# DuMapNet: An End-to-End Vectorization System for City-Scale Lane-Level Map Generation

Deguo Xia[*]
Tsinghua University
Beijing, China
Baidu Inc.
Beijing, China
xiadeguo@baidu.com

Weiming Zhang[*]
Baidu Inc.
Beijing, China
zhangweiming@baidu.com

Xiyan Liu[*]
Baidu Inc.
Beijing, China
liuxiyan@baidu.com

Wei Zhang[*]
Baidu Inc.
Beijing, China
zhangwei99@baidu.com

Chenting Gong[*]
Baidu Inc.
Beijing, China
gongchenting@baidu.com

Jizhou Huang[†]
Baidu Inc.
Beijing, China
huangjizhou01@baidu.com

Mengmeng Yang[†]
Tsinghua University
Beijing, China
yangmm_qh@tsinghua.edu.cn

Diange Yang
Tsinghua University
Beijing, China
ydg@mail.tsinghua.edu.cn

## ABSTRACT

Generating city-scale lane-level maps faces significant challenges due to the intricate urban environments, such as blurred or absent lane markings. Additionally, a standard lane-level map requires a comprehensive organization of lane groupings, encompassing lane direction, style, boundary, and topology, yet has not been thoroughly examined in prior research. These obstacles result in labor-intensive human annotation and high maintenance costs. This paper overcomes these limitations and presents an industrial-grade solution named DuMapNet that outputs standardized, vectorized map elements and their topology in an end-to-end paradigm. To this end, we propose a group-wise lane prediction (GLP) system that outputs vectorized results of lane groups by meticulously tailoring a transformer-based network. Meanwhile, to enhance generalization in challenging scenarios, such as road wear and occlusions, as well as to improve global consistency, a contextual prompts encoder (CPE) module is proposed, which leverages the predicted results of spatial neighborhoods as contextual information. Extensive experiments conducted on large-scale real-world datasets demonstrate the superiority and effectiveness of DuMapNet. Additionally, DuMapNet has already been deployed in production at Baidu Maps since June 2023, supporting lane-level map generation tasks for over 360 cities while bringing a 95% reduction in costs. This demonstrates

that DuMapNet serves as a practical and cost-effective industrial solution for city-scale lane-level map generation.

## CCS CONCEPTS

• **Applied computing** → *Transportation*.

## KEYWORDS

Lane-Level Map Generation; End-to-End; Lane Group; Baidu Maps

## 1 INTRODUCTION

Lane-level map, as a crucial layer of the high-definition map, offers critical prior information for autonomous driving, facilitating beyond visual line of sight (BVLOS) perception, path planning, and decision-making with globally consistent road data. Specifically, lane-level map models the real world within decimeter level. This indicates that they not only depict the fundamental structure and layout of roads but also provide lane-level details, including lane line geometry, lane marking style, as well as connection topology, *etc.* Building on this, lane-level navigation, a pioneering advancement in precise travel guidance, has been extensively deployed to assist public travel by providing detailed, high-precision map elements and suggested routes (see Figure 1). Additionally, with its higher precision and advanced route planning capabilities, lane-level map will potentially benefit a wide range of geographical tasks, such as geo-object change detection [36], traffic condition prediction [37], estimated time of arrival prediction (ETA prediction, a.k.a., travel time estimation) [4, 5, 11], and road extraction [38] at Baidu Maps.

---

[*]These authors contributed equally to this work.
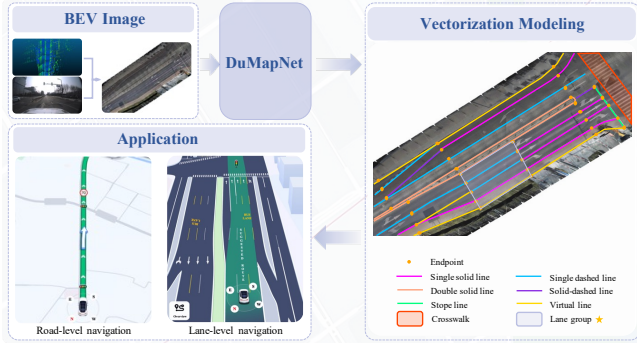[†]Corresponding authors.

**Figure 1: DuMapNet introduces a learning-based methodology for lane-level map vectorization. Our proposed method incorporates a scheme based on contextual prompts and components dedicated to group-wise lane prediction. With these advancements, DuMapNet achieves cost-effective generation of city-scale vectorized maps and significantly supports various applications, such as lane-level navigation in Baidu Maps.**

The task of lane-level map generation can be formulated as constructing and updating the core geographic elements that meet lane-level precision, given abundant road data. These geographic elements primarily include open-shape elements (*e.g.*, lane lines, stop lines, *etc.*) and closed-shape elements (*e.g.*, crosswalks). Moreover, as an efficient and standardized management unit in high-definition map, lane group plays a crucial role in onboard navigation and autonomous driving. It can be defined as a set of one or more lanes on a road segment perpendicular to the direction of travel [15, 29] (see Figure 1). Within a lane group, the number of lanes remains constant, and all lanes belong to the same road segment with the same direction of travel. Given this definition, the lane group emerges as an exceptionally convenient and efficient unit for driving guidance, markedly improving vehicle interaction with urban environments. Consequently, this paper introduces a learning-based solution that directly generates final standardized results in an end-to-end manner. This method effectively supersedes existing techniques that rely on manual post-processing to construct lane groups.

Traditional map generation solutions are often costly and labor-intensive as they require trained experts to manually annotate geographic elements. To improve efficiency with less human effort, leveraging advancements in computer vision for map generation has become a viable approach. These algorithms can be roughly categorized into segmentation-based methods [16, 26, 31], lane detection-based methods [6, 21], as well as vectorization-based methods [18, 23, 40]. Specifically, segmentation-based methods are suboptimal since they often require a series of post-processing strategies, such as thinning and fitting, to convert mask into vectorized map. Lane detection-based methods are usually limited in terms of extensibility and flexibility regarding map element types. While the vectorization-based solutions have achieved commendable results, they still exhibit limitations in prediction accuracy, post-processing logic and handling complex road scenarios, such as road wear and vehicle occlusion. Moreover, such onboard approaches are often constrained by computational power and local

construction patterns, preventing them from meeting the precision and global consistency required for city-scale lane-level map.

