

000 001 002 003 004 005 UNI-MAP: UNIFIED CAMERA-LIDAR PERCEPTION 006 FOR ROBUST HD MAP CONSTRUCTION 007 008 009

010 **Anonymous authors**
 011 Paper under double-blind review
 012
 013
 014
 015
 016
 017
 018
 019
 020
 021
 022
 023
 024
 025
 026
 027
 028
 029
 030
 031
 032

ABSTRACT

033 High-definition (HD) map construction methods play a vital role in providing pre-
 034 cise and comprehensive static environmental information essential for autonomous
 035 driving systems. The primary sensors used are cameras and LiDAR, with input
 036 configurations varying among camera-only, LiDAR-only, or camera-LiDAR fusion
 037 based on cost-performance considerations, while fusion-based methods typically
 038 perform the best. However, current methods face two major issues: high costs
 039 due to separate training and deployment for each input configuration, and low
 040 robustness when sensors are missing or corrupted. To address these challenges, we
 041 propose the Unified Robust HD Map Construction Network (Uni-Map), a single
 042 model designed to perform well across all input configurations. Our approach
 043 designs a novel **Mixture Stack Modality (MSM)** training scheme, allowing the
 044 map decoder to learn effectively from camera, LiDAR, and fused features. We
 045 also introduce a projector module to align Bird's Eye View features from different
 046 modalities into a shared space, enhancing representation learning and overall model
 047 performance. During inference, our model utilizes a switching modality strategy
 048 to adapt seamlessly to any input configuration, ensuring compatibility across vari-
 049 ous modalities. To evaluate the robustness of HD map construction methods, we
 050 designed 13 different sensor corruption scenarios and conducted extensive exper-
 051 iments comparing Uni-Map with state-of-the-art methods. Experimental results
 052 show that Uni-Map outperforms previous methods by a significant margin across
 053 both normal and corrupted modalities, demonstrating superior performance and ro-
 054 bustness. Notably, our unified model surpasses independently trained camera-only,
 055 LiDAR-only, and camera-LiDAR MapTR models with a gain of 4.6, 5.6, and 5.6
 056 mAP on the nuScenes dataset, respectively. The source code will be released.
 057
 058

059 1 INTRODUCTION

060 Online high-definition (HD) map provides abundant and precise static environmental information
 061 about the driving scenes, which is fundamental for planning and navigation in autonomous driving
 062 systems. Cameras and LiDAR are the predominant sensors, offering semantic-rich image data and
 063 explicit geometric information from point clouds, respectively. HD map construction models can be
 064 categorized into three groups based on input configurations: camera-only Qiao et al. (2023); Ding et al.
 065 (2023); Yuan et al. (2024); Liu et al. (2024a); Li et al. (2024), LiDAR-only Li et al. (2022a); Liu et al.
 066 (2023a), and camera-LiDAR fusion Liao et al. (2023a;b); Zhou et al. (2024) models. As illustrated in
 067 Fig. 1 (a)-(c), HD map construction methods with different input configurations have been widely
 068 studied and deployed in real-world systems based on different cost-effective considerations.

069 However, existing methods entail the training and deployment of separate models for each input
 070 configuration, resulting in substantial development, maintenance, and deployment overheads. To
 071 address this problem, we propose a novel **Unified Robust HD Map Construction Network (Uni-Map)**,
 072 where one trained model can perform well under all input configurations, depicted in Fig. 1(d).
 073 Our approach elaborates a novel Mixture Stack Modality (MSM) training scheme during the training
 074 phase, allowing the map decoder to glean rich knowledge from the camera, LiDAR, or fused
 075 features. Furthermore, we introduce a novel projector module to map Bird's Eye View (BEV) features
 076 of different modalities into a shared space. During inference, we present a switching modality
 077 strategy enabling precise predictions by Uni-Map when utilizing arbitrary modality inputs. Extensive

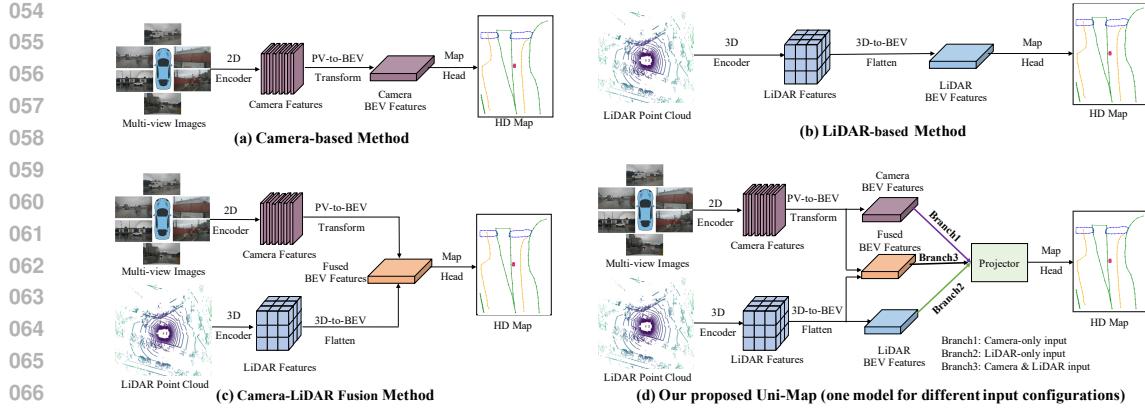


Figure 1: Illustration of the Camera-based method, LiDAR-based method, Camera-LiDAR Fusion method, and the proposed Uni-Map (one model for different input configurations).

experiments demonstrate that Uni-Map can achieve high performance in different input configurations while reducing the training and deployment costs of the model.

Another critical concern of HD map construction methods for autonomous driving is the model’s robustness Kong et al. (2024). While Camera-LiDAR fusion methods have shown promising performance by incorporating information from both modalities Liao et al. (2023a); Zhou et al. (2024); Hao et al. (2024b), existing fusion methods often assume access to complete sensor information, leading to low robustness and potential collapse when sensors are corrupted or missing. To comprehensively evaluate the robustness of the Camera-LiDAR fusion model, we design 13 types of camera-LiDAR corruption combinations that perturb both camera and LiDAR inputs separately or concurrently. These combinations are summarized into 6 cases and illustrated in Fig. 2 (left). We compare Uni-Map with state-of-the-art MapTR Liao et al. (2023a) method, Uni-Map performs more robustly as depicted in Fig. 2 (right), benefiting from the comprehensive feature representations learned by our proposed MSM and aligned by the projector module. Quantitatively, when facing missing camera sensors, Uni-Map still achieves 61.2 mAP, which outperforms the original MapTR Liao et al. (2023a) by +38.7 mAP (61.2 vs. 22.5). Experimental results show that Uni-Map exhibits stronger robustness on various multi-sensor corruption types. Importantly, the core components of Uni-Map, *i.e.*, MSM training scheme, projector module, and the switching modality strategy are simple yet effective plug-and-play techniques compatible with existing pipelines.

In summary, the main contributions of this paper are threefold:

- We propose a novel Unified Robust HD Map Construction Network (Uni-Map), which stands out as an **all-in-one model** to operate on arbitrary input configurations.
- We design a novel **Mixture Stack Modality training scheme** with a simple yet effective projector module to project the BEV features of different modalities into a shared space, allowing the map decoder to learn strong representation from different modalities and a switching modality strategy to utilize arbitrary modality inputs during inference.
- Our single Uni-Map model beats the state-of-the-art MapTR models independently trained on camera-only, LiDAR-only, and camera-LiDAR fusion modalities with a gain of 4.6, 5.6, and 5.6 mAP, respectively. Moreover, Uni-Map shows much better robustness on 13 types of camera-LiDAR corruption combinations. These benefits extend to various map construction models due to our simple, task-independent designs.

2 RELATED WORK

HD Map Construction. HD map construction is a prominent and extensively researched area within the field of autonomous driving. According to the input sensor modality, HD map construction models can be categorized into camera-only Liao et al. (2023a); Zhang et al. (2024c); Ding et al. (2023); Liao et al. (2023b); Yuan et al. (2024), LiDAR-only Li et al. (2022a); Liu et al. (2023a) and camera-LiDAR fusion Liao et al. (2023a;b); Zhou et al. (2024); Hao et al. (2024c) models.

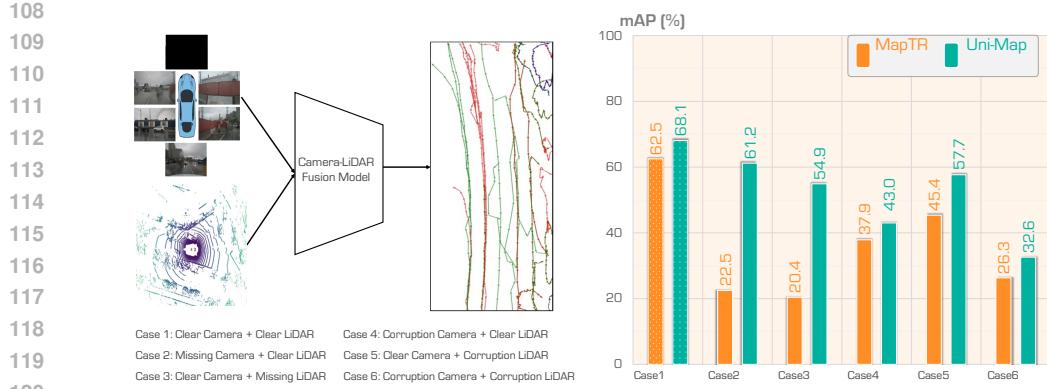


Figure 2: **Uni-Map shows stronger robustness on various multi-sensor corruption types.** We show mAP results for MapTR and Uni-Map models on clean data and each type of multi-sensor corruption. Results show Uni-Map can mitigate the performance drop on sensor missing or corruptions.

Recently, camera-only methods Liu et al. (2024b); Zhang et al. (2024a); Shi et al. (2024); Wang et al. (2024); Zhang et al. (2024b); Peng et al. (2024); Choi et al. (2024); Chen et al. (2024); Jiang et al. (2024); Hu et al. (2024); Hao et al. (2024a) have increasingly employed the Bird’s-eye view (BEV) representation as an ideal feature space for multi-view perception due to its remarkable ability to mitigate scale-ambiguity and occlusion challenges. Various techniques have been proposed and utilized to project perspective view (PV) features into the BEV space by leveraging geometric priors, such as LSS Philion & Fidler (2020), Deformable Attention Li et al. (2022b) and GKT Chen et al. (2022). However, camera-only methods suffer from a lack of explicit depth information. LiDAR-only methods Wang et al. (2023); Li et al. (2022a); Liu et al. (2023a); Liao et al. (2023b;a) benefit from the accurate 3D geometric information from the LiDAR input. However, they struggle to deal with data sparsity and sensing noise problems robustly. Recently, camera-LiDAR feature fusion in the unified BEV space has attracted much attention Liao et al. (2023a;b); Zhou et al. (2024); Dong et al. (2024). BEV-level fusion uses two independent streams that encode the raw inputs from the camera and LiDAR sensors into features within the same BEV space. This fusion at the BEV level incorporates complementary modality features, surpassing unimodal input approaches in performance.

