# SHIFT: A Synthetic Driving Dataset for Continuous Multi-Task Domain Adaptation

Tao Sun[1*]    Mattia Segu[1*]    Janis Postels[1]    Yuxuan Wang[1]
Luc Van Gool[1]    Bernt Schiele[2]    Federico Tombari[3,4]    Fisher Yu[1]

[1]ETH Zürich    [2]MPI Informatics    [3]Google    [4]Technical University of Munich

{taosun47, segum, jpostels, yuxuwang}@ethz.ch
vangool@vision.ee.ethz.ch, schiele@mpi-inf.mpg.de, tombari@in.tum.de, i@yf.io

## Abstract

*Adapting to a continuously evolving environment is a safety-critical challenge inevitably faced by all autonomous driving systems. Existing image and video driving datasets, however, fall short of capturing the mutable nature of the real world. In this paper, we introduce the largest multi-task synthetic dataset for autonomous driving, SHIFT. It presents discrete and continuous shifts in cloudiness, rain and fog intensity, time of day, and vehicle and pedestrian density. Featuring a comprehensive sensor suite and annotations for several mainstream perception tasks, SHIFT allows investigating the degradation of a perception system performance at increasing levels of domain shift, fostering the development of continuous adaptation strategies to mitigate this problem and assess model robustness and generality. Our dataset and benchmark toolkit are publicly available at* www.vis.xyz/shift.
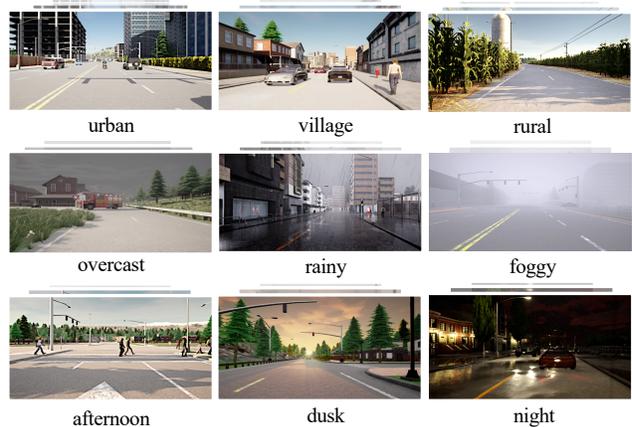
## 1. Introduction

Recent years have witnessed the remarkable progress of perception systems for autonomous driving. Betting on the role that autonomous driving will serve for society, industry and academia have joined forces to collect and release several large-scale driving datasets, raising hopes for a forthcoming successful deployment of self-driving cars.

Large-scale driving datasets have played a pivotal role in the prosperity of perception algorithms and provide a playground for different techniques to compete and thrive on multiple tasks. However, while the algorithm accuracy surges, progress in terms of generalization to unforeseen environmental conditions has been underwhelming [10, 43].

To achieve full autonomy, self-driving cars must adapt to new environments and identify life-threatening failure cases

---

*Equal contribution.

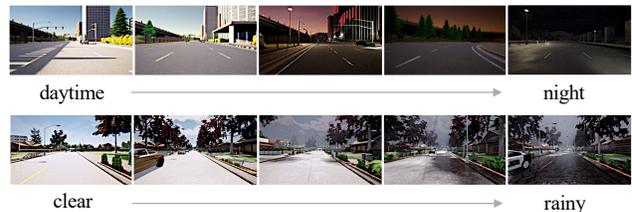**Discrete domain shifts**



**Continuous domain shifts**



Figure 1. SHIFT provides: (a) **discrete domain shifts**, a set of sequences each collected using different domain parameters and initial states; (b) **continuous domain shifts**, a set of sequences where domain parameters change continuously during driving.

to promptly prevent crashes. Examples of domain shifts affecting driving are changes in weather and lighting conditions, scenery, and behavior, appearance, and quantity of agents on the road. Domain shift [2] is a well-known problem for learning algorithms, causing unforeseeable performance drops under conditions different from the training ones. Techniques to prevent, counteract or assess its im-