To fully explore the paradigm of large-scale lane-level map generation, we propose an automatic industrial-grade offboard solution termed DuMapNet. Given a bird's-eye-view (BEV) image, DuMapNet can unify the modeling of polyline-style and polygon-style map elements as a set of points. To significantly improve the prediction results for difficult scenarios such as road wear, occlusions, and complex intersections, as well as the connections of vectorization results among frames, we propose the **contextual prompts encoder (CPE)** module. By using the spatial prediction results of the current BEV image's neighborhood as prompt information, CPE significantly enhances the geometric and category consistency of the prediction results in a larger receptive field. To avoid the error accumulation effect and weak generalization of traditional multi-stage map-making methods, and considering the requirements for standardized map construction, we design a **group-wise lane prediction (GLP)** to output the vectorized results of lane groups through mutual constraints between lane group polygons and lane lines, without the need for complex post-processing logic. Finally, to achieve an end-to-end large-scale map generation mode, we develop the **topology prediction** module, which predicts the lane line topological relationships between BEV images, enabling large-scale map correlation. Our key contributions to both the research and industrial communities are as follows:
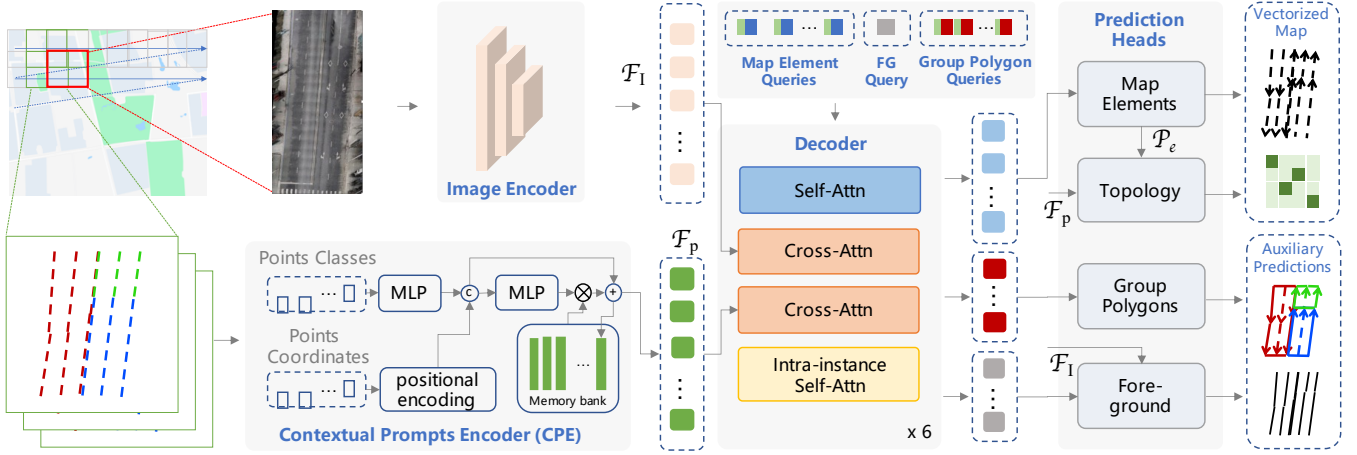
- **Potential impact:** We introduce DuMapNet, an end-to-end vectorization modeling framework, as an industrial-grade solution for city-scale lane-level map generation. DuMapNet has been successfully deployed in production at Baidu Maps, supporting lane-level map generation for over 360 cities and realizing a 95% reduction in costs.
- **Novelty:** DuMapNet represents a new paradigm for city-scale lane-level map generation task, achieving end-to-end predictions from bird's-eye-view (BEV) images to vectorized results that meet cartographic standards. The novelty lies in each stage, from the unified vectorization modeling, the group-wise lane prediction system, the contextual prompts encoder, to the topology prediction module, making the lane-level map generation task highly automatic and cost-effective.
- **Technical quality:** Extensive qualitative and quantitative experiments are performed on large-scale, real-world datasets collected from Baidu Maps, which demonstrate the superiority of DuMapNet. The successful deployment of DuMapNet at Baidu Maps further shows that it is a practical and robust solution for city-scale lane-level map generation.

## 2 DuMapNet

### 2.1 Preliminaries

The task of lane-level map generation from BEV image is defined as follows: given a BEV image $I$ collected from vehicle-mounted sensors as input, the network is supposed to predict the vectorized map elements $V$. Next, we will describe the data preparation and unified vectorization.

**Data Preparation.** Different from most onboard methods that operate on BEV features using multi-view images as input, following [1, 13], our offboard approach is built upon the BEV image

**Figure 2: Overall architecture of DuMapNet. DuMapNet processes the entire city-scale land area using a sliding window approach. For each local area, an image encoder is utilized to extract image features from the BEV image. Meanwhile, we propose a novel Contextual Prompts Encoder (CPE) to encode the predictions of adjacent scanned areas. To achieve Group-wise Lane Prediction (GLP), we meticulously tailor key network components, including the query, decoder, and prediction heads. Consequently, the network is capable of generating a vectorized map, which encompasses vectorized elements and their topology. Additionally, two auxiliary predictions are generated: the use of group polygons aids in the organization of lane groups, while foreground segmentation enhances lane point localization. For detailed illustrations, please refer to Section 2.**

produced using multi-view images, point cloud data, and vehicle poses information. The advantages primarily lie in two aspects: first, the regional global information can be fully utilized, such as geometric smoothness constraints, semantic correlations, and global precision consistency; second, conducting multi-trip data collection can alleviate the inevitable challenges such as precision bias and dynamic occlusion. Instead of appealing to heavy labeling manpower, we solve the large-scale annotation problem using the Baidu Map Database in an automatic fashion. Specifically, given a BEV image $I$ with $H \times W$ resolution, with the spatial resolution $4cm \times 4cm$ for each pixel, it covers $H/25$ meters by $W/25$ meters region with a certain geographic coordinate range. First, we index the instance geometries, labels, and lane group IDs within that range from the database. Second, based on the lane group IDs, the instance geometries are organized into a list format at the granularity of lane groups and mapped to pixel coordinate system. Simultaneously, we compute the minimum bounding rectangle of all instances contained within each lane group to create the group polygon. Finally, based on the spatial relationship between BEV images, neighborhood IDs are added to the ground truth. For better understanding, we have released a demo for reference on GitHub at https://github.com/XiyanLiu/DuMapNet.

**Unified Vectorization.** We define a unified vectorized representation for the core geographic elements in each local land area. Formally, given a BEV image $I$, we denote the corresponding lane-level vectorization as $V = \{G_i\}_{i=0}^{N_g}$, where $N_g$ denotes the number of the lane groups in the local land area. Each lane group $G_i$ is composed of a set of geographic elements $P_j$ and element style $C_j \in \mathbb{R}^{N_{cls}}$, where $C_j$ is a one-hot vector with $N_{cls}$ element styles in total. We thus denote a lane group as $G = \{P_j, C_j\}_{j=0}^{N_l}$, where $N_l$ is the element number in a lane group. Next, the point set of each element instance is denoted as $P = \{p_k\}_{k=0}^{N_p}$, where $N_p$ is the

number of the points and $p_k \in \mathbb{R}^2$ denotes the coordinates of each point.

## 2.2 Overall Architecture

The city-scale lane-level map generation is inherently complex, requiring a comprehensive organization of lane groupings and existing methods only generate partial elements of lane groupings. Meanwhile, DuMapNet is the first end-to-end solution to achieve city-scale lane-level map generation, realizing practical and effective industrial gains.