While significant progress has been made using various methods with different input configurations (camera-only, LiDAR-only, camera-LiDAR fusion) chosen based on cost-performance considerations, a common challenge persists. Current methods necessitate training and deploying separate models for each input configuration, leading to considerable costs in development, maintenance, and deployment. In this paper, we introduce a novel Unified Robust HD map construction approach to address this issue. This method enables training a single model capable of operating on any input configuration, thereby streamlining the process.

Robustness Under Sensor Failures. Sensor failures can significantly impact the accuracy of HD map tasks, thereby jeopardizing the safety of autonomous driving. While Camera-LiDAR fusion methods have shown promising performance, which can make use of both the semantic-rich information from cameras and the explicit geometric information from LiDAR, existing fusion methods often assume access to complete sensor information from both cameras and LiDAR, leading to low robustness in the face of sensor missing or corruptions. This means that their performance may degrade significantly or even fail entirely when sensor data is incomplete or corrupted. Recently, there have been a few studies that focus on benchmarking and improving the robustness under natural corruptions, particularly in various BEV perception algorithms such as 3D object detection Li et al. (2022b); Liu et al. (2023b); Ge et al. (2023), BEV segmentation Zhang et al. (2022); Zhou & Krähenbühl (2022), occupancy prediction Wei et al. (2023b); Huang et al. (2023), and depth estimation Wei et al. (2023a). However, approaches addressing sensor failures for HD map construction are still under exploration.

In this paper, to explore the camera-LiDAR fusion model robustness on the HD map construction task, we design 13 types of camera-LiDAR corruption combinations that perturb both camera and LiDAR inputs separately or concurrently. Our proposed Uni-Map model demonstrates improved robustness under various sensor failure scenarios. To the best of our knowledge, Uni-Map is the first study to explore the robustness of HD map construction task under multi-sensor corruptions.

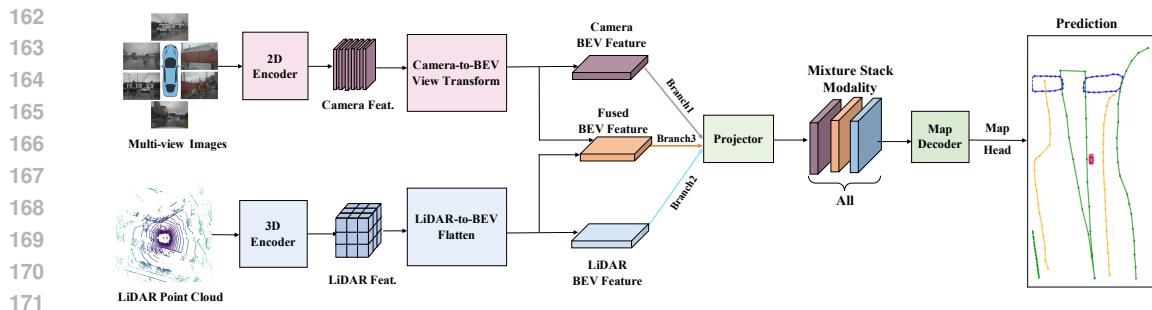


Figure 3: An overview of Uni-Map framework. First, we extract features from multi-modal sensor inputs and convert them into a unified bird’s-eye view (BEV) space efficiently using view transformations. Then, we design a novel Mixture Stack Modality (MSM) scheme with a projector module to re-project the BEV features of different modalities into a shared space. Finally, the mixture stack BEV features are fed into a shared decoder and prediction heads for HD Map construction.

3 METHODOLOGY

Uni-Map pursues a novel Unified Robust HD Map construction approach, which can train an all-in-one model capable of operating with various input configurations. For this purpose, we feed the model decoder with the features from all input configurations at the training stage and process one specific feature based on the deployed input configuration during inference. The overview framework of Uni-Map is shown in Fig. 3. Given different sensory inputs, we first apply modality-specific encoders to extract their features. These multi-modal features are then transformed into a unified BEV representation that preserves both geometric and semantic information. Then, we incorporate a projector module to align BEV features from different modalities into a shared space, thereby enhancing representation learning. Additionally, we introduce a novel Mixture Stack Modality training scheme, enabling the map decoder module to glean rich knowledge from the camera, LiDAR, or fused features. Specifically, the mixture stack BEV features are fed into the decoder and prediction heads for the HD Map construction task. During inference, we employ a switching modality strategy, enabling Uni-Map to make precise predictions using arbitrary modality inputs.

3.1 PRELIMINARIES

For notation clarity, we first introduce some symbols and definitions used throughout this paper. Our goal is to design a novel Unified Robust HD map construction framework taking arbitrary modal sensor data χ as input and predicting vectorized map elements in BEV space, and the types of the map elements (supported types are road boundary, lane divider, and pedestrian crossing). Formally, assume that we have a set of inputs, $\chi = \{Camera, LiDAR\}$, containing multi-view RGB camera images in perspective view, $Camera \in \mathbb{R}^{B \times N^{cam} \times H^{cam} \times W^{cam} \times 3}$, where $B, N^{cam}, H^{cam}, W^{cam}$ denote batch, number of cameras, image height, and image width, respectively, as well as a LiDAR point cloud, $LiDAR \in \mathbb{R}^{B \times P \times 5}$, with number of points P . Each point consists of its 3-dimensional coordinates, reflectivity, and beam index. The detailed architectural designs are described as follows.

3.2 MAP ENCODER

We build our Map Encoder based on the state-of-the-art HD map construction method MapTR Liao et al. (2023a), which applies modality-specific encoders to extract their features and transforms multi-modal features into a unified BEV representation that preserves both geometric and semantic information. Note that our approach is compatible with other Map Encoders that can also be employed to generate camera-only, LiDAR-only, and camera-LiDAR fusion BEV features.

Camera to BEV. For camera images, we first utilize Resnet50 He et al. (2016b) as the backbone to extract the multi-view features. Then we adopt GKT Chen et al. (2022) as the 2D-to-BEV transformation module to convert the multi-view features into BEV space. The generated BEV features can be denoted as $F_{Camera}^{BEV} \in \mathbb{R}^{B \times H \times W \times C}$, where H, W, C represent the height, width, and the number of channels of BEV features, respectively.

LiDAR to BEV. For the LiDAR points, we follow SECOND Yan et al. (2018) in using voxelization and a sparse LiDAR encoder. The LiDAR features are projected to BEV space using a flattening operation as in Liu et al. (2023b), to obtain the unified LiDAR BEV representation $F_{LiDAR}^{BEV} \in \mathbb{R}^{B \times H \times W \times C}$.

Fused BEV. We utilize a convolution-based fusion method Liao et al. (2023a); Zhou et al. (2024) to effectively fuse the BEV features from both camera and LiDAR sensors. More specifically, we utilize concatenation followed by convolution to fuse features from multi-modal BEV feature inputs, $F_{Camera}^{BEV} \in \mathbb{R}^{B \times H \times W \times C}$ and $F_{LiDAR}^{BEV} \in \mathbb{R}^{B \times H \times W \times C}$, resulting in the aggregated features $F_{Fused}^{BEV} \in \mathbb{R}^{B \times H \times W \times C}$.

3.3 MIXTURE STACK MODALITY AND PROJECTOR

In this section, we first introduce the projector module that aims to align BEV features from different modalities into a shared space, thereby enhancing representation learning and overall model performance. Then, we offer the details of the Mixture Stack Modality (MSM) training scheme, which enables the map decoder module to learn rich knowledge from the camera, LiDAR, or fused features.

Projector Module. After input sensor features converted to the shared BEV representation, we can easily obtain the BEV features of the three modalities, *i.e.*, $F_{Camera}^{BEV} \in \mathbb{R}^{B \times H \times W \times C}$, $F_{LiDAR}^{BEV} \in \mathbb{R}^{B \times H \times W \times C}$ and $F_{Fused}^{BEV} \in \mathbb{R}^{B \times H \times W \times C}$. While in the same space, camera BEV features, LiDAR BEV features, and fused BEV features can still be misaligned to some extent due to the inaccurate depth in the view transformer and the large modality gap (See Fig. 7 (a)). Existing works Liang et al. (2022); Liu et al. (2023b) show the phenomenon of modal gaps, *i.e.*, the features of different BEV modalities usually focus on completely separate regions in BEV space. Thus, we propose a projector module to align BEV features from different modalities into a shared space (see the *Remarks* below), thereby enhancing representation learning. To address this issue, we project BEV features of different modalities into a new shared space via a learnable projector $\text{projector}(\cdot)$, *i.e.*,

$$\hat{F}_{camera}^{BEV} = \text{projector}(F_{camera}^{BEV}), \quad (1)$$

$$\hat{F}_{LiDAR}^{BEV} = \text{projector}(F_{LiDAR}^{BEV}), \quad (2)$$

$$\hat{F}_{Fused}^{BEV} = \text{projector}(F_{Fused}^{BEV}), \quad (3)$$

where $\text{projector}(\cdot)$ is the multi-layer linear perceptron (MLP) function. Note that, the BEV features of different modalities use a shared projector, and the details are discussed in the ablation experiments.

Mixture Stack Modality Training Scheme. The map decoder module in existing HD map construction methods is typically trained using BEV features from a single mode, limiting it to one input configuration. To address this limitation and ensure that a single trained model can perform well across all input configurations, we introduce a novel Mixture Stack Modality training scheme after the projector module. Specifically, it can be formulated as:

$$\hat{F}_{Stack}^{BEV} = \text{Stack}(\hat{F}_{camera}^{BEV}, \hat{F}_{LiDAR}^{BEV}, \hat{F}_{Fused}^{BEV}). \quad (4)$$

Using the MSM scheme, we obtain the stacked multi-modal BEV feature $\hat{F}_{Stack}^{BEV} \in \mathbb{R}^{3B \times H \times W \times C}$, which serves as input for the HD map construction task. Notably, the stacking operation preserves the feature map shape as $H \times W \times C$ by stacking along the batch dimension. This design choice enables seamless integration with the subsequent Map Decoder module in existing methods, such as MapTR Liao et al. (2023a). Consequently, our method operates in a plug-and-play manner, ensuring easy implementation and compatibility.

Remarks: The MSM scheme offers three key advantages. First, by stacking BEV features from different modalities that share the *same* map decoder and ground truth labels, the projector module is supervised (via gradient back-propagation) to implicitly align BEV features from different modalities in the shared feature space. Second, inputting stacked BEV features into the same map decoder increases the diversity of the BEV feature space accessible to the decoder module, thereby improving the model’s generalization ability and robustness across different input configurations. Third, this scheme allows the map decoder module to process BEV features of different modalities. As a result, Uni-Map can flexibly handle various input configurations during inference.