1

pact have been developed in the form of domain generalization [28, 44, 76, 84], domain adaptation [14, 38, 80, 87], uncertainty estimation [13, 32, 39, 54] and out-of-distribution (OOD) detection [24, 53, 64, 89]. However, such approaches are typically deployed and tested on toy datasets [35, 65, 82] or synthetically corrupted ones [22]. Although there are preliminary attempts at providing driving datasets with different domains [5, 10, 51, 63, 68, 69, 79, 90], each only covers a limited amount of perception tasks (*e.g.* only semantic segmentation [68, 69]) and a narrow selection of domain shift directions (*e.g.* only rain [79] or snow [51]). Consequently, current solutions to domain shift cannot undergo scrutiny in controlled autonomous driving scenarios, making it difficult to verify their safety without risking real-world car crashes.

Given their short length, sequences from existing driving datasets are captured under approximately stationary conditions, and only *discrete shifts* are witnessed among sets of sequences presenting different homogeneous conditions from one set to another (*e.g.* clear weather and rainy). However, nothing in this world is constant except change and becoming. *Continuous shifts* - the intra-sequence shifts from one domain into another - are a certainty in the real world, where a sunny day can rapidly turn into a rainy one, or a quiet road can quickly become busy. Moreover, continuous distributional shift has recently been shown to represent a critical challenge for current learning systems [55].

An adequate dataset design is thus needed to quantify and address domain shift both at discrete and continuous levels. Consequently, we set the goal of overcoming the outdated paradigm of previous driving datasets and introduce SHIFT, a new synthetic dataset capturing the continuously evolving nature of the real world through realistic discrete and continuous shifts along safety-critical environmental directions: time of day, cloudiness, rain, fog strength, and vehicle and pedestrian density. Collected in the CARLA simulator [12], SHIFT includes a comprehensive sensor suite and covers the most important perception tasks. Counting 4,800+ sequences captured from a multi-view sensor suite in 8 different locations, our dataset supports 13 perception tasks for multi-task driving systems: semantic/instance segmentation, monocular/stereo depth regression, 2D/3D object detection, 2D/3D multiple object tracking (MOT), optical flow estimation, point cloud registration, visual odometry, trajectory forecasting and human pose estimation.

Our dataset aims to foster research in several under-explored fields related to the generality and reliability of perception systems for autonomous driving, *e.g.* domain generalization, domain adaptation, and uncertainty estimation. Moreover, by collecting incremental discrete shifts from one domain to another, we hope to foster research in the field of continual learning [18, 83, 86] for autonomous driving, so far only studied on discrete levels of synthetic corruptions [22] of traditional image classification datasets [11, 31]. Finally, by collecting sequences with realistic intra-sequence continuous domain shifts, we provide the first driving dataset allowing research on continuous test-time learning and adaptation [52, 73, 77, 78, 86].

The main contributions of this work are:

- We introduce SHIFT, a multi-task driving dataset featuring the most important perception tasks under a variety of conditions and with a comprehensive sensor setup. To the best of our knowledge, it is the largest synthetic dataset for autonomous driving and provides the most inclusive set of annotations and conditions.

- Using SHIFT, we analyze the importance of modeling discrete and continuous domain shifts, and demonstrate new findings on different adaptation and uncertainty estimation methods under continuous shifts.

## 2. Related Work

During the past decade, a large variety of realistic and synthetic driving datasets emerged, providing a playground for researchers to develop novel algorithms. Domain shift is a common threat to the performance and safety of learning-based methods.

We here introduce the most notable driving datasets and the techniques to mitigate the domain shift effect. For an overview of the current driving datasets, refer to Tab. 1.