Specifically, Figure 2 illustrates the overall architecture of our proposed DuMapNet. To obtain a city-scale lane-level vectorized map, DuMapNet processes the entire land area using a sliding window approach following a zig-zag scan sequence. The model's inputs contain two parts: a BEV image and contextual prompts. Specifically, the BEV image is obtained by the aforementioned data preparation process, providing abundant appearance features of the local land area. We employ an image encoder that contains a backbone network and a Feature Pyramid Network (FPN) to extract BEV feature $F_I \in \mathbb{R}^{384 \times 384}$ from the BEV image $I$. Meanwhile, to enable spatial coherent lane-level vectorized predictions over adjacent land areas, we propose to take the predicted vectorized map of adjacent scanned areas as the additional input of DuMapNet. We further tailor a contextual prompts encoder (CPE) to realize an effective encoding for the predictions of adjacent scanned areas. Moreover, we devise a query combination that contains a set of hierarchical queries for lane line prediction, one foreground segmentation query, and a set of queries for lane group polygon prediction. In the Decoder, the proposed query combination interacts with both the BEV feature and the contextual prompt embeddings derived from CPE. Finally, we construct multiple tasking prediction heads to facilitate various predictions, where the predictions include vectorized

results for lane lines and the topology that indicates the connection of lane lines between different lane areas. Furthermore, we propose two auxiliary predictions, where the group polygons are adopted to facilitate lane group organization while the foreground segmentation is to help improve lane point localization.

## 2.3 Contextual Prompts Encoder (CPE)

Inspired by the recent success of prompt-based vision models [12, 14, 20], our proposed CPE adopts a simple yet effective architecture to encode both the geometric and semantic information of the vectorization results of adjacent land areas, providing contextual cues for the vectorization of the current land area during the sliding window operation. Formally, we define

$$F_p = \mathbb{CPE}(\{V_s\}_{s \in A}) \tag{1}$$

where $A$ denotes a set of adjacent land areas. $F_p \in R^{M_g \times M_l \times N_p \times 256}$ denotes the prompt embeddings that interact with the features of intermediate layers in the decoder. Where $M_g, M_l, N_p$ denote the number of total predicted groups, the number of element instances in each group, and the number of points in each instance, respectively.

The architecture of CPE is illustrated in Figure 2. Specifically, we adopt a shared MLP (Multi-Layer Percseptron) network to encode the predicted style type for each element to serve as the semantic encoding. For the geometric information, we perform a shared positional encoding sub-network to the coordinates of element points. This sub-network consists of sine and cosine functions with different frequencies as well as a subsequent MLP. Finally, the geometric and semantic cues are concatenated and then fed to an MLP, generating the final prompt embeddings $F_p$. Furthermore, we introduce a memory mechanism into CPE to realize a long-term feature dependence. Specifically, we adopt FIFO queue as a memory bank to store the $F_p$ prompt embeddings of previous $T$ frames of local land areas. The memory bank efficiently stores prompt embeddings of the remaining $T - 1$ neighboring frames. These embeddings are then aggregated using a weighted sum operation within the CPE. The weights assigned to each stored frame are learnable parameters generated by intermediate layer in CPE, allowing the model to adaptively focus on relevant information. This design effectively reduces noise in the prompted information (*e.g.* prompts may contain prediction errors) while maintaining a lightweight architecture with minimal computational overhead. The memory bank is the core module of the CPE and directly reflects its performance. Finally, $F_p$ is obtained via learning an aggregation of the stored embeddings.

## 2.4 Group-wise Lane Prediction (GLP)

A lane group refers to a collection of lanes that share common characteristics *e.g.* same style or are directed towards a common goal or destination. Practically, lane groups are essential for path planning and navigation as they help in understanding complex road structures. However, predicting an accurate group-wise lane is challenging since it demands sophisticated semantic analysis and geometrical reasoning. Particularly, localizing the endpoint of a lane instance requires the knowledge of the style and topology changing of the other lanes in the same group. To address this, we propose to use a polygon, namely group polygon to outline

the boundary of a lane group. We further introduce an auxiliary task in the network architecture to predict the group polygons. Since all the points of the predicted map elements are located in the group polygons, we propose an additional point-in-polygon loss to facilitate group-wise lane prediction.

In this section, we introduce the key components regarding to group-wise lane prediction (GLP), including the query design, decoder architecture, and the prediction heads.

**Queries.** We design a query combination to flexibly encode structured map information and perform hierarchical bipartite matching for both map element and group polygon learning. Specifically, we extend the hierarchical query scheme in MapTR [18] and customize two set of queries, namely, element queries $\{q_i^l\}_{i=0}^{N_l}$ and group polygon queries $\{q_i^g\}_{i=0}^{N_g}$. These two types of queries adopt the same hierarchical query scheme that efficiently encodes instance-level and point-level information. Moreover, we introduce an additional foreground-background (FG) query $q^s$ for the auxiliary task of semantic segmentation.
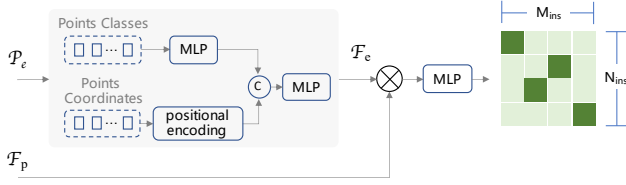
**Decoder.** All map elements, group polygons, and segmentation masks are simultaneously predicted using a unified Transformer structure. The decoder is composed of several cascaded layers, each incorporating a self-attention module, two cross-attention modules, and an intra-instance self-attention module. The initial self-attention module is designed to enable hierarchical queries to exchange information across the entire feature space. The subsequent cross-attention module facilitates interaction between hierarchical queries and BEV features. To enhance prediction accuracy and spatial consistency, an additional cross-attention module has been innovatively introduced, with contextual prompt embeddings serving as input keys and values to interact with hierarchical queries. Lastly, the intra-instance self-attention module allows for interactions between points within the same instance, thereby improving geometric smoothness. Ultimately, after processing through the decoder, the hierarchical queries are effectively encoded into group-level query embeddings $F_g \in \mathbb{R}^{N_g \times N_p \times 256}$, line-level query embeddings $F_l \in \mathbb{R}^{N_l \times N_p \times 256}$, and a foreground embedding $F_s \in \mathbb{R}^{1 \times N_p \times 256}$.

**Predictions.** For predicting lane lines and lane group polygons, we input both map element and group polygon query embeddings into a shared classification branch and a shared regression branch to facilitate type classification and geometric property regression, respectively. For each predicted instance, the regression branch outputs a vector of dimension $\mathbb{R}^{N_p \times 2}$, representing the normalized coordinates of $N_p$ points.

Furthermore, to enhance the performance of the classification and regression branches, thereby improving prediction accuracy and accelerating training convergence, we propose a foreground segmentation branch. Instead of directly utilizing the BEV features for segmentation, an individual foreground query $q^{seg} \in \mathbb{R}^{1 \times 256}$ is introduced alongside the hierarchical queries. Following processing by the conventional decoder network and MLP encoding, the foreground query embedding interacts with the BEV feature to generate a segmentation map.

## 2.5 Topology Prediction

Considering our objective to generate city-scale lane-level map in an end-to-end manner, it is insufficient to solely predict the

**Figure 3: Topology prediction. The topology matrix is produced as an additional output of the decoder to indicate the connections between $N_{ins}$ element instances in the current land area and $M_{ins}$ element instances in the contextual land areas.**

vectorization of single-frame BEV images. Predicting the topological relationships between frames is indispensable for our task. To address this, we propose to directly predict a topology matrix $\mathcal{M} \in \mathbb{R}^{M_{ins} \times N_{ins}}$ which indicates the connections between $N_{ins}$ element instances in current land area and $M_{ins}$ element instances in the contextual land areas.

