| PIT | MapUpdate | 1.5908 | 7.3283 | 3.0888 | 0.7593 | 3.2997 | 0.8397 |
|     | MapExpand | 2.6654 | 16.9768 | 8.0468 | 3.9482 | 4.2949 | 3.9699 |
| MIA | MapUpdate | 1.4747 | 6.8693 | 3.4033 | 1.0948 | 5.5333 | 1.1910 |
|     | MapExpand | 2.0637 | 11.1354 | 4.2605 | 1.4922 | 4.7318 | 1.5940 |

Table 6. **Map update and expansion evaluation.** As can be seen, map expansion to novel areas can be much harder than updating previously-seen areas.

We qualitatively evaluate the graph retrieval results in Fig. 6 in the main paper. Please see the caption for a detailed description, but generally speaking, we find *Pix2Map* returns reasonable graphs similar to the ground truth. In Tab. 6, we evaluate the performance of map update and map expansion in two cities (Pittsburgh and Miami). We do so by comparing expanded/updated maps with ground-truth maps using metrics. As shown above, map expansion to novel areas is harder than updating previously-seen areas.
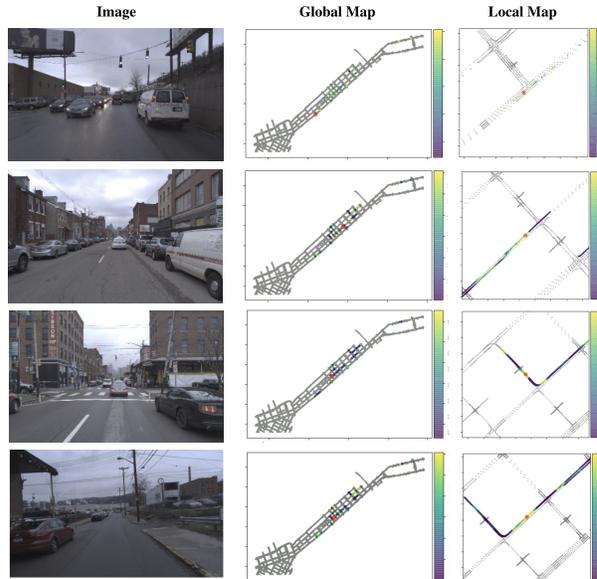


Figure 7. **Visual localization via Pix2Map.** We overlay retrieval scores on the corresponding local graphs from the original city map, generating a graph "heatmap" of possible locations given instantaneous ego-view images. We plot the ground-truth location as a red dot. In general, ground-truth locations tend to lie in high-scoring (yellow) regions. For example, the top ground truth corresponds to an intersection, while other high-scoring regions also tend to be graph intersections as well. Given a sequence of images, one may be able to reduce the ambiguity over time [8]

### A.2. Localization

Furthermore, our method demonstrates great visual localization ability based on visual and geometric understanding. We use the cosine similarities of retrieved graphs to generate a heatmap of possible ego-vehicle locations over a city-level map, showing the locations where their corresponding graphs are assigned a high likelihood, relative to the ground truth location shown as a red dot. While the ground truth is usually assigned a high likelihood, which indicates the promising performance of localization ability, the distribution becomes less sharp with respect to position when farther away from intersections. See Fig. 7 for more details.

### A.3. Map2Pix

We further show that it is also possible to retrieve ego-centric camera data using street maps. Such techniques could be used in the future to synthesize virtual worlds consistent with the query road geometry. We provide a few example image retrievals in Figure 8, visualizing the front
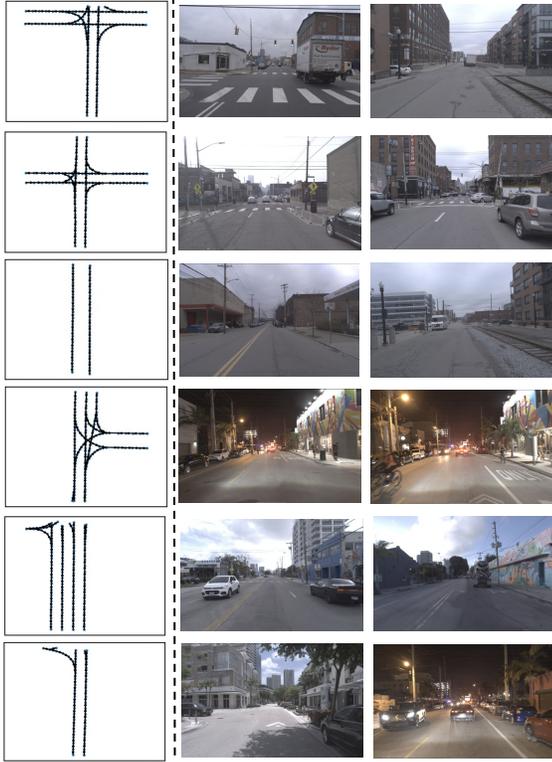
Figure 8. **Qualitative results for Map2Pix.** The goal is to retrieve ego-camera images given a street map. Such image retrieval may be useful for simulator-based training and validation of autonomous stacks. A single street geometry might retrieve multiple consistent, realistic imagery.

views of the top $K = 2$ images for each street map. As can be seen, the retrieved images correspond to rough geometric layouts encoded in the query graphs.