**Real-world driving datasets** typically focus on a specific subset of perception tasks due to the high cost of data collection and annotation. After almost a decade of development, the pioneering real-world dataset KITTI [16] supports almost all the perception tasks for autonomous driving, including semantic / instance segmentation, depth estimation, 2D and 3D object detection and tracking, optical flow, scene flow, and visual odometry. However, its small scale represents an obvious problem and its diversity is severely limited compared to modern large-scale datasets. CamVid [4], Cityscapes [9], and Mapillary [46] are image-based driving datasets for segmentation, A*3D [50] for 3D object detection, and HD1K [30] for optical flow estimation. Recently, many large-scale datasets, *e.g.*, BDD100K [90], Waymo Open [72], H3D [48], and nuScenes [5], have been released with multi-task annotations, although mainly focusing on object detection and tracking. Our dataset offers a complete set of annotations for all the frames, comprehensive of all the most important perception tasks supported by other datasets, and enabling multi-task learning on a broader range of tasks and conditions.

**Synthetic driving datasets** are collected using graphic engines and physical simulators. SYNTHIA [63] contains images and segmentation annotations generated by its simulator. AIODrive [88] is produced using CARLA Simula-

2

| | Dataset | Cities | Tracking sequences | Max length for sequence | Domain shift† | Annotated frames for | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | Seg. | 2D Det. | 3D Det. | MOT | Depth | Flow | Pose◇ |
| Real-world | KITTI [16] | 1 | 22 | 106 sec | no | 200 | 15k | 15k | 15k | 93k | 389 | - |
| | CamVid [4] | 4 | - | - | no | 700 | - | - | - | - | - | - |
| | Cityscapes [9] | 27 | - | - | no | 25k | 25k | 25k | - | - | - | - |
| | Cityscapes-C‡ [43] | 27 | - | - | discrete | 25k | 25k | 25k | - | - | - | - |
| | H3D [48] | 4 | 160 | 20 sec | discrete | - | - | 27k | 27k | - | - | - |
| | HD1K [30] | 1 | - | - | discrete | - | - | - | - | - | 1k | - |
| | A*3D [49] | 1 | - | - | discrete | - | - | 39k | - | - | - | - |
| | nuScenes [5] | 2 | 1,000 | 20 sec | discrete | - | - | 40k | 40k | - | - | - |
| | Waymo Open [72] | 3 | 1,150 | 20 sec | discrete | - | 200k | 230k | 230k | - | - | 230k |
| | BDD100K [90] | -§ | 2,000 | 40 sec | discrete | 10k | 100k | - | 318k | - | - | - |
| Synthetic | SYNTHIA [63] | 3 | - | - | discrete | 9,000 | 200k | 200k | - | - | - | - |
| | GTA-V [61] | 1 | - | - | no | 25k | - | - | - | - | - | - |
| | VIPER [60] | 1 | 184 | 10 min | discrete | 320k | 320k | - | 320k | - | 320k | - |
| | AIODrive [88] | 8 | 100 | 100 sec | discrete | 100k | 100k | 100k | 100k | 100k | - | - |
| | **SHIFT (ours)** | 8 | 4,850 | 33 min | discrete + continuous | 2.5M | 2.5M | 2.5M | 2.5M | 2.5M | 2.5M | 2.5M |

Table 1. Comparison of size and supported tasks of existing driving datasets. SHIFT is the largest synthetic dataset and, most notably, the only dataset providing realistic continuous domain shifts, diverse annotations, and longer annotated sequences. † indicates whether the dataset presents domain annotations. ‡ artificially corrupted. § multiple cities; exact number not known. ◇ key points for human pose.

tor with multiple sensor support, focusing on high-density long-range LiDAR sets. Compared to ours, these datasets present sequences of limited length and are restricted to discrete domain labels (Tab. 1). Further, video games have also been used for data generation. GTA-V [26, 61] provides images and segmentation masks captured from a popular game. VIPER [60] extends GTA-V by providing optical flow masks and discrete environmental labels. However, low-level control of video game engines is hardly accessible, impeding fine-grained environmental control and the collection of continuous shifts.