Inspired by [35], we formulate topology prediction as a classification task, where the topology matrix is produced as an additional output of the decoder. The architecture is outlined in Figure 2 and specified in Figure 3. Specifically, we adopt a similar sub-network as the CPE to encoder the predicted map elements, producing an embedding $F_e \in \mathbb{R}^{M_{ins} \times N_p \times 256}$ that encapsulate both the predicted coordinates and class information of instances. Subsequently, we calculate the correlation of $F_e$ and the prompt embeddings $F_p$ than adopt a MLP to produce the topology matrix. During inference, lane lines and lane group polygons are aggregated based on their geometric relationships, facilitating the generation of coherent lane group configurations.

## 2.6 End-to-End Training

In the training stage, for each frame, we apply the hierarchical matching scheme in MapTR [18] to obtain pairs of the map-elements predictions and the ground truths, denoted as $\{\hat{P}_i^l, P_i^l\}$. Meanwhile, we adopt the same matching scheme to obtain pairs of the predicted group polygons and the ground truths, denoted as $\{\hat{P}_j^g, P_j^g\}$.

Based on the matching results and the correspondence between map elements and their group polygons, we employ multiple task-specific losses to train our proposed DuMapNet in an end-to-end manner:

$$\mathcal{L} = \alpha \mathcal{L}_l + \beta \mathcal{L}_t + \lambda \mathcal{L}_g + \eta \mathcal{L}_{gl} + \mu \mathcal{L}_s \tag{2}$$

where $\mathcal{L}_l$, $\mathcal{L}_t$ denotes the loss for learning map element and topology, respectively. Besides, $\mathcal{L}_g$ denotes the loss for learning group polygons while $\mathcal{L}_{gl}$ is an additional point-in-polygon loss to facilitate group-wise elements organization. Finally, we introduce a foreground segmentation loss $\mathcal{L}_s$ to enhance lane point localization. $\alpha, \beta, \lambda, \eta, \mu$ are hyperparameters that strike a balance between difference losses. Next, we provide detailed illustrations for each type of loss function.

**Map Elements Learning.** For each map element instance, we employ aligned Focal Loss $\mathcal{L}_{cls}$ for style classification and an L1 regression loss $\mathcal{L}_{L1}$ for point localization, respectively. Here we follow MapTR [18] and further apply a direction loss $\mathcal{L}_{dir}$ to align

the direction of the predicted lane segments with the ground truth. The loss for training the map elements is thus denoted as:

$$\mathcal{L}_l = \mathcal{L}_{cls} + \mathcal{L}_{reg} + \mathcal{L}_{dir} \tag{3}$$

Furthermore, inspired by Stable DINO [22], we adapt the aligned Focal Loss to enhance the alignment between classification score and localization quality in $\mathcal{L}_{cls}$. As shown in Equation. (4), L1 distance of the $i$-th matched pairs between prediction $\hat{P}_i^l$ and corresponding ground truth $P_i^l$ is used as a positional metric to supervise the training probabilities of positive examples. The classification loss is thus formulated as:

$$\mathcal{L}_{cls} = \sum_{i=1}^{N_{pos}} (|d_i - s_i|^\gamma) BCE(s_i, d_i) + \sum_{i=1}^{N_{neg}} s_i^\gamma BCE(s_i, 0) \tag{4}$$

$$d_i = ||\hat{P}_i^l - P_i^l||_1 \tag{5}$$

where $s_i$ is the probability for the $i$-th predicted map element. $N_{pos}$ and $N_{neg}$ denote the number of positive and negative elements, respectively. Moreover, as for $\mathcal{L}_t$, we apply the same loss function as $\mathcal{L}_l$ for learning group polygons.

**Topology Learning.** We define the topological relationship as a 2-class classification task *i.e.* connected or not. As the number of connected ones is significantly less than the number of unconnected ones, we apply a focal loss to supervise the topological association matrix prediction.

**Group-guided Auxiliary Supervision.** We leverage the group polygons to provide auxiliary supervision for learning high-quality group-wise lane lines. Specifically, this auxiliary supervision is designed based on the following observations: all the points of the lane lines should be located inside or on the boundary of their corresponding group polygons. We thus propose a point-in-polygon loss as:

$$\mathcal{L}_{gl}(\hat{P}_i^l, P_j^g) = \sum_{\substack{p \in \hat{P}_i^l \ and \\ outside \ P_j^g}} D(p, P_j^g) \tag{6}$$

where $D(p, P_j^g)$ is the closest distance from point $p$ to any edge of the lane group polygon $P_j^g$. During the training stage, we adopt the ground-truth group polygons to punish predicted points located outside their group polygon. During the inference stage, we simply employ the predicted group polygons to check map elements' locations.

**Segmentation-guided Auxiliary Supervision.** To improve the accuracy of vectorized predictions, we choose a combination of Binary Cross Entropy and Dice loss to calculate the loss between the predicted foreground mask $F_s \in R^{H \times W}$ and the ground truth mask $F_m \in R^{H \times W}$ generated from ground truth lane lines:

$$\mathcal{L}_s = BCE(F_s, F_m) + \mathcal{L}_{dice}(F_s, F_m) \tag{7}$$

# 3 EXPERIMENTS

## 3.1 Experimental Settings

**Datasets.** As aforementioned, DuMapNet has been deployed at Baidu Maps, supporting over 360 cities. To evaluate the effectiveness of DuMapNet, we have collected a large-scale real-world dataset,

DuLD, consisting of bird's-eye view (BEV) images and ground truth data from six cities: Beijing, Guangzhou, Changchun, Changzhou, Chongqing, and Leshan. These cities were chosen for their varied urban scales and geographic features. The dataset from Beijing, Guangzhou, Changchun, and Changzhou was divided into a training set and a validation set in a 9 : 1 ratio. Meanwhile, data from Chongqing and Leshan were used as the test set to evaluate the model's performance. Statistically, DuLD contains 134,524 images, spanning 8,072 kilometers, with each image at a resolution of 1536 × 1536 pixels. More details can be found in Table 1. Importantly, to investigate the benefits of larger-scale data, we introduce DuLD-L, a dataset with one million paired images and corresponding ground truths, and evaluate DuMapNet's performance on this expanded dataset.

**Evaluation Metrics.** We adopt recall (R) and precision (P) to assess the quality of map construction at the instance level. The evaluation considers category consistency, endpoint distance, and overlap to determine if a pair of lane instances from ground truth and prediction match. Moreover, category consistency and IoU ($IoU > 0.5$) are used for closed-shape elements such as crosswalks. Specifically, category consistency mandates that the instances belong to the same category. Endpoint distance requires the L2 distance between the start points and end points of the instances should be less than 3 meters, respectively. Overlap considers the parallel distance between instances. When calculating overlap, both prediction instances and ground truth are first divided into multiple segments with 1-meter intervals. Then, the projection distance between segments is calculated. If the proportion of segments with projection distance smaller than the threshold $d = \{0.5, 1\}m$ exceeds the threshold $r = \{0.5, 0.8\}$, the prediction instance will be considered as true positive (TP). In the following experiments, we use $R@P_{d,r} = p$ to represent the recall at $p$ precision with threshold $d$ and $r$. Lower values of threshold $d$ and higher values of threshold $r$ signify more stringent precision requirements.