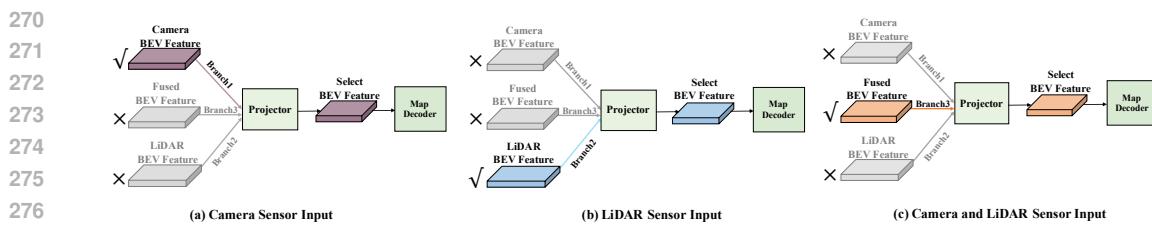


Figure 4: Illustration of the switching modality strategy.

3.4 FULL OBJECTIVE AND INFERENCE

Overall Training. We follow the MapTR Liao et al. (2023a) model’s training loss function, which is composed of three parts, including the classification loss \mathcal{L}_{cls} , the point2point loss \mathcal{L}_{p2p} , and the edge direction loss \mathcal{L}_{dir} . Combining these loss terms, the overall objective function can be formulated as:

$$\mathcal{L} = \lambda_1 \mathcal{L}_{cls} + \lambda_2 \mathcal{L}_{p2p} + \lambda_3 \mathcal{L}_{dir}, \quad (5)$$

where λ_1 , λ_2 and λ_3 are hyper parameters for balancing these terms. For all experiments, λ_1 is set to 2, λ_2 is set to 5, and λ_3 is set to $5e^{-3}$.

Inference Phase. During inference, our model utilizes a switching modality strategy to seamlessly adapt to arbitrary modality inputs, ensuring compatibility across various input configurations. The switching modality strategy can be formulated as:

$$\hat{F}_{Select}^{BEV} = \begin{cases} \hat{F}_{camera}^{BEV}, & \text{if Camera only sensor input,} \\ \hat{F}_{lidar}^{BEV}, & \text{if LiDAR only sensor input,} \\ \hat{F}_{fused}^{BEV}, & \text{if Camera and LiDAR are both obtained.} \end{cases} \quad (6)$$

This switching strategy simulates real-world scenarios where sensors may be missing during the inference phase. As shown in Fig. 4, if the LiDAR data is not available due to uninstallation or damage, we select camera BEV feature \hat{F}_{camera}^{BEV} as the input of the map decoder, and vice versa. Moreover, if Camera and LiDAR data are both obtained, we select Fused BEV features \hat{F}_{fused}^{BEV} as the input for the map decoder. As a result, Uni-Map is compatible with any of the above three input configurations, thereby enhancing its practicality in autonomous driving.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETTINGS

Datasets. We evaluate our method on the widely-used challenging nuScenes Caesar et al. (2020) dataset following the standard setting of previous methods Liao et al. (2023a); Zhou et al. (2024). The nuScenes dataset contains 1,000 sequences of recordings collected by autonomous driving cars. Each sample is annotated at 2Hz and contains 6 camera images covering 360° horizontal FOV of the ego-vehicle. Following Liao et al. (2023a); Li et al. (2022a); Gao et al. (2024), three kinds of map elements are chosen for fair evaluation – pedestrian crossing, lane divider, and road boundary.

Evaluation Metrics. We adopt the evaluation metrics consistent with previous works Liao et al. (2023a); Li et al. (2022a); Zhang et al. (2024c), where average precision (AP) is used to evaluate the map construction quality and Chamfer distance $D_{Chamfer}$ determines the matching between predictions and ground truth. We calculate the AP_τ under several $D_{Chamfer}$ thresholds ($\tau \in T = \{0.5m, 1.0m, 1.5m\}$), and then average across all thresholds as the final mean AP (mAP) metric. The perception ranges are $[-15.0m, 15.0m]/[-30.0m, 30.0m]$ for X/Y-axes.

Implementation Details. Uni-Map is trained with 4 NVIDIA RTX A6000 GPUs. During the training phase, the GT labels are duplicated twice and stacked to form $3B$ batch dimension, matching with the stacked feature map from Eq. 4. The design choice of the MSM scheme is discussed in the ablation studies. For the projector module, we use a two-layer perceptron whose dimension is C->C/2->C. We adopt the AdamW optimizer Loshchilov & Hutter (2019) for all our experiments. We set the mini-batch size to 16, and use a step-decayed learning rate with an initial value of $4e^{-3}$. The inference time is measured on a single NVIDIA RTX A6000 GPU with batch size 1.

Table 1: Comparisons with state-of-the-art methods on nuScenes val set. “L” and “C” represent LiDAR and camera, respectively. “Effi-B0”, “R50”, “PP”, and “Sec” are short for EfficientNet-B0 Tan & Le (2019), ResNet50 He et al. (2016a), PointPillars Lang et al. (2019) and SECOND Yan et al. (2018), respectively. Note that Uni-Map (MapModel) means our method is integrated into an existing MapModel. Best viewed in color. nfan

Method	Modality	BEV Encoder	Backbone	Epoch	AP _{ped.}	AP _{div.}	AP _{bou.}	mAP ↑
HDMapNet Li et al. (2022a)	C	NVT	Effi-B0	30	14.4	21.7	33.0	23.0
VectorMapNet Liu et al. (2023a)	C	IPM	R50	110	36.1	47.3	39.3	40.9
PivotNet Ding et al. (2023)	C	PersFormer	R50	30	53.8	58.8	59.6	57.4
BeMapNet Qiao et al. (2023)	C	IPM-PE	R50	30	57.7	62.3	59.4	59.8
MapVR Zhang et al. (2024c)	C	GKT	R50	24	47.7	54.4	51.4	51.2
MapTrv2 Liao et al. (2023b)	C	BEVPoolv2	R50	24	59.8	62.4	62.4	61.5
StreamMapNet Yuan et al. (2024)	C	BEVFormer	R50	30	61.7	66.3	62.1	63.4
MapTR Liao et al. (2023a)	C	GKT	R50	24	46.3	51.5	53.1	50.3
HIMap Zhou et al. (2024)	C	BEVFormer	R50	24	62.2	66.5	67.9	65.5
Uni-Map (MapTR)	C	GKT	R50	24	52.1	57.5	55.2	54.9
Uni-Map (HIMap)	C	BEVFormer	R50	24	64.5	68.2	68.3	67.0
VectorMapNet Liu et al. (2023a)	L	-	PP	110	25.7	37.6	38.6	34.0
MapTrv2 Liao et al. (2023b)	L	-	Sec	24	56.6	58.1	69.8	61.5
MapTR Liao et al. (2023a)	L	-	Sec	24	48.5	53.7	64.7	55.6
HIMap Zhou et al. (2024)	L	-	Sec	24	54.8	64.7	73.5	64.3
Uni-Map (MapTR)	L	-	Sec	24	56.5	57.8	69.4	61.2
Uni-Map (HIMap)	L	-	Sec	24	65.3	69.5	77.8	70.8
MapTrv2 Liao et al. (2023b)	C & L	BEVPoolv2	R50 & Sec	24	65.6	66.5	74.8	69.0
MapTR Liao et al. (2023a)	C & L	GKT	R50 & Sec	24	55.9	62.3	69.3	62.5
HIMap Zhou et al. (2024)	C & L	BEVFormer	R50 & Sec	24	71.0	72.4	79.4	74.3
Uni-Map (MapTR)	C & L	GKT	R50 & Sec	24	64.4	66.8	73.2	68.1
Uni-Map (HIMap)	C & L	BEVFormer	R50 & Sec	24	73.6	75.3	81.2	76.7

Table 2: Comparison of MapTR Liao et al. (2023a) and Uni-Map in terms of accuracy, model size, training epochs and training time on nuScenes dataset. Note that only one Uni-Map model is trained while three MapTR models (MapTR-C, MapTR-L, and MapTR-F) are trained for different input configurations. † represents using the total time of training three MapTR models to train our Uni-Map model.

Method	Camera-only (mAP)	LiDAR-only (mAP)	Camera & LiDAR (mAP)	Params(MB)	Epoch	Training Time
MapTR-C	50.3	—	—	35.9	24	13h55m
MapTR-L	—	55.6	—	14.3	24	9h7m
MapTR-F	—	—	62.5	39.8	24	15h44m
Uni-Map (MapTR)	54.9	61.2	68.1	39.9	24	21h57m
Uni-Map (MapTR)†	57.2	64.5	70.4	39.9	42	38h44m

4.2 COMPARISON WITH THE STATE-OF-THE-ARTS

With the same settings, we compare our method with several state-of-the-art methods across three categories, *i.e.*, camera-only methods, LiDAR-only methods, and camera-LiDAR fusion methods. Specifically, we integrate our Uni-Map into two recent methods, MapTR Liao et al. (2023a) and HIMap Zhou et al. (2024), where we insert the projector module into these models and apply the MSM training scheme. Moreover, to fairly evaluate the effectiveness, we train the same epochs as the original model. It’s noteworthy that while three MapTR/HIMap models need to be trained for different input configurations, our Uni-Map model only requires training once. As shown in Tab. 1, our Uni-Map significantly improves the performance compared to the original models. Specifically, Uni-Map (MapTR) outperforms independently trained camera-only, LiDAR-only, and camera-LiDAR MapTR models on NuScenes with a large gain of 4.6, 5.6, and 5.6 mAP, under the respective input configurations, respectively. Based on the previous state-of-the-art HIMap, our all-in-one model surpasses HIMap-C, HIMap-L, and HIMap-F by 1.5, 6.5, and 2.4 mAP respectively, establishing a new state-of-the-art in vectorized map reconstruction. Results for more datasets like Argoverse2 Wilson et al. (2021) are shown in the supplementary material A.3. All these results prove the effectiveness of our design.