**Adverse conditions datasets** support the evaluation of robustness under different OOD conditions. A recent work [40] collects meteorological and air temperature measurements under discrete real-world shifts. Image-based datasets, *e.g.* CIFAR10/100-C [43], ImageNet-R [21] and Cityscapes-C [22], have been generated by applying artificial corruptions such as blurring, additive Gaussian noise and addition of specific patterns on the original dataset. Though carefully designed, such ad-hoc corruptions cannot fully represent the challenges presented by visual shifts in the real world. To this end, recent driving datasets [5, 41, 49, 72, 90] provide manually labeled tags for various weather conditions, scene categories, and day periods. However, each only covers a limited amount of perception tasks (see Tab. 1) and a narrow selection of domain shift directions. Moreover, ad-hoc datasets have been collected for specific underrepresented domains, *e.g.* rain [27, 79], fog [67, 68, 74], night [10], snow [51]. However, domain tags remain coarse-grained and only certain tasks and domain shift directions are supported. Recently, the ACDC dataset [69] has been proposed, featuring images evenly

distributed between fog, nighttime, rain, and snow. However, it supports only semantic segmentation. Interestingly, the India Driving Dataset [81] is the only dataset to provide extremely busy roads as adverse conditions. Overall, BDD100K [90] is the large-scale real-world dataset presenting the largest diversity of perception driving tasks and discrete domain labels for the time of day and weather conditions. For this reason, we use it as a reference to validate empirical observations drawn from our dataset. Nevertheless, compared to our dataset, BDD100K only provides annotated images from single cameras, does not provide 3D bounding boxes and optical flow annotations, distribution of domains is highly imbalanced and the domain is stationary within each sequence. In contrast, our dataset provides a full sensor suite, annotations for multiple tasks, balanced domain distribution and sets of sequences with continuously changing time of day, weather conditions (cloudiness, rain and fog strength), and vehicle and pedestrians density.

**Unsupervised domain adaptation (UDA)** means simultaneously learning on a labeled source and an unlabeled target domain to find transferable features across domains. UDA is mainly achieved via feature-space alignment [56, 71], domain-consistent regularization [14, 15, 25] and minimization of surrogate functions of domain gaps [66, 85]. The discrete shifts provided in our dataset can be directly used for training and evaluating UDA approaches.

**Continual domain adaptation** aims at performing consecutive discrete adaptation steps from one domain to multiple others. Incremental domain adaptation (IncDA) is a subset of continual DA that requires the source data and assumes availability of intermediate domains where domain shifts
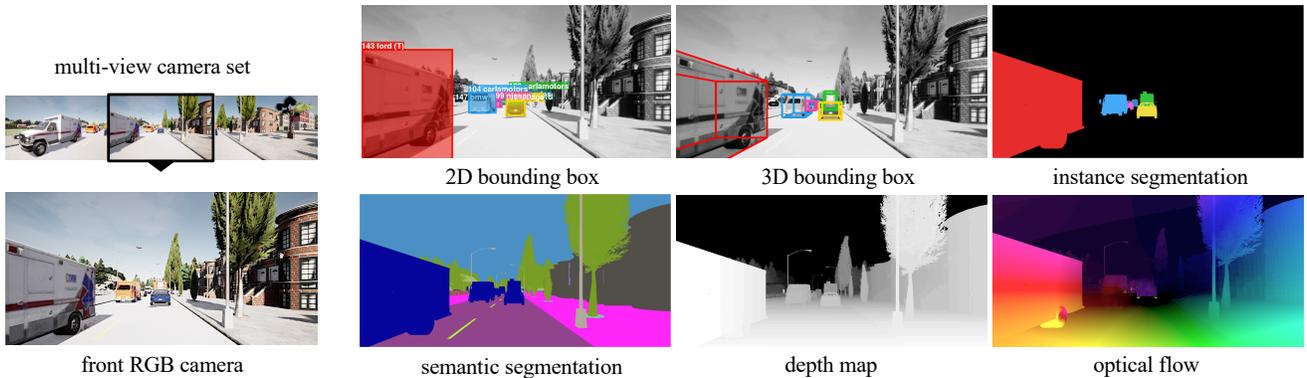
Figure 2. The annotation set of the RGB camera in our dataset. Each frame is associated with annotations of 2D/3D bounding boxes with tracking identities (visualized by different colors), semantic/instance segmentation, depth map and optical flow label.

occur gradually [33, 83, 86], allowing to minimize the gap between adaptation steps and performing adaptation from the source to the final target domain more effectively than with direct UDA. Providing different strengths of variations along natural axes, our dataset is suitable for IncDA.