**Implementation Details.** Our model is trained using 16 NVIDIA Tesla V100 GPUs, with a batch size of 16. The AdamW [24] optimizer is employed with a weight decay of 0.01 and the initial learning rate is set to $6 \times 10^{-4}$ with cosine decay. The input images have a resolution of 768 × 768 pixels. For our architecture, we employ ResNet50 [8] and HRNet48 [34] as backbones. The default number of instance queries, point queries and decoder layers are 50, 50 and 6, respectively. As for hyper-parameters of loss weight, we set $\alpha$, $\beta$, $\lambda$, $\eta$ and $\mu$ to 1, 1, 0.2, 0.15 and 100, respectively. The inference time is measured on a single NVIDIA Tesla V100 GPU with batch size 1.

## 3.2 Evaluation

**Comparison with Baselines.** We compare our DuMapNet with segmentation-based method [32] and other vectorization-based methods [18, 40]. As shown in Table 2, vectorization-based methods achieve better results without complex post-processing logic. In particular, our DuMapNet outperforms the existing state-of-the-art method by a large margin (+2.66% when $P_{1,0.8} = 90\%$) under the same setting of ResNet50, indicating the effectiveness of our method. Surprisingly, our method achieves a further improvement of 4.99% by replacing the backbone with HRNet48 to obtain enhanced feature representation. Furthermore, our method achieves 73.28% recall on

**Table 1: Statistics of DuLD dataset.**

| Dataset | Mileage (km) | Image | City | |
|---|---|---|---|---|
| Train | 7,922 | 118,822 | Beijing | Guangzhou |
| Val | | 13,202 | Changchun | Changzhou |
| Test | 150 | 2,500 | Chongqing | Leshan |
| All | 8,072 | 134,524 | - | |

one million of training data, demonstrating that as the training data volume increases, our method can achieve greater advantages.

The quantitative comparisons of different evaluation thresholds are summarized in Table 3. From the results, we observe that DuMapNet consistently brings significant improvements. Taking $R@P_{0.5,0.8} = 90\%$ as an example, DuMapNet achieves better performance with +3.00% ~ 5.04% recall gains on DuLD, which underscores our method's ability to achieve superior geometric precision and maintain category consistency. Particularly, as the projection distance decreases, the performance of all approaches experiences a significant decline, yet our DuMapNet shows a slighter drop, indicating that our method is more robust and maintains superior performance at higher precision levels.

**Ablation Studies.** In this section, extensive ablation studies are conducted to systematically evaluate the key designs of our DuMapNet. As shown in Table 4, Group I is the baseline without a series of designs. From Group I and II, it is proved that adding the task-aligned supervision can bring slight improvements by fostering better synergy between geometric learning and category identification. Further analysis between Groups II and III reveals that incorporating intra-instance self-attention results in a 0.99 improvement under $P_{1,0.8} = 95\%$. The transition from Group III to Group IV examines the impact of adding segmentation-guided auxiliary supervision, showing marked improvements in all metrics, especially at higher precision levels. This outcome is expected as the segmentation branch contributes to more fine-grained pixel-level modeling, enhancing semantic understanding and refining prediction accuracy. The results from Group V highlight the significant role of the proposed contextual prompts encoder (CPE) module, showing a notable increase of 2.74% in recall at a high accuracy level ($P_{1,0.8} = 95\%$). These findings demonstrate that CPE, by leveraging spatial prediction results from the area surrounding the current BEV image, significantly enhances the geometric and category consistency of predictions across a broader receptive field.

**Analysis of generalization.** To further demonstrate the generalization of our method, five cities are additionally selected as the test set. These cities are distributed across various regions, such as Harbin in northeastern China and Xi'an in northwestern China, and exhibit diverse sizes, with Shanghai being a large first-tier city, while Zhongshan is a second-tier city. Finally, a total of 5,000 images were collected for evaluation. The experimental results are presented in Table 5. On the one hand, DuMapNet outperforms the existing state-of-the-art methods on all city test sets, demonstrating the effectiveness of our approach. On the other hand, DuMapNet shows superior generalization with less fluctuation in

**Table 2: Comparisons with state-of-the-art methods on DuLD test set at different precision levels. R50 and HR48 correspond to ResNet50 [8] and HRNet48 [34], respectively. FPSs are measured on the same machine with NVIDIA Tesla V100.**

| Method | Backbone | Training Set | $R@P_{1,0.8} = 80\%$ | $R@P_{1,0.8} = 90\%$ | $R@P_{1,0.8} = 95\%$ | FPS |
|---|---|---|---|---|---|---|
| HMSA [32] + post-processing | HR48 | DuLD | 58.48 | - | - | - |
| MapTR [18] | R50 | DuLD | 69.56 | 58.58 | 39.49 | 29.7 |
| GeMap [40] | R50 | DuLD | 71.96 | 60.45 | 39.32 | 29.3 |
| DuMapNet | R50 | DuLD | 74.61 | 63.11 | 43.27 | 27.9 |
| DuMapNet | HR48 | DuLD | 77.34 | 68.10 | 54.21 | 26.6 |
| DuMapNet | HR48 | DuLD-L | **83.40** | **73.28** | **61.24** | 26.6 |

**Table 3: Comparisons with state-of-the-art methods on DuLD test set under different thresholds $d$ and $r$, where $d$ represents the projection distance threshold and $r$ represents the proportion of segments threshold.**

| Method | Backbone | Training Set | $R@P_{1,0.8} = 90\%$ | $R@P_{1,0.5} = 90\%$ | $R@P_{0.5,0.8} = 90\%$ | $R@P_{0.5,0.5} = 90\%$ |
|---|---|---|---|---|---|---|
| MapTR [18] | R50 | DuLD | 58.58 | 59.29 | 54.29 | 56.5 |
| GeMap [40] | R50 | DuLD | 60.45 | 61.0 | 56.33 | 58.15 |
| DuMapNet | R50 | DuLD | 63.11 | 63.56 | 59.33 | 60.87 |
| DuMapNet | HR48 | DuLD | 68.10 | 68.91 | 65.47 | 66.91 |
| DuMapNet | HR48 | DuLD-L | **73.28** | **73.98** | **71.66** | **72.41** |

**Table 4: Effects of core components of DuMapNet.**

| Group | Task-aligned Supervision | Intra-instance Self-attention | Segmentation-guided Auxiliary Supervision | CPE | $R@P_{1,0.8} = 80\%$ | $R@P_{1,0.8} = 90\%$ | $R@P_{1,0.8} = 95\%$ |
|---|---|---|---|---|---|---|---|
| I | | | | | 73.93 | 63.11 | 49.42 |
| II | ✓ | | | | 74.04 | 63.61 | 49.72 |
| III | ✓ | ✓ | | | 76.95 | 64.60 | 49.87 |
| IV | ✓ | ✓ | ✓ | | 77.18 | 65.81 | 51.47 |
| V | ✓ | ✓ | ✓ | ✓ | 77.34 | 68.10 | 54.21 |