Model Size, Training Time, GPU Memory and Inference Speed. To systematically evaluate the effectiveness of our proposed Uni-Map model, we comprehensively analyze it in terms of accuracy, model size, training time, and inference speed. The experimental results are shown in Tab. 2 and Appendix Tab. 6–Tab. 7. The experimental results reveal some interesting findings: (1) Compared with MapTR, Uni-Map performs much better in all input configurations in both single-class APs

378 Table 3: Ablation study on the MSM training scheme. The mAP values on nuScenes val set are
 379 reported. ‘Mean’ represents the average mAP of three input configurations.
 380

Random Select	Mixture Stack	Projector	Camera-only	LiDAR-only	Camera & LiDAR	Mean
✗	✗	✗	20.4	22.5	62.5	35.1
✓	✗	✗	36.9	47.5	62.9	49.1
✗	✓	✗	53.7	59.4	67.9	60.3
✓	✗	✓	45.6	55.3	61.2	54.0
✗	✓	✓	54.9	61.2	68.1	61.4

387 Table 4: Ablation study on Projector Module. The mAP values on nuScenes val set are reported.
 388 ‘Mean’ represents the average mAP of three input configurations.
 389

Method	Camera-only	LiDAR-only	Camera&LiDAR	Mean
Baseline (w/o projector)	53.7	59.4	67.9	60.3
Variant 1: Independent Projector	53.6	62.2	67.6	61.1
Variant 2:Partially Shared Projector	53.3	61.5	68.0	60.9
Variant 3: Skip Shared Projector	53.4	61.7	68.0	61.0
Variant 4: Shared Projector (Ours)	54.9	61.2	68.1	61.4

396 and the overall mAP. Note that only one Uni-Map model is trained while three MapTR models
 397 (MapTR-C, MapTR-L, and MapTR-F) are trained for different input configurations. Thus, we use
 398 the same computational budget of training three MapTR models to train our Uni-Map model, and
 399 the resulting Uni-Map model (last row of Tab. 2) beats independently trained camera-only, LiDAR-
 400 only, and camera-LiDAR fusion MapTR models with a larger gain of 6.9, 8.9, 7.9 mAP, under the
 401 respective input configurations. (2) In terms of model size, our Uni-Map model only increases the
 402 number of parameters by 0.1MB compared to the MapTR-F model, as shown in Tab. 2. It is more
 403 parameter-efficient than deploying the three models simultaneously in practice. (3) In terms of GPU
 404 Memory and inference speed, the quantities of our Uni-Map and MapTR are almost the same, as
 405 shown in Appendix Tab. 6-Tab. 7. All in all, the Uni-Map model achieves significant performance
 406 improvements over the strong MapTR baseline with less training time and fewer parameters (for
 407 various input configurations), while maintaining the same inference speed and memory footprint.

4.3 ABLATION STUDIES

410 **Analysis of the MSM training scheme.** To systematically evaluate the effectiveness of the MSM
 411 training scheme, we train the model using different schemes and report the mAP results in Tab. 3. In
 412 addition to MSM, we also introduce the Random Select Modality (RSM) training scheme that receives
 413 inputs from one BEV feature map randomly selected among \hat{F}_{camera}^{BEV} , \hat{F}_{LiDAR}^{BEV} , \hat{F}_{Fused}^{BEV} . In the main
 414 ablation study, we design the following model variants: (1) We train the model without the projector
 415 module and any of the RSM and MSM training schemes. (2) We train the model without the projector
 416 module using RSM or MSM training schemes, respectively. (3) We train the model with the projector
 417 module using RSM or MSM training schemes, respectively. The experimental results reveal some
 418 interesting findings: (1) The results of both RSM and MSM schemes are significantly better than the
 419 Baseline model (only learned/seen the BEV features of one modality), verifying the effectiveness of
 420 learning with rich knowledge from different BEV features to improve the generalization ability of the
 421 map decoder. (2) The results of the RSM training scheme are inferior to the MSM training scheme
 422 under both settings (with and without the Projector). This demonstrates the MSM training scheme’s
 423 advantage in enhancing the map decoder’s effective use of camera, LiDAR, and fused features. This
 424 increases the diversity of the BEV feature space, resulting in a high-performance integrated model.

425 **Analysis on projector module.** We investigate the design choice of the projector module in our
 426 method. The ablation variants include Independent Projector, Partially Shared Projector, Skip
 427 Shared Projector, and Shared Projector (the default setting). The detailed formulation of the variant
 428 projector module is in the supplementary material A.1. As shown in Tab. 4, the experimental results
 429 reveal some interesting findings: (1) Using different projector variants consistently outperforms the
 430 baseline model, implying that using the simple projector module can facilitate learning better feature
 431 representations. This can be owing to the fact that our model uses the same map decoder and ground
 432 truth labels to promote feature alignment in this latent space. (2) Using a shared projector module
 433 consistently outperforms other projector variants. It is reasonable that using BEV feature information

from different modalities to perform gradient updates on a shared projector, rather than on multiple projectors, aligns BEV features from different modalities more effectively. These observations validate the effectiveness of the projector module in aligning BEV features from different modalities into a shared space, thereby enhancing representation learning and overall model performance.

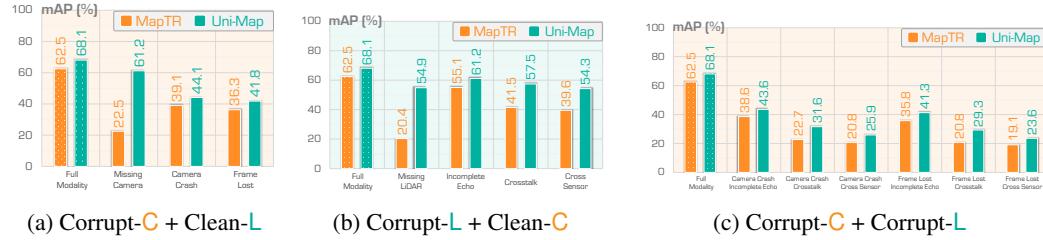


Figure 5: The result of multi-sensor corruption on MapTR vs. Uni-Map (MapTR) fusion model.

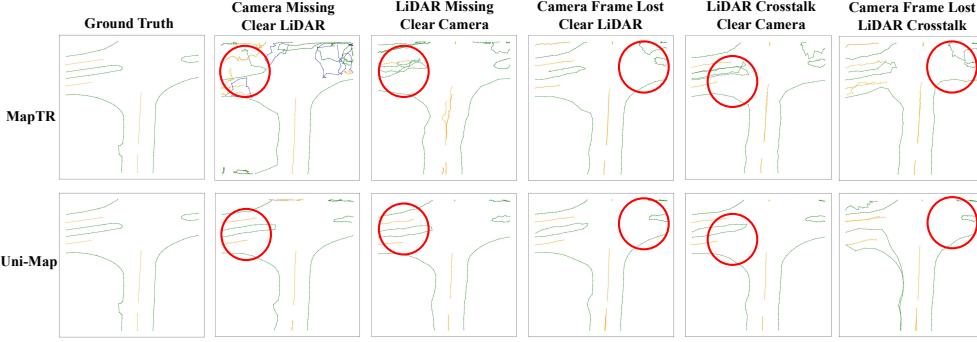


Figure 6: Qualitative results of the nuScenes val set on the MapTR and UniMap models respectively.

4.4 ROBUSTNESS OF MULTI-SENSOR CORRUPTIONS

To explore the camera-LiDAR fusion model robustness, we design 13 types of camera-LiDAR corruption combinations that perturb both camera and LiDAR inputs separately or concurrently. Camera-LiDAR corruption combinations are grouped into camera-only corruptions, LiDAR-only corruptions, and their combinations, covering the majority of real-world corruption cases. The definition of multi-sensor corruption is detailed in A.2. Fig. 5 shows the results of three Camera-LiDAR corruption combinations. We have the following observations. (1) In the sensor missing , Uni-Map can prevent the model from collapsing owing to the switching modality strategy. Quantitatively, when facing a missing LiDAR sensor, Uni-Map still achieves 54.9 mAP, which outperforms the original MapTR Liao et al. (2023a) by 34.5 mAP. (2) In case of the corruption of the camera and LiDAR sensor individually or simultaneously, Uni-Map shows stronger robustness. For example, in the face of camera frame lost and LiDAR crosstalk, compared to the MapTR fused model, the Uni-Map model achieved significant improvements in 8.5 mAP (29.3 vs. 20.8). These results demonstrate that the MSM training scheme enhances the generalization ability of the map decoder. By stacking BEV features from different modalities into the same map decoder, the diversity of the BEV feature space accessible to the decoder increases, thereby improving the model’s robustness. All in all, Uni-Map shows stronger robustness on our designed 13 types of camera-LiDAR corruption combinations.

4.5 VISUALIZATION

Qualitative Results. To further analyze the effectiveness of our Uni-Map model, we compare it with MapTR Liao et al. (2023a) and present the qualitative results in Fig. 6. We compare the predicted vectorized HD map results of different settings, including the camera sensor missing, LiDAR sensor missing, camera frame lost and clear LiDAR, LiDAR crosstalk and clear camera, and camera frame lost with LiDAR crosstalk. We observe that the baseline MapTR predictions are highly erroneous,

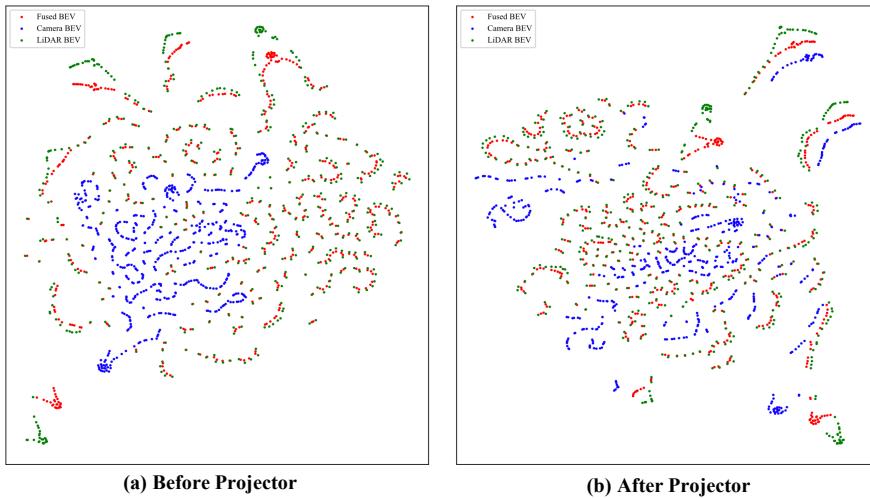


Figure 7: The t-SNE visualizations of (a) Before Projector module and (b) After Projector module. Red/Blue/green denotes fused BEV feature/camera BEV feature/LiDAR BEV feature. After the projector module, the BEV features from different modalities are aligned in a shared space, e.g., red, blue, and green circles are close together after the projector module (best viewed in color).

whereas our Uni-Map model can already correct significant errors in the baseline predictions in all settings. All in all, our model shows significant advantages in clear and various corruption situations.

t-SNE. We randomly choose 500 samples on the nuScenes dataset and show the tSNE Van der Maaten & Hinton (2008) visualizations of (a) Before Projector module and (b) After Projector module in Fig. 7. Red/Blue/green denotes fused BEV feature/camera BEV feature/LiDAR BEV feature. As can be seen, Fig. 7 (a) Before Projector module shows that blue and red/green features are clearly separated, indicating that although in the same space, camera BEV features, LiDAR BEV features, and fused BEV features can still be misaligned to some extent due to the inaccurate depth in the view transformer and the large modality gap. Fig. 7 (b) After the projector module, the BEV features from different modalities are aligned in a shared space, *i.e.*, red, blue, and green circles are close together after the projector module.