**Continuous test-time adaptation** (ContinuousTTA) assumes that gradual domain shifts occur within the same test sequence, and adaptation is performed at test time on the incoming data stream. ContinuousTTA is a suitable choice for any scenario where a model is required to adapt on the go to a shifting domain and no large labeled or unlabeled collection of data from the target domain is available in advance. Recent works [45, 73, 86] show the efficiency of TTA when applied to artificial corruptions in the image-based datasets ImageNet-C/-R [21, 22]. The continuously shifting video sequences in our dataset provide instead realistic domain shift along natural directions, facilitating the development of ContinuousTTA methods transferable to the real world.

**Uncertainty Estimation** is a fundamental task for safety-critical vision applications. Quantifying the confidence about a model's prediction allows avoiding dangerous failures in autonomous driving. However, current uncertainty estimation techniques [13, 32, 36, 53] mainly focus on classification on toy datasets [31, 34], while recent work [55] has observed poor calibration, *i.e.* uncertainty uncorrelated with prediction's error, when such techniques are extended to more difficult datasets [23] and tasks under distributional shift. We hope that the domain shifts and multiple tasks supported in SHIFT will enable the study of uncertainty estimation methods on a wide variety of tasks for autonomous driving and their calibration under distributional shift.

## 3. The SHIFT Dataset

We provide a driving dataset with a comprehensive sensor suite (Sec. 3.1) and a rich set of annotations (Sec. 3.2), supporting multiple image- and video-based perception and

forecasting tasks against environmental changes. We detail our design choices regarding domain shifts in Sec. 3.3.

### 3.1. Sensor Suite

We collect the data through a comprehensive sensor suite. Our sensor suite features 11 different sensors, including a multi-view RGB camera set with 5 cameras, a stereo RGB camera set, an optical flow sensor, a depth camera, a GNSS sensor, and an IMU. All the cameras have a field-of-view of $90°$ and resolution of $1280 \times 800$ pixel. Moreover, we provide point clouds captured by a 128-channel LiDAR sensor. All sensors are synchronized and captured at a 10Hz rate. We follow the Scalabel [1] format and right-hand coordinate systems for storing all the annotations. More details are in the Appendix.

### 3.2. Annotations

We provide annotations for multiple mainstream perception tasks in autonomous driving, including 2D/3D bounding box trajectories, instance/semantic segmentation, optical flow and dense depth. Unlike real-world datasets, whose annotations are often limited to a group of keyframes due to prohibitive labeling cost, we offer full annotations for each frame in the sequences. More details are in the Appendix.

### 3.3. Dataset Design

Given their short sequence length, existing driving datasets are captured under approximately stationary conditions, and only discrete shifts are witnessed among sets of sequences presenting different homogeneous conditions (*e.g.* clear weather and rainy). We set the goal of overcoming the outdated paradigm of previous driving datasets and introduce SHIFT, a new synthetic dataset capturing the continuously evolving nature of the real world through realistic discrete and continuous shifts along safety-critical environmental directions: time of day, cloudiness, rain, fog strength, and vehicle and pedestrian density. We collect 5,250 sequences, of which 4,250 contain stationary environ-
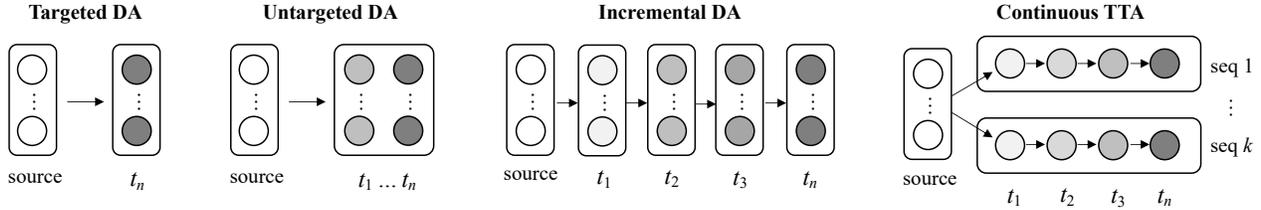
4