**Table 5: Comparisons with state-of-the-art methods on the additional test sets are conducted to assess the generalization performance. The additional test sets consist of data from five cities, including Harbin, Hangzhou, Xi'an, Zhongshan, and Shanghai. The metric is $R@P_{1,0.8} = 80\%$.**
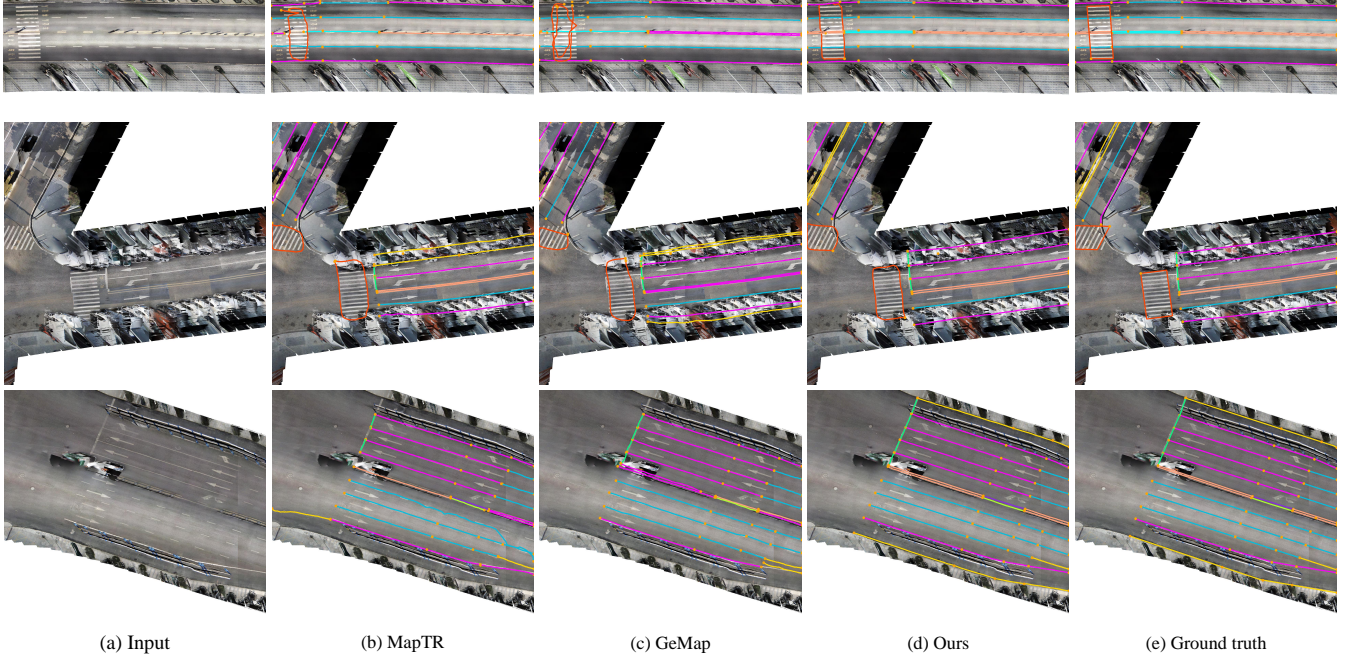
| Method | Backbone | Training Set | Harbin | Hangzhou | Xi'an | Zhongshan | Shanghai |
|---|---|---|---|---|---|---|---|
| MapTR[18] | R50 | DuLD | 73.30 | 66.30 | 71.76 | 73.55 | 69.66 |
| GeMap[40] | R50 | DuLD | 75.80 | 69.07 | 72.03 | 75.19 | 71.81 |
| DuMapNet | R50 | DuLD | 76.12 | 73.35 | 76.11 | 77.30 | 73.31 |
| DuMapNet | HR48 | DuLD | 78.73 | 76.19 | 77.50 | 77.33 | 76.66 |
| DuMapNet | HR48 | DuLD-L | **83.99** | **82.10** | **82.06** | **82.55** | **84.38** |

performance across the five cities. For example, the maximum deviation of DuMapNet-R50 across the five cities is 3.99%, while the second-best model, *i.e.*, GeMap has a maximum deviation of 6.73%.
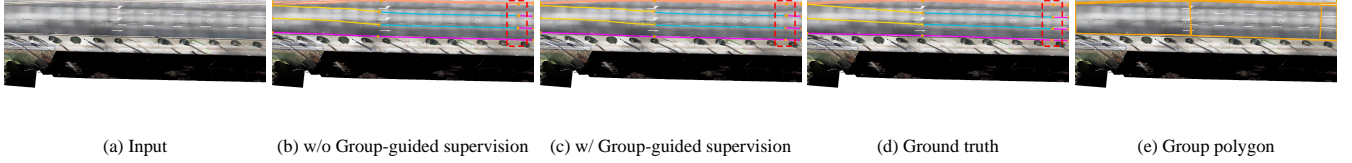
## 3.3 Visualization

Qualitative results from the DuLD dataset are presented in Figure 4 and Figure 5. DuMapNet not only performs well in simple scenes but also predicts high-quality vectorized map elements in complex scena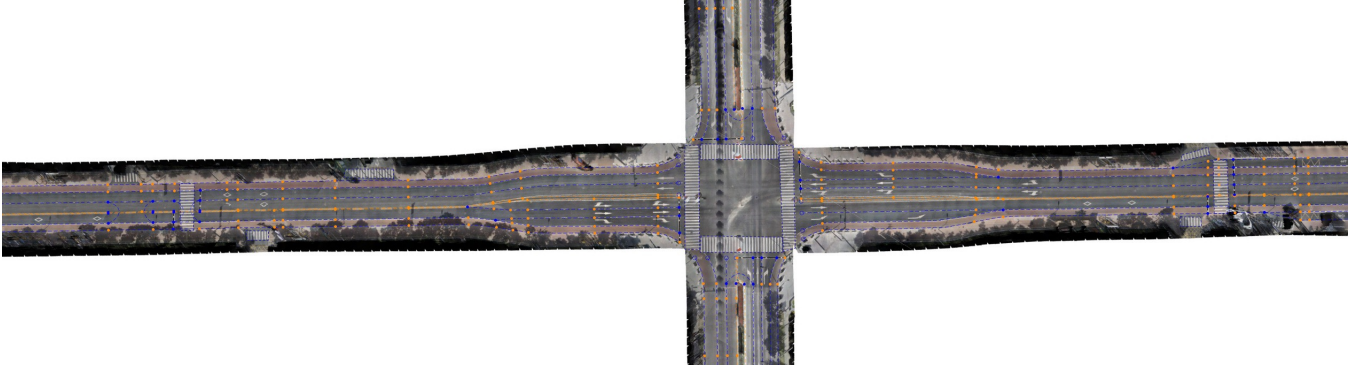rios like intersections, road wear and occlusions. As shown in Figure 4, DuMapNet exhibits significant advantages in terms of lane recall, lane accuracy, and endpoint accuracy. For example, as demonstrated in the second row of Figure 4, DuMapNet precisely captures both the geometry and category of lane lines in occlusion scenes, avoiding unnecessary lane line predictions. In addition, Figure 5 provides a visual comparison that underscores the effectiveness of group-guided supervision in enabling accurate prediction of endpoint positions, even in scenarios with subtle visual differences, such as the dashed line and the segment of solid line. Group-guided

(a) Input　　　(b) MapTR　　　(c) GeMap　　　(d) Ours　　　(e) Ground truth

**Figure 4: Comparisons of our method with state-of-the-art models in lane-level map generation.**



(a) Input　　(b) w/o Group-guided supervision　　(c) w/ Group-guided supervision　　(d) Ground truth　　(e) Group polygon

**Figure 5: Qualitative visualization of the proposed group-guided supervision.**



**Figure 6: Qualitative visualization of urban scene.**

supervision also ensures that the endpoints of lane lines within a lane group are correctly aligned. Furthermore, as illustrated in Figure 6, the implementation of the topology prediction module enables the end-to-end generation of comprehensive urban scenes at a global level.