5 CONCLUSION

In this paper, we propose a novel Unified Robust HD Map Construction Network (Uni-Map), which can train an all-in-one model to operate on arbitrary input configurations. The core components of Uni-Map, *i.e.* MSM training scheme, projector module, and the switching modality strategy, are simple yet effective plug-and-play techniques compatible with existing pipelines. Extensive experiments demonstrate that Uni-Map can achieve high performance in different input configurations while reducing the training and deployment costs of the model. Moreover, Uni-Map shows stronger robustness on our designed 13 types of camera-LiDAR corruption combinations. We hope that our method can be applied to more autonomous driving perception tasks.

Ethics Statement. Our work can boost the performance and robustness of HD map construction task. Although our method significantly improves the robustness of the HD map model, the overall robustness is still low. Special caution is needed in deploying our methods onto vehicles on the road to ensure safety. Therefore, future research is necessary to further investigate more advanced robustness methods.

Reproducibility. To ensure the reproducibility of our work, we have included a comprehensive Reproducibility Statement. Specifically, for the novel model and algorithms presented in this work, we will make them open source upon paper acceptance. Additionally, all multi-sensor corruption details and more experimental results can be found in Appendix A. For the datasets used in our experiments, we follow the standard protocol of the open source work MapTR Liao et al. (2023a) . This Reproducibility Statement is intended to guide readers to the relevant resources that will aid in replicating our work, ensuring transparency and clarity throughout.

540 REFERENCES
541

- 542 Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush
543 Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for
544 autonomous driving. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
545 11618–11628, 2020.
- 546 Jiacheng Chen, Yuefan Wu, Jiaqi Tan, Hang Ma, and Yasutaka Furukawa. Maptracker: Tracking with
547 strided memory fusion for consistent vector hd mapping. *arXiv preprint arXiv:2403.15951*, 2024.
- 548 Shaoyu Chen, Tianheng Cheng, Xinggang Wang, Wenming Meng, Qian Zhang, and Wenyu Liu.
549 Efficient and robust 2d-to-bev representation learning via geometry-guided kernel transformer.
550 *arXiv preprint arXiv:2206.04584*, 2022.
- 551 Sehwan Choi, Jungho Kim, Hongjae Shin, and Jun Won Choi. Mask2map: Vectorized hd map
552 construction using bird’s eye view segmentation masks. *arXiv preprint arXiv:2407.13517*, 2024.
- 553 Wenjie Ding, Limeng Qiao, Xi Qiu, and Chi Zhang. Pivotnet: Vectorized pivot learning for end-to-end
554 hd map construction. In *Proceedings of the IEEE/CVF International Conference on Computer
555 Vision*, pp. 3672–3682, 2023.
- 556 Hao Dong, Weihao Gu, Xianjing Zhang, Jintao Xu, Rui Ai, Huimin Lu, Juho Kannala, and Xieyuanli
557 Chen. Superfusion: Multilevel lidar-camera fusion for long-range hd map generation. In *2024
558 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 9056–9062. IEEE, 2024.
- 559 Wenjie Gao, Jiawei Fu, Yanqing Shen, Haodong Jing, Shitao Chen, and Nanning Zheng. Complementing
560 onboard sensors with satellite maps: A new perspective for hd map construction. In *2024
561 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 11103–11109, 2024.
- 562 Chongjian Ge, Junsong Chen, Enze Xie, Zhongdao Wang, Lanqing Hong, Huchuan Lu, Zhenguo
563 Li, and Ping Luo. Metabev: Solving sensor failures for 3d detection and map segmentation. In
564 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 8721–8731, 2023.
- 565 Xiaoshuai Hao, Ruihai Li, Hui Zhang, Dingzhe Li, Rong Yin, Sangil Jung, Seung-In Park, ByungIn
566 Yoo, Haimei Zhao, and Jing Zhang. Mapdistill: Boosting efficient camera-based hd map
567 construction via camera-lidar fusion model distillation. In *European Conference on Computer Vision*,
568 2024a.
- 569 Xiaoshuai Hao, Mengchuan Wei, Yifan Yang, Haimei Zhao, Hui Zhang, Yi Zhou, Qiang Wang,
570 Weiming Li, Lingdong Kong, and Jing Zhang. Is your hd map constructor reliable under sensor
571 corruptions? In *Advances in Neural Information Processing System*, 2024b.
- 572 Xiaoshuai Hao, Hui Zhang, Yifan Yang, Yi Zhou, Sangil Jung, Seung-In Park, and ByungIn Yoo.
573 Mbfusion: A new multi-modal bev feature fusion method for hd map construction. In *2024 IEEE
574 International Conference on Robotics and Automation (ICRA)*, pp. 15922–15928, 2024c.
- 575 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
576 recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778,
577 2016a.
- 578 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
579 recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
580 Recognition*, 2016b.
- 581 Haotian Hu, Fanyi Wang, Yaonong Wang, Laifeng Hu, Jingwei Xu, and Zhiwang Zhang. Admap:
582 Anti-disturbance framework for reconstructing online vectorized hd map. *arXiv preprint
583 arXiv:2401.13172*, 2024.
- 584 Yuanhui Huang, Wenzhao Zheng, Yunpeng Zhang, Jie Zhou, and Jiwen Lu. Tri-perspective view for
585 vision-based 3d semantic occupancy prediction. In *Proceedings of the IEEE/CVF Conference on
586 Computer Vision and Pattern Recognition*, pp. 9223–9232, 2023.
- 587 Zhou Jiang, Zhenxin Zhu, Pengfei Li, Huan-ang Gao, Tianyuan Yuan, Yongliang Shi, Hang Zhao,
588 and Hao Zhao. P-mapnet: Far-seeing map generator enhanced by both sdmap and hdmap priors.
589 *arXiv preprint arXiv:2403.10521*, 2024.

- 594 Lingdong Kong, Shaoyuan Xie, Hanjiang Hu, Yaru Niu, Wei Tsang Ooi, Benoit R Cottereau, Lai Xing
 595 Ng, Yuexin Ma, Wenwei Zhang, Liang Pan, et al. The robodrive challenge: Drive anytime anywhere
 596 in any condition. *arXiv preprint arXiv:2405.08816*, 2024.
- 597 Alex H Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. Point-
 598 pillars: Fast encoders for object detection from point clouds. In *Proceedings of the IEEE/CVF*
 599 *conference on computer vision and pattern recognition*, pp. 12697–12705, 2019.
- 600 601 Qi Li, Yue Wang, Yilun Wang, and Hang Zhao. Hdmapnet: An online hd map construction and
 602 evaluation framework. In *IEEE International Conference on Robotics and Automation*, pp. 4628–
 603 4634, 2022a.
- 604 605 Siyu Li, Kailun Yang, Hao Shi, Song Wang, You Yao, and Zhiyong Li. Genmapping: Unleashing the
 606 potential of inverse perspective mapping for robust online hd map construction. *arXiv preprint*
 607 *arXiv:2409.08688*, 2024.
- 608 609 Zhiqi Li, Wenhui Wang, Hongyang Li, Enze Xie, Chonghao Sima, Tong Lu, Yu Qiao, and Jifeng Dai.
 610 Bevformer: Learning bird’s-eye-view representation from multi-camera images via spatiotemporal
 611 transformers. In *European Conference on Computer Vision*, pp. 1–18, 2022b.
- 612 613 Tingting Liang, Hongwei Xie, Kaicheng Yu, Zhongyu Xia, Zhiwei Lin, Yongtao Wang, Tao Tang,
 614 Bing Wang, and Zhi Tang. Bevfusion: A simple and robust lidar-camera fusion framework. In
 615 *Advances in Neural Information Processing Systems*, pp. 10421–10434, 2022.
- 616 617 Bencheng Liao, Shaoyu Chen, Xinggang Wang, Tianheng Cheng, Qian Zhang, Wenyu Liu, and
 618 Chang Huang. Maptr: Structured modeling and learning for online vectorized hd map construction.
 619 In *International Conference on Learning Representations*, 2023a.
- 620 621 Bencheng Liao, Shaoyu Chen, Yunchi Zhang, Bo Jiang, Qian Zhang, Wenyu Liu, Chang Huang, and
 622 Xinggang Wang. Maptrv2: An end-to-end framework for online vectorized HD map construction.
 623 *arXiv preprint arXiv:2308.05736*, 2023b.
- 624 625 Hui Liu, Faliang Chang, Chunsheng Liu, Yansha Lu, and Minhang Liu. Camera-based online
 626 vectorized hd map construction with incomplete observation. *IEEE Robotics and Automation
 Letters*, 2024a.
- 627 628 Xiaolu Liu, Song Wang, Wentong Li, Ruizi Yang, Junbo Chen, and Jianke Zhu. Mgmap: Mask-guided
 629 learning for online vectorized hd map construction. In *Proceedings of the IEEE/CVF Conference
 630 on Computer Vision and Pattern Recognition*, pp. 14812–14821, 2024b.
- 631 632 Yicheng Liu, Tianyuan Yuan, Yue Wang, Yilun Wang, and Hang Zhao. Vectormapnet: End-to-end
 633 vectorized hd map learning. In *International Conference on Machine Learning*, pp. 22352–22369,
 634 2023a.
- 635 636 Zhijian Liu, Haotian Tang, Alexander Amini, Xinyu Yang, Huizi Mao, Daniela L Rus, and Song Han.
 637 Bevfusion: Multi-task multi-sensor fusion with unified bird’s eye view representation. In *IEEE
 International Conference on Robotics and Automation*, pp. 2774–2781, 2023b.
- 638 639 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Confer-
 ence on Learning Representations*, 2019.
- 640 641 Nan Peng, Xun Zhou, Mingming Wang, Xiaojun Yang, Songming Chen, and Guisong Chen. Pre-
 642 vpredmap: Exploring temporal modeling with previous predictions for online vectorized hd map
 643 construction. *arXiv preprint arXiv:2407.17378*, 2024.
- 644 645 Jonah Philion and Sanja Fidler. Lift, splat, shoot: Encoding images from arbitrary camera rigs by
 646 implicitly unprojecting to 3d. In *European Conference on Computer Vision*, pp. 194–210, 2020.
- 647 648 Limeng Qiao, Wenjie Ding, Xi Qiu, and Chi Zhang. End-to-end vectorized hd-map construction with
 649 piecewise bezier curve. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
 650 Pattern Recognition*, pp. 13218–13228, 2023.
- Anqi Shi, Yuze Cai, Xiangyu Chen, Jian Pu, Zeyu Fu, and Hong Lu. Globalmapnet: An online
 framework for vectorized global hd map construction. *arXiv preprint arXiv:2409.10063*, 2024.