## 4 RELATED WORK

Here we briefly review the closely related work in the fields of map construction and lane detection.

## 4.1 Map Construction

With the development of deep learning and BEV perception [10], map construction is transitioning from a labor-intensive annotation task to a model-based dense prediction challenge. Segmentation-based methods [7, 9, 17, 25, 26] generate rasterized map by performing BEV semantic segmentation. To build vectorized maps, HDMapNet [16] adopts a two-stage approach of segmentation followed by post-processing to generate vectorized instances. As the first end-to-end framework, VectorMapNet [23] utilizes an auto-regressive decoder to predict points sequentially. MapTR [18] proposes a unified shape modeling method based on a parallel end-to-end framework, which has been followed by many works [3, 27, 28]. MapVR [39] applies differentiable rasterization to vectorized outputs to performs precise and geometry-aware supervision. MapTRv2 [19] further introduces auxiliary one-to-many matching and auxiliary dense supervision to speedup convergence. BeMapNet [27] adopts a parameterized paradigm and constructs map elements as piecewise Bezier curves. PivotNet [3] utilizes a dynamic number of pivotal points to model map elements, preventing the loss of essential details. Different from the existing works, our proposed DuMapNet leverages neighboring map elements as prompts to guide the generation of map elements in the current frame, which can enhance the spatial consistency of the map elements.

## 4.2 Lane Detection

Lane detection plays a critical role in detecting lane elements in road scenes and can be considered as a subtask of map construction. LaneATT [30] utilizes an anchor-based deep lane detection model. CondLaneNet [21] adopts a conditional lane detection model based on conditional convolution and row-wise formulation. GANet [33] formulates lane detection as a keypoint estimation and association problem. BezierLaneNet [6] proposes a parametric Bezier curve-based method, which can model the geometric shapes of lane lines. PersFormer [2] utilizes a transformer-based spatial feature transformation module and unify 2D and 3D lane detection simultaneously. Different from these methods that primarily focus on lane elements, our proposed DuMapNet models map elements in a unified vectorized form, which can detect open-shape map elements such as lanes, as well as closed-shape elements like crosswalks.

## 5 DISCUSSION

Before the deployment of DuMapNet, Baidu Maps relied heavily on labor-intensive manual annotation processes that involved segmentation techniques and complex post-processing logic. This approach significantly increased operational costs and decreased efficiency. With the introduction of DuMapNet, now operational in over 360 cities, production efficiency has seen a twenty-fold improvement, leading to a remarkable 95% reduction in costs.

Despite DuMapNet's impressive achievements, several challenging issues remain unresolved and require further investigation. For instance, the model struggles in scenarios with extensive static obstructions, such as long stretches of road with parked vehicles. Such conditions disrupt performance because the lack of visible road surface markings compromises the effectiveness of the contextual prompts encoder. To address this challenge, integrating

multi-source data may be an effective approach. In addition, generating qualified map data from low-precision sources, such as crowdsourced data, presents an intriguing challenge that merits deeper exploration in future work. Currently, leveraging the high timeliness, broad coverage, and low cost of crowdsourced data for map updates represents a more reasonable paradigm. For example, crowdsourced data can provide timely updates for elements with lower accuracy requirements, such as style changes, or for dynamic changes like construction or temporary road closures.

## 6 CONCLUSIONS

In this paper, we present an effective industrial solution for city-scale lane-level map generation. Specifically, we reformulate this task as a vectorization modeling task that takes bird's-eye-view (BEV) images as input and outputs standardized, vectorized map elements and their topology in an end-to-end paradigm. We pioneer organize the lane group using a learning-based methodology and address it through the proposed group-wise lane prediction (GLP) system that outputs vectorized results of lane groups by applying mutual constraints between lane group polygons and lane lines, thereby eliminating the need for intricate post-processing logic. To improve the generalization in challenging scenarios, such as road wear and occlusions, as well as to improve the continuity of vectorization results across frames, we present the contextual prompts encoder (CPE) module, which leverages the spatial prediction results from the surrounding area of the current BEV image as contextual information. Extensive experiments conducted on the collected large-scale real-world dataset from Baidu Maps demonstrate the superiority of DuMapNet. The successful deployment of DuMapNet at Baidu Maps has significantly improved its performance. Since its launch in June 2023, DuMapNet served over 360 cities while bringing a 95% reduction in costs.

## 7 ACKNOWLEDGMENTS

# REFERENCES

[1] Syed Ammar Abbas and Andrew Zisserman. 2019. A geometric approach to obtain a bird's eye view from an image. In *Proceedings of the IEEE/CVF international conference on computer vision workshops*. 0–0.

[2] Li Chen, Chonghao Sima, Yang Li, Zehan Zheng, Jiajie Xu, Xiangwei Geng, Hongyang Li, Conghui He, Jianping Shi, Yu Qiao, et al. 2022. Persformer: 3d lane detection via perspective transformer and the openlane benchmark. In *European Conference on Computer Vision*. Springer, 550–567.

[3] Wenjie Ding, Limeng Qiao, Xi Qiu, and Chi Zhang. 2023. Pivotnet: Vectorized pivot learning for end-to-end hd map construction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 3672–3682.

[4] Xiaomin Fang, Jizhou Huang, Fan Wang, Lihang Liu, Yibo Sun, and Haifeng Wang. 2021. SSML: Self-Supervised Meta-Learner for En Route Travel Time Estimation at Baidu Maps. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 2840–2848.

[5] Xiaomin Fang, Jizhou Huang, Fan Wang, Lingke Zeng, Haijin Liang, and Haifeng Wang. 2020. ConSTGAT: Contextual Spatial-Temporal Graph Attention Network for Travel Time Estimation at Baidu Maps. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2697–2705.

[6] Zhengyang Feng, Shaohua Guo, Xin Tan, Ke Xu, Min Wang, and Lizhuang Ma. 2022. Rethinking efficient lane detection via curve modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 17062–17070.

[7] Adam W Harley, Zhaoyuan Fang, Jie Li, Rares Ambrus, and Katerina Fragkiadaki. 2022. A simple baseline for bev perception without lidar. *arXiv e-prints* (2022), arXiv–2206.

[8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.

[9] Anthony Hu, Zak Murez, Nikhil Mohan, Sofía Dudas, Jeffrey Hawke, Vijay Badrinarayanan, Roberto Cipolla, and Alex Kendall. 2021. Fiery: Future instance prediction in bird's-eye view from surround monocular cameras. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 15273–15282.

[10] Junjie Huang, Guan Huang, Zheng Zhu, Yun Ye, and Dalong Du. 2021. Bevdet: High-performance multi-camera 3d object detection in bird-eye-view. *arXiv preprint arXiv:2112.11790* (2021).