- 648 Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks.
 649 In *International conference on machine learning*, pp. 6105–6114, 2019.
- 650
- 651 Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine*
 652 *learning research*, 9(11), 2008.
- 653 Rongxuan Wang, Xin Lu, Xiaoyang Liu, Xiaoyi Zou, Tongyi Cao, and Ying Li. Priormapnet:
 654 Enhancing online vectorized hd map construction with priors. *arXiv preprint arXiv:2408.08802*,
 655 2024.
- 656
- 657 Song Wang, Wentong Li, Wenyu Liu, Xiaolu Liu, and Jianke Zhu. Lidar2map: In defense of lidar-
 658 based semantic map construction using online camera distillation. In *Proceedings of the IEEE/CVF*
 659 *Conference on Computer Vision and Pattern Recognition*, pp. 5186–5195, 2023.
- 660 Yi Wei, Linqing Zhao, Wenzhao Zheng, Zheng Zhu, Yongming Rao, Guan Huang, Jiwen Lu, and
 661 Jie Zhou. Surrounddepth: Entangling surrounding views for self-supervised multi-camera depth
 662 estimation. In *Conference on Robot Learning*, pp. 539–549, 2023a.
- 663
- 664 Yi Wei, Linqing Zhao, Wenzhao Zheng, Zheng Zhu, Jie Zhou, and Jiwen Lu. Surroundocc: Multi-
 665 camera 3d occupancy prediction for autonomous driving. In *Proceedings of the IEEE/CVF*
 666 *International Conference on Computer Vision*, pp. 21729–21740, 2023b.
- 667 Benjamin Wilson, William Qi, Tanmay Agarwal, John Lambert, Jagjeet Singh, Siddhesh Khandelwal,
 668 Bowen Pan, Ratnesh Kumar, Andrew Hartnett, Jhony Kaesemeyer Pontes, Deva Ramanan, Peter
 669 Carr, and James Hays. Argoverse 2: Next generation datasets for self-driving perception and
 670 forecasting. In *Proceedings of the Neural Information Processing Systems Track on Datasets and*
 671 *Benchmarks*, 2021.
- 672 Yan Yan, Yuxing Mao, and Bo Li. SECOND: sparsely embedded convolutional detection. *Sensors*,
 673 pp. 3337, 2018.
- 674
- 675 Tianyuan Yuan, Yicheng Liu, Yue Wang, Yilun Wang, and Hang Zhao. Streammapnet: Streaming
 676 mapping network for vectorized online hd map construction. In *IEEE/CVF Winter Conference on*
 677 *Applications of Computer Vision*, pp. 7356–7365, 2024.
- 678 Chi Zhang, Qi Song, Feifei Li, Yongquan Chen, and Rui Huang. Hybrimap: Hybrid clues utilization
 679 for effective vectorized hd map construction. *arXiv preprint arXiv:2404.11155*, 2024a.
- 680
- 681 Xiaoyu Zhang, Guangwei Liu, Zihao Liu, Ningyi Xu, Yunhui Liu, and Ji Zhao. Enhancing vectorized
 682 map perception with historical rasterized maps. *arXiv preprint arXiv:2409.00620*, 2024b.
- 683 Yunpeng Zhang, Zheng Zhu, Wenzhao Zheng, Junjie Huang, Guan Huang, Jie Zhou, and Jiwen
 684 Lu. Beverse: Unified perception and prediction in birds-eye-view for vision-centric autonomous
 685 driving. *arXiv preprint arXiv:2205.09743*, 2022.
- 686
- 687 Zhixin Zhang, Yiyuan Zhang, Xiaohan Ding, Fusheng Jin, and Xiangyu Yue. Online vectorized hd
 688 map construction using geometry. In *European Conference on Computer Vision*, 2024c.
- 689
- 690 Brady Zhou and Philipp Krähenbühl. Cross-view transformers for real-time map-view semantic
 691 segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*
 692 *recognition*, pp. 13760–13769, 2022.
- 693
- 694 Yi Zhou, Hui Zhang, Jiaqian Yu, Yifan Yang, Sangil Jung, Seung-In Park, and ByungIn Yoo. Himap:
 695 Hybrid representation learning for end-to-end vectorized hd map construction. In *Proceedings of*
 696 *the IEEE/CVF conference on computer vision and pattern recognition*, 2024.
- 697
- 698
- 699
- 700
- 701

702 **A APPENDIX / SUPPLEMENTAL MATERIAL**
 703

704 This supplementary material provides additional details on the proposed method and experimental
 705 results that could not be included in the main manuscript due to page limitations.
 706

- 707 • Section A.1 discusses details of different variant projector modules.
- 708 • Section A.2 provides additional details of the multi-sensor corruptions.
- 709 • Section A.3 complements Argoverse2 dataset experiment results and corresponding analysis.
- 710 • Section A.4 presents the results of the switching modality strategy on the original MapTR
- 711 fusion model.
- 712 • Section A.5 offers more experimental results regarding model robustness.
- 713 • Section A.6 offers 3D object detection results to prove the generalization ability of the
- 714 Uni-Map.
- 715 • Section A.7 includes more visualization results to prove the effectiveness of the Uni-Map.

718 **A.1 VARIANT PROJECTOR MODULE**
 719

720 After input sensor features converted to the shared BEV representation, we can easily obtain the
 721 BEV features of the three modalities, *i.e.*, $F_{\text{Camera}}^{\text{BEV}} \in \mathbb{R}^{B \times H \times W \times C}$, $F_{\text{LiDAR}}^{\text{BEV}} \in \mathbb{R}^{B \times H \times W \times C}$ and
 722 $F_{\text{Fused}}^{\text{BEV}} \in \mathbb{R}^{B \times H \times W \times C}$. While in the same space, camera BEV features, LiDAR BEV features, and
 723 fused BEV features can still be misaligned to some extent due to the inaccurate depth in the view
 724 transformer and the large modality gap (See Fig. 7 (a)). Existing works Liang et al. (2022); Liu
 725 et al. (2023b) show the phenomenon of modal gaps, *i.e.*, the features of different BEV modalities
 726 usually focus on completely separate regions in BEV space. Thus, we propose a projector module to
 727 align BEV features from different modalities into a shared space, thereby enhancing representation
 728 learning. To address this issue, we project BEV features of different modalities into a new shared
 729 space via a learnable projector $\text{projector}(\cdot)$.

730 **Shared Projector.** The Shared Projector formula can be written as:

$$\hat{F}_{\text{camera}}^{\text{BEV}} = \text{projector}(F_{\text{camera}}^{\text{BEV}}), \quad (7)$$

$$\hat{F}_{\text{LiDAR}}^{\text{BEV}} = \text{projector}(F_{\text{LiDAR}}^{\text{BEV}}), \quad (8)$$

$$\hat{F}_{\text{Fused}}^{\text{BEV}} = \text{projector}(F_{\text{Fused}}^{\text{BEV}}), \quad (9)$$

735 where $\text{projector}(\cdot)$ is the two-layer linear perceptron function. Note that, the BEV features of
 736 different modalities use a shared projector module.
 737

738 **Partially Shared Projector.** The main difference from the shared projector is that the first linear
 739 layer of the partially shared projector learns three modes independently, and the second linear layer is
 740 shared.

741 **Independent Projector.** The Independent Projector formula can be written:

$$\hat{F}_{\text{camera}}^{\text{BEV}} = \text{projector}_1(F_{\text{camera}}^{\text{BEV}}), \quad (10)$$

$$\hat{F}_{\text{LiDAR}}^{\text{BEV}} = \text{projector}_2(F_{\text{LiDAR}}^{\text{BEV}}), \quad (11)$$

$$\hat{F}_{\text{Fused}}^{\text{BEV}} = \text{projector}_3(F_{\text{Fused}}^{\text{BEV}}), \quad (12)$$

747 where $\text{projector}(\cdot)$ is the multi-layer linear perceptron function. Note that, the BEV features of
 748 different modalities use different projector modules.
 749

Skip Shared Projector. The Skip Shared Projector formula can be written as:

$$\hat{F}_{\text{camera}}^{\text{BEV}} = \text{projector}(F_{\text{camera}}^{\text{BEV}}) + F_{\text{camera}}^{\text{BEV}}, \quad (13)$$

$$\hat{F}_{\text{LiDAR}}^{\text{BEV}} = \text{projector}(F_{\text{LiDAR}}^{\text{BEV}}) + F_{\text{LiDAR}}^{\text{BEV}}, \quad (14)$$

$$\hat{F}_{\text{Fused}}^{\text{BEV}} = \text{projector}(F_{\text{Fused}}^{\text{BEV}}) + F_{\text{Fused}}^{\text{BEV}}, \quad (15)$$

755 where $\text{projector}(\cdot)$ is the two-layer linear perceptron function. Note that, the BEV features of
 756 different modalities use a shared skip projector module.

756
757 Table 5: Description and severity level setups in camera/LiDAR corruption simulations. Camera
758 Crash (Camera), Frame Lost (Frame), Crosstalk, Incomplete Echo (Echo), and Cross-Sensor (Sensor).
759

Corruption	Description	Parameter	Easy	Moderate	Hard
Camera	dropping view images	number of dropped camera	2	4	5
Frame	dropping temporal frames	probability of frame dropping	2/6	4/6	5/6
Crosstalk	light impulses interference	percentage	0.03	0.07	0.12
Echo	imcomplete LiDAR readings	drop ratio	0.75	0.85	0.95
Sensor	cross sensor data	beam number to drop	8	16	20

Camera Corruption

LiDAR Corruption

765
766 Figure 8: Visualization results of camera/LiDAR sensor corruptions.
767
768
769
770
771

772 A.2 MULTI-SENSOR CORRUPTIONS

773 To explore the camera-LiDAR fusion model robustness, we design 13 types of camera-LiDAR
774 corruption combinations that perturb both camera and LiDAR inputs separately or concurrently.
775 Camera-LiDAR corruption combinations are grouped into camera-only corruptions, LiDAR-only
776 corruptions, and their combinations, covering the majority of real-world corruption cases. Specifically,
777 we design 3 types of camera-only corruptions by utilizing the clean LiDAR point data and three
778 camera failure cases such as Unavailable Camera (*all pixel values are set to zero for all RGB images*),
779 Camera Crash, and Frame Lost. Moreover, we design 4 types for LiDAR-only corruptions by utilizing
780 the clean camera data and the corrupted LiDAR data as the input. The LiDAR corruption types
781 include complete LiDAR failure which means LiDAR data are unavailable (*Since no model can*
782 *work when all points are absent, we approximate this scenario by only retaining a single point*
783 *as input*), LiDAR Incomplete Echo, LiDAR Crosstalk, and LiDAR Cross-Sensor. Note that our
784 implementation of complete LiDAR failure is close to the real-world situation. Lastly, we design 6
785 types of camera-LiDAR corruption combinations that perturb both sensor inputs concurrently, using
786 the previously mentioned image/LiDAR sensor failure types. We establish several corruption severity
787 levels (*i.e.*, three levels including easy, moderate, and hard) for each type of corruption. Furthermore,
788 for a comprehensive evaluation, we report metrics for each corruption type by averaging over three
789 severity levels. Description and severity level setups in 2 types of camera corruption and 3 types of
790 LiDAR corruption are shown in Tab. 5. Visualization results of camera/LiDAR sensor corruptions
791 are shown in Fig. 8.

792 A.3 RESULTS ON ARGOVERSE2 DATASET

793 There are 1000 logs in the Argoverse2 dataset Wilson et al. (2021). Each log contains 15s of 20Hz
794 RGB images from 7 cameras, 10Hz LiDAR sweeps, and a 3D vectorized map. The train, validation,
795 and test sets contain 700, 150, and 150 logs, respectively. Following previous works Liao et al.
796 (2023a); Zhou et al. (2024), we report results on its validation set and focus on the same three map
797 categories as the nuScenes dataset.

803 Tab. 8 and Tab. 9 show the overall performance of Uni-Map and all the baselines on the Argoverse2
804 dataset. Compared with MapTR, Uni-Map outperforms all input configurations in both single-class
805 APs and the overall mAP by a significant margin on the Argoverse2 dataset. Note that only one
806 Uni-Map model is trained while three MapTR models (MapTR-C, MapTR-L, and MapTR-F) are
807 trained for different input configurations. Thus, we use the total time of training three MapTR models
808 to train our Uni-Map model, and the resulting Uni-Map model (last row of Tab. 9) beats independently
809 trained camera-only, LiDAR-only, and camera-LiDAR fusion MapTR models with gains of 5.0,
4.8, 6.6 mAP, under the respective input configurations. In a nutshell, Uni-Map shows significant

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
17210
17211
17212
17213
17214
17215
17216
17217
17218
17219
17220
17221
17222
17223
17224
17225
17226
17227
17228
17229
17230
17231
17232
17233
17234
17235
17236
17237
17238
17239
17240
17241
17242
17243
17244
17245
17246
17247
17248
17249
17250
17251
17252
17253
17254
17255
17256
17257
17258
17259
17260
17261
17262
17263
17264
17265
17266
17267
17268
17269
17270
17271
17272
17273
17274
17275
17276
17277
17278
17279
17280
17281
17282
17283
17284
17285
17286
17287
17288
17289
17290
17291
17292
17293
17294
17295
17296
17297
17298
17299
172100
172101
172102
172103
172104
172105
172106
172107
172108
172109
172110
172111
172112
172113
172114
172115
172116
172117
172118
172119
172120
172121
172122
172123
172124
172125
172126
172127
172128
172129
172130
172131
172132
172133
172134
172135
172136
172137
172138
172139
172140
172141
172142
172143
172144
172145
172146
172147
172148
172149
172150
172151
172152
172153
172154
172155
172156
172157
172158
172159
172160
172161
172162
172163
172164
172165
172166
172167
172168
172169
172170
172171
172172
172173
172174
172175
172176
172177
172178
172179
172180
172181
172182
172183
172184
172185
172186
172187
172188
172189
172190
172191
172192
172193
172194
172195
172196
172197
172198
172199
172200
172201
172202
172203
172204
172205
172206
172207
172208
172209
172210
172211
172212
172213
172214
172215
172216
172217
172218
172219
172220
172221
172222
172223
172224
172225
172226
172227
172228
172229
172230
172231
172232
172233
172234
172235
172236
172237
172238
172239
172240
172241
172242
172243
172244
172245
172246
172247
172248
172249
172250
172251
172252
172253
172254
172255
172256
172257
172258
172259
172260
172261
172262
172263
172264
172265
172266
172267
172268
172269
172270
172271
172272
172273
172274
172275
172276
172277
172278
172279
172280
172281
172282
172283
172284
172285
172286
172287
172288
172289
172290
172291
172292
172293
172294
172295
172296
172297
172298
172299
172300
172301
172302
172303
172304
172305
172306
172307
172308
172309
172310
172311
172312
172313
172314
172315
172316
172317
172318
172319
172320
172321
172322
172323
172324
172325
172326
172327
172328
172329
172330
172331
172332
172333
172334
172335
172336
172337
172338
172339
172340
172341
172342
172343
172344
172345
172346
172347
172348
172349
172350
172351
172352
172353
172354
172355
172356
172357
172358
172359
172360
172361
172362
172363
172364
172365
172366
172367
172368
172369
172370
172371
172372
172373
172374
172375
172376
172377
172378
172379
172380
172381
172382
172383
172384
172385
172386
172387
172388
172389
172390
172391
172392
172393
172394
172395
172396
172397
172398
172399
172400
172401
172402
172403
172404
172405
172406
172407
172408
172409
172410
172411
172412
172413
172414
172415
172416
172417
172418
172419
172420
172421
172422
172423
172424
172425
172426
172427
172428
172429
172430
172431
172432
172433
172434
172435
172436
172437
172438
172439
172440
172441
172442
172443
172444
172445
172446
172447
172448
172449
172450
172451
172452
172453
172454
172455
172456
172457
172458
172459
172460
172461
172462
172463
172464
172465
172466
172467
172468
172469
172470
172471
172472
172473
172474
172475
172476
172477
172478
172479
172480
172481
172482
172483
172484
172485
172486
172487
172488
172489
172490
172491
172492
172493
172494
172495
172496
172497
172498
172499
172500
172501
172502
172503
172504
172505
172506
172507
172508
172509
172510
172511
172512
172513
172514
172515
172516
172517
172518
172519
172520
172521
172522
172523
172524
172525
172526
172527
172528
172529
172530
172531
172532
172533
172534
172535
172536
172537
172538
172539
172540
172541
172542
172543
172544
172545
172546
172547
172548
172549
172550
172551
172552
172553
172554
172555
172556
172557
172558
172559
172560
172561
172562
172563
172564
172565
172566
172567
172568
172569
172570
172571
172572
172573
172574
172575
172576
172577
172578
172579
172580
172581
172582
172583
172584
172585
172586
172587
172588
172589
172590
172591
172592
172593
172594
172595
172596
172597
172598
172599
172600
172601
172602
172603
172604
172605
172606
172607
172608
172609
172610
172611
172612
172613
172614
172615
172616
172617
172618
172619
172620
172621
172622
172623
172624
172625
172626
172627
172628
172629
172630
172631
172632
172633
172634
172635
172636
172637
172638
172639
172640
172641
172642
172643
172644
172645
172646
172647
172648
172649
172650
172651
172652
172653
172654
172655
172656
172657
172658
172659
172660
172661
172662
172663
172664
172665
172666
172667
172668
172669
172670
172671
172672
172673
172674
172675
172676
172677
172678
172679
172680
172681
172682
172683
172684
172685
172686
172687
172688
172689
172690
172691
172692
172693
172694
172695
172696
172697
172698
172699
172700
172701
172702
172703
172704
172705
172706
172707
172708
172709
172710
172711
172712
172713
172714
172715
172716
172717
172718
172719
172720
172721
172722
172723
172724
172725
172726
172727
172728
172729
172730
172731
172732
172733
172734
172735
172736
172737
172738
172739
172740
172741
172742
172743
172744
172745
172746
172747
172748
172749
172750
172751
172752
172753
172754
172755
172756
172757
172758
172759
172760
172761
172762
172763
172764
172765
172766
172767
172768
172769
1

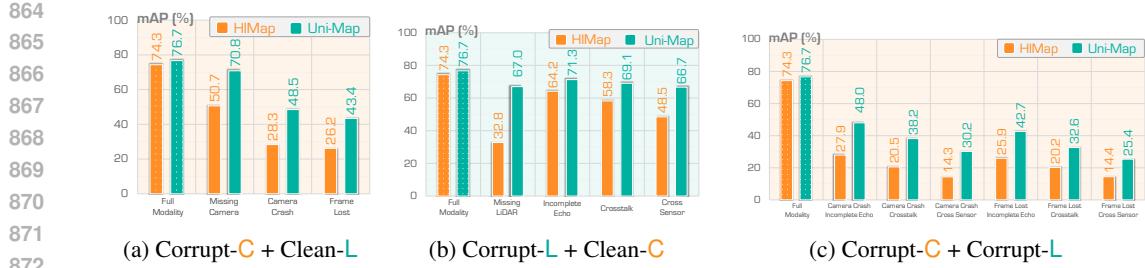


Figure 9: The result of multi-sensor corruption on HIMap vs. Uni-Map (HIMap) fusion model.

corruptions, and their combinations, covering the majority of real-world corruption cases. Fig. 9 shows the results of three Camera-LiDAR corruption combinations on HIMap Zhou et al. (2024) fusion model. We can find that: (1) In the sensor missing scenario, Uni-Map can still keep the model from collapsing based on our switching modality strategy. Quantitatively, when facing the camera sensor missing case, Uni-Map still achieves 70.8 mAP, which outperforms the original HIMap Zhou et al. (2024) by +20.1 mAP. (2) In case of corruption of camera and LiDAR sensor individually or simultaneously, Uni-Map still shows stronger robustness. For example, in the face of camera crash and LiDAR crosstalk, compared to the MapTR fused model, the Uni-Map model achieved significant improvements in 17.7 mAP (38.2 vs. 20.5). All in all, Uni-Map shows stronger robustness on our designed 13 types of camera-LiDAR corruption combinations. Experimental results for all corruption types for MapTR and Uni-Map (MapTR) are shown in Tab. 12-Tab. 14. And, experimental results for all corruption types for HIMap and Uni-Map (HIMap) are shown in Tab. 15-Tab. 17.

Table 11: Comparison of BEVFusion Liu et al. (2023b) and Uni-Map in terms of accuracy on the nuScenes dataset. Note that only one Uni-Map model is trained while three BEVFusion models (BEVFusion-C, BEVFusion-L and BEVFusion-F) are trained for different input configurations.