[11] Jizhou Huang, Zhengjie Huang, Xiaomin Fang, Shikun Feng, Chen Xuyi, Liu Jiaxiang, Yuan Haitao, and Haifeng Wang. 2022. DuETA: Traffic Congestion Propagation Pattern Modeling via Efficient Graph Learning for ETA Prediction at Baidu Maps. In *Proceedings of the 31st ACM International Conference on Information and Knowledge Management*.

[12] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. 2022. Visual prompt tuning. In *European Conference on Computer Vision*. Springer, 709–727.

[13] Giseop Kim, Sunwook Choi, and Ayoung Kim. 2021. Scan context++: Structural place recognition robust to rotation and lateral variations in urban environments. *IEEE Transactions on Robotics* 38, 3 (2021), 1856–1874.

[14] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. 2023. Segment anything. *arXiv preprint arXiv:2304.02643* (2023).

[15] Nikol Krausz, Vivien Potó, János Máté Lógó, and Árpád Barsi. 2022. Comparison of complex traffic junction descriptions in automotive standard formats. *Periodica Polytechnica Civil Engineering* 66, 1 (2022), 282–290.

[16] Qi Li, Yue Wang, Yilun Wang, and Hang Zhao. 2022. Hdmapnet: An online hd map construction and evaluation framework. In *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 4628–4634.

[17] Zhiqi Li, Wenhai Wang, Hongyang Li, Enze Xie, Chonghao Sima, Tong Lu, Yu Qiao, and Jifeng Dai. 2022. Bevformer: Learning bird's-eye-view representation from multi-camera images via spatiotemporal transformers. In *European conference on computer vision*. Springer, 1–18.

[18] Bencheng Liao, Shaoyu Chen, Xinggang Wang, Tianheng Cheng, Qian Zhang, Wenyu Liu, and Chang Huang. 2022. Maptr: Structured modeling and learning for online vectorized hd map construction. *arXiv preprint arXiv:2208.14437* (2022).

[19] Bencheng Liao, Shaoyu Chen, Yunchi Zhang, Bo Jiang, Qian Zhang, Wenyu Liu, Chang Huang, and Xinggang Wang. 2023. Maptrv2: An end-to-end framework for online vectorized hd map construction. *arXiv preprint arXiv:2308.05736* (2023).

[20] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. *arXiv preprint arXiv:2304.08485* (2023).

[21] Lizhe Liu, Xiaohao Chen, Siyu Zhu, and Ping Tan. 2021. Condlanenet: a top-to-down lane detection framework based on conditional convolution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 3773–3782.

[22] Shilong Liu, Tianhe Ren, Jiayu Chen, Zhaoyang Zeng, Hao Zhang, Feng Li, Hongyang Li, Jun Huang, Hang Su, Jun Zhu, et al. 2023. Detection Transformer with Stable Matching. *arXiv preprint arXiv:2304.04742* (2023).

[23] Yicheng Liu, Tianyuan Yuan, Yue Wang, Yilun Wang, and Hang Zhao. 2023. Vectormapnet: End-to-end vectorized hd map learning. In *International Conference on Machine Learning*. PMLR, 22352–22369.

[24] Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101* (2017).

[25] Bowen Pan, Jiankai Sun, Ho Yin Tiga Leung, Alex Andonian, and Bolei Zhou. 2020. Cross-view semantic segmentation for sensing surroundings. *IEEE Robotics and Automation Letters* 5, 3 (2020), 4867–4873.

[26] Lang Peng, Zhirong Chen, Zhangjie Fu, Pengpeng Liang, and Erkang Cheng. 2023. BEVSegFormer: Bird's Eye View Semantic Segmentation From Arbitrary Camera Rigs. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 5935–5943.

[27] Limeng Qiao, Wenjie Ding, Xi Qiu, and Chi Zhang. 2023. End-to-End Vectorized HD-Map Construction With Piecewise Bezier Curve. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13218–13228.

[28] Limeng Qiao, Yongchao Zheng, Peng Zhang, Wenjie Ding, Xi Qiu, Xing Wei, and Chi Zhang. 2023. MachMap: End-to-End Vectorized Solution for Compact HD-Map Construction. *arXiv preprint arXiv:2306.10301* (2023).

[29] NDS Open Lane Model 1.0 Release. 2019. http://www.openlanemodel.org/.

[30] Lucas Tabelini, Rodrigo Berriel, Thiago M Paixao, Claudine Badue, Alberto F De Souza, and Thiago Oliveira-Santos. 2021. Keep your eyes on the lane: Real-time attention-guided lane detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 294–302.

[31] Andrew Tao, Karan Sapra, and Bryan Catanzaro. 2020. Hierarchical multi-scale attention for semantic segmentation. *arXiv preprint arXiv:2005.10821* (2020).

[32] Andrew Tao, Karan Sapra, and Bryan Catanzaro. 2020. Hierarchical multi-scale attention for semantic segmentation. *arXiv preprint arXiv:2005.10821* (2020).

[33] Jinsheng Wang, Yinchao Ma, Shaofei Huang, Tianrui Hui, Fei Wang, Chen Qian, and Tianzhu Zhang. 2022. A keypoint-based global association network for lane detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1392–1401.

[34] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingkui Tan, Xinggang Wang, et al. 2020. Deep high-resolution representation learning for visual recognition. *IEEE transactions on pattern analysis and machine intelligence* 43, 10 (2020), 3349–3364.

[35] Dongming Wu, Fan Jia, Jiahao Chang, Zhuoling Li, Jianjian Sun, Chunrui Han, Shuailin Li, Yingfei Liu, Zheng Ge, and Tiancai Wang. 2023. The 1st-place Solution for CVPR 2023 OpenLane Topology in Autonomous Driving Challenge. *arXiv preprint arXiv:2306.09590* (2023).

[36] Deguo Xia, Jizhou Huang, Jianzhong Yang, Xiyan Liu, and Haifeng Wang. 2022. DuARUS: Automatic Geo-object Change Detection with Street View Imagery for Updating Road Database at Baidu Maps. In *Proceedings of the 31st ACM International Conference on Information and Knowledge Management*.

[37] Deguo Xia, Xiyan Liu, Wei Zhang, Hui Zhao, Chengzhou Li, Weiming Zhang, Jizhou Huang, and Haifeng Wang. 2022. DuTraffic: Live Traffic Condition Prediction with Trajectory Data and Street Views at Baidu Maps. In *Proceedings of the 31st ACM International Conference on Information and Knowledge Management*.

[38] Jianzhong Yang, Xiaoqing Ye, Bin Wu, Yanlei Gu, Ziyu Wang, Deguo Xia, and Jizhou Huang. 2022. DuARE: Automatic Road Extraction with Aerial Images and Trajectory Data at Baidu Maps. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 4321–4331.

[39] Gongjie Zhang, Jiahao Lin, Shuang Wu, Yilin Song, Zhipeng Luo, Yang Xue, Shijian Lu, and Zuoguan Wang. 2023. Online Map Vectorization for Autonomous Driving: A Rasterization Perspective. *arXiv preprint arXiv:2306.10502* (2023).

[40] Zhixin Zhang, Yiyuan Zhang, Xiaohan Ding, Fusheng Jin, and Xiangyu Yue. 2023. Online Vectorized HD Map Construction using Geometry. *arXiv preprint arXiv:2312.03341* (2023).