Method	Camera-only (mAP/NDS)	LiDAR-only (mAP/NDS)	Camera & LiDAR (mAP/NDS)
BEVFusion-C	35.6/41.2	—	—
BEVFusion-L	—	64.7/69.3	—
BEVFusion-F	—	—	68.5/71.4
Uni-Map (BEVFusion)	39.2/46.1	67.3/71.6	71.1/73.5

A.6 GENERALIZATION TO 3D OBJECT DETECTION TASK

In order to verify the universality of the Uni-Map method, we thereby generalize our method to the 3D object detection task, to further show its effectiveness on other perception tasks. We select the popular 3D object detection method BEVFusion Liu et al. (2023b) as the baseline model. As shown in the Tab. 11, our Uni-Map consistently improves the performance, compared to the original model. For example, our Uni-Map beats independently trained camera-only, LiDAR-only, and camera-LiDAR fusion models with gains of 3.6/4.9, 2.6/2.3, 2.6/2.1 mAP/NDS, under the respective input configurations. Obviously, our method can be directly utilized in the 3D objection detection task, demonstrating the generalization ability of our method.

A.7 MORE VISUALIZATION RESULTS

We provide more visualization results of qualitative results. Visualization results of qualitative results are shown in Fig. 10. We observe that in the case of multi-sensor corruption, the source MapTR model predictions are highly incorrect. However, our Uni-Map model can already correct significant errors in the baseline predictions in all settings. Qualitative results demonstrate the superiority of the UniMap model in various corruption scenarios.



972

973

974

Table 12: The result of camera-only corruptions on MapTR vs Uni-Map (MapTR) fusion model.

975

Method	Modality	Camera	LiDAR	AP _{ped.}	AP _{div.}	AP _{bou.}	mAP↑
MapTR Liao et al. (2023a)	C & L	✓	✓	55.9	62.3	69.3	62.5
MapTR Liao et al. (2023a)	C & L	✗	✓	15.0	18.2	34.4	22.5 _{-40.0}
MapTR Liao et al. (2023a)	C & L	Camera Crash	✓	32.5	36.5	48.4	39.1 _{-23.4}
MapTR Liao et al. (2023a)	C & L	Frame Lost	✓	29.1	33.7	46.1	36.3 _{-26.2}
Uni-Map (MapTR)	C & L	✓	✓	64.4	66.8	73.2	68.1
Uni-Map (MapTR)	C & L	✗	✓	56.5	57.8	69.4	61.2 _{-6.9}
Uni-Map (MapTR)	C & L	Camera Crash	✓	40.3	40.3	51.5	44.1 _{-24.0}
Uni-Map (MapTR)	C & L	Frame Lost	✓	37.0	38.6	49.9	41.8 _{-26.3}

981

982

983

984

985

986

987

988

989

990

Table 13: The result of LiDAR-only corruptions on MapTR vs Uni-Map (MapTR) fusion model.

991

Method	Modality	Camera	LiDAR	AP _{ped.}	AP _{div.}	AP _{bou.}	mAP↑
MapTR Liao et al. (2023a)	C & L	✓	✓	55.9	62.3	69.3	62.5
MapTR Liao et al. (2023a)	C & L	✓	✗	20.7	27.4	13.1	20.4 _{-42.1}
MapTR Liao et al. (2023a)	C & L	✓	Incomplete Echo	47.9	55.2	62.2	55.1 _{-7.4}
MapTR Liao et al. (2023a)	C & L	✓	Crosstalk	36.7	42.5	45.3	41.5 _{-21.0}
MapTR Liao et al. (2023a)	C & L	✓	Cross-Sensor	33.9	42.9	42.0	39.6 _{-22.9}
Uni-Map (MapTR)	C & L	✓	✓	64.4	66.8	73.2	68.1
Uni-Map (MapTR)	C & L	✓	✗	52.1	57.5	55.2	54.9 _{-13.2}
Uni-Map (MapTR)	C & L	✓	Incomplete Echo	56.5	61.3	65.9	61.2 _{-6.9}
Uni-Map (MapTR)	C & L	✓	Crosstalk	53.3	58.2	60.9	57.5 _{-10.6}
Uni-Map (MapTR)	C & L	✓	Cross-Sensor	50.5	55.4	57.2	54.3 _{-13.8}

1003

1004

1005

1006

1007

1008

1009

1010

Table 14: The result of camera and LiDAR corruptions on MapTR vs Uni-Map (MapTR) fusion model.

1011

Method	Modality	Camera	LiDAR	AP _{ped.}	AP _{div.}	AP _{bou.}	mAP↑
MapTR Liao et al. (2023a)	C & L	✓	✓	55.9	62.3	69.3	62.5
MapTR Liao et al. (2023a)	C & L	Camera Crash	Incomplete Echo	32.4	35.6	47.8	38.6 _{-23.9}
MapTR Liao et al. (2023a)	C & L	Camera Crash	Crosstalk	19.7	21.6	26.9	22.7 _{-39.8}
MapTR Liao et al. (2023a)	C & L	Camera Crash	Cross-Sensor	18.4	20.8	23.2	20.8 _{-41.7}
MapTR Liao et al. (2023a)	C & L	Frame Lost	Incomplete Echo	28.9	32.8	45.5	35.8 _{-26.7}
MapTR Liao et al. (2023a)	C & L	Frame Lost	Crosstalk	16.9	19.9	25.5	20.8 _{-41.7}
MapTR Liao et al. (2023a)	C & L	Frame Lost	Cross-Sensor	15.8	19.4	22.2	19.1 _{-43.4}
Uni-Map (MapTR)	C & L	✓	✓	64.4	66.8	73.2	68.1
Uni-Map (MapTR)	C & L	Camera Crash	Incomplete Echo	40.3	39.7	50.8	43.6 _{-24.5}
Uni-Map (MapTR)	C & L	Camera Crash	Crosstalk	29.8	28.7	36.4	31.6 _{-36.5}
Uni-Map (MapTR)	C & L	Camera Crash	Cross-Sensor	24.5	24.6	28.8	25.9 _{-42.2}
Uni-Map (MapTR)	C & L	Frame Lost	Incomplete Echo	36.9	37.8	49.2	41.3 _{-26.8}
Uni-Map (MapTR)	C & L	Frame Lost	Crosstalk	26.3	27.3	34.3	29.3 _{-38.8}
Uni-Map (MapTR)	C & L	Frame Lost	Cross-Sensor	20.9	23.3	26.6	23.6 _{-44.5}

1024

1025

1026

1027

1028

1029

Table 15: The result of camera-only corruptions on HIMap vs Uni-Map (HIMap) fusion model.

Method	Modality	Camera	LiDAR	AP _{ped.}	AP _{div.}	AP _{bou.}	mAP↑
HIMap Zhou et al. (2024)	C & L	✓	✓	71.0	72.4	79.4	74.3
HIMap Zhou et al. (2024)	C & L	✗	✓	40.9	46.4	74.7	50.7 _{-23.6}
HIMap Zhou et al. (2024)	C & L	Camera Crash	✓	36.3	27.7	20.9	28.3 _{-46.0}
HIMap Zhou et al. (2024)	C & L	Frame Lost	✓	29.9	25.0	23.8	26.2 _{-48.1}
Uni-Map (HIMap)	C & L	✓	✓	73.6	75.3	81.2	76.7
Uni-Map (HIMap)	C & L	✗	✓	65.3	69.5	77.8	70.8 _{-5.9}
Uni-Map (HIMap)	C & L	Camera Crash	✓	42.5	47.6	55.5	48.5 _{-28.2}
Uni-Map (HIMap)	C & L	Frame Lost	✓	36.7	42.3	51.1	43.4 _{-33.3}

1039

1040

1041

1042

1043

1044

Table 16: The result of LiDAR-only corruptions on HIMap vs Uni-Map (HIMap) fusion model.

Method	Modality	Camera	LiDAR	AP _{ped.}	AP _{div.}	AP _{bou.}	mAP↑
HIMap Zhou et al. (2024)	C & L	✓	✓	71.0	72.4	79.4	74.3
HIMap Zhou et al. (2024)	C & L	✓	✗	30.7	38.7	29.0	32.8 _{-41.5}
HIMap Zhou et al. (2024)	C & L	✓	Incomplete Echo	59.1	63.7	69.9	64.2 _{-10.1}
HIMap Zhou et al. (2024)	C & L	✓	Crosstalk	54.1	57.5	63.4	58.3 _{-16.0}
HIMap Zhou et al. (2024)	C & L	✓	Cross-Sensor	44.2	50.7	50.8	48.5 _{-25.8}
Uni-Map (HIMap)	C & L	✓	✓	73.6	75.3	81.2	76.7
Uni-Map (HIMap)	C & L	✓	✗	64.5	68.2	68.3	67.0 _{-9.7}
Uni-Map (HIMap)	C & L	✓	Incomplete Echo	68.0	70.8	75.0	71.3 _{-5.4}
Uni-Map (HIMap)	C & L	✓	Crosstalk	65.9	68.9	72.6	69.1 _{-7.6}
Uni-Map (HIMap)	C & L	✓	Cross-Sensor	63.8	67.4	69.1	66.7 ₋₁₀

1057

1058

1059

1060

1061

1062

Table 17: The result of camera and LiDAR corruptions on HIMap vs Uni-Map (HIMap) fusion model.

Method	Modality	Camera	LiDAR	AP _{ped.}	AP _{div.}	AP _{bou.}	mAP↑
HIMap Zhou et al. (2024)	C & L	✓	✓	71.0	72.4	79.4	74.3
HIMap Zhou et al. (2024)	C & L	Camera Crash	Incomplete Echo	36.2	26.9	20.5	27.9 _{-46.4}
HIMap Zhou et al. (2024)	C & L	Camera Crash	Crosstalk	29.2	19.3	12.9	20.5 _{-53.8}
HIMap Zhou et al. (2024)	C & L	Camera Crash	Cross-Sensor	23.1	13.8	5.9	14.3 _{-60.0}
HIMap Zhou et al. (2024)	C & L	Frame Lost	Incomplete Echo	29.9	24.4	23.5	25.9 _{-48.4}
HIMap Zhou et al. (2024)	C & L	Frame Lost	Crosstalk	23.6	18.9	18.0	20.2 _{-54.1}
HIMap Zhou et al. (2024)	C & L	Frame Lost	Cross-Sensor	17.7	14.3	11.2	14.4 _{-59.9}
Uni-Map (HIMap)	C & L	✓	✓	73.6	75.3	81.2	76.7
Uni-Map (HIMap)	C & L	Camera Crash	Incomplete Echo	42.4	46.7	54.8	48.0 _{-28.7}
Uni-Map (HIMap)	C & L	Camera Crash	Crosstalk	35.1	36.6	42.8	38.2 _{-38.5}
Uni-Map (HIMap)	C & L	Camera Crash	Cross-Sensor	28.9	30.8	31.0	30.2 _{-46.5}
Uni-Map (HIMap)	C & L	Frame Lost	Incomplete Echo	36.6	41.2	50.3	42.7 _{-34.0}
Uni-Map (HIMap)	C & L	Frame Lost	Crosstalk	29.2	31.3	37.5	32.6 _{-44.1}
Uni-Map (HIMap)	C & L	Frame Lost	Cross-Sensor	23.9	25.9	26.4	25.4 _{-51.3}

1078

1079