Figure 3. We evaluate four adaptation strategies: targeted domain adaptation (Targeted DA), untargeted domain adaptation (Untargeted DA), incremental domain adaptation (Incremental DA) and continuous test-time adaptation (Continuous TTA). The dots in the same row represent frames from the same sequence; their grayscale marks the degree of domain shift (white dots = source, dark gray dots = target.)
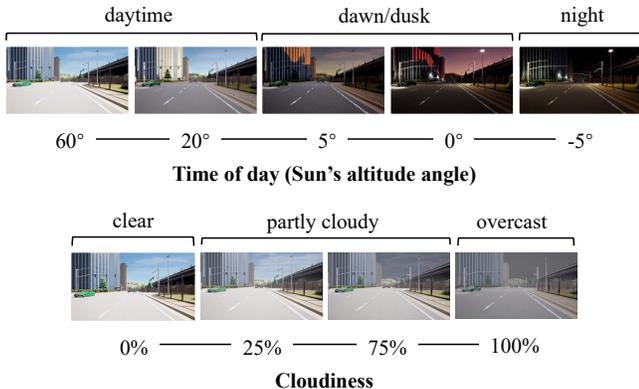


Figure 4. Examples of the two-level structure for domain labels. Each discrete label (tag above images) corresponds to an interval of continuous labels (*i.e.*, severity, axis below images).

mental conditions, *i.e.* inter-sequence domain shift. Each sequence is composed of 500 frames collected at 10 Hz, equivalent to 50 seconds of driving time. The remaining 600 sequences have continuously shifting conditions, *i.e.* inter-sequence domain shift. Totalling 70+ hours of driving and 2,500,000 annotated frames, SHIFT is the largest synthetic driving dataset available.

**Domain shift types.** We consider the most-frequent real-world environmental changes. SHIFT provides domain shifts in (a) weather conditions, including cloudiness, rain, and fog intensity, (b) time of day, (c) the density of vehicles and pedestrians, and (d) camera orientation.

**Domain shifts level.** To facilitate research on domain adaptation in different scenarios, SHIFT provides two levels of domain shifts, namely discrete domain shifts and continuous domain shifts. The *discrete* set contains 4,250 sequences generated with fixed environmental parameters and random initial states. We group these sequences into different domains, according to their severity. Fig. 4 shows grouping examples. All possible domain combinations are uniformly distributed across all sequences. The *continuous* set contains additional 600 sequences with continuous domain variations. In particular, each sequence presents a gradual shift from one domain to another, where the shift

happens through the intermediate domains that would naturally occur in the real world. In total, we collect 500 sequences of a basic 40 seconds length (1x), 80 sequences 10x longer than the basic length, and 20 100x longer. Each set is uniformly divided among the following shifts, each of which also loops back to the source domain: day → night, clear → rain, clear → foggy, clear → overcast. Given a domain shift direction, *e.g.* day to night, all other domain parameters are uniformly distributed across all sequences. Different sequence lengths allow analyzing the impact of domain shift speed on continuous TTA strategies (Sec. 4.2).

## 4. Experiments

SHIFT allows studying the robustness of perception systems for driving under both discrete and continuous distributional shifts. We first (Sec. 4.1) assess the impact of discrete domain shifts on model performance for multiple perception tasks available in our dataset and empirically demonstrate that observations from our simulation dataset transfer to real-world datasets. Moreover, we compare different discrete adaptation strategies and assess the calibration of uncertainty estimation methods under domain shifts. In Sec. 4.2 we extend the analysis to continuous domain shifts and investigate properties of continuous domain adaptation methods [86] against incremental adaptation and unsupervised domain adaptation [85]. Further experiments, implementation details, and ablations on the data collection choices are reported in the Appendix, together with additional experiments on multitask learning.

**Domain adaptation strategies.** To analyze the impact of our dataset design choice, we examine the four domain adaptation strategies allowed by our dataset (Fig. 3). As *Baseline*, we consider the model trained on the source domain only and directly tested on the other domains. *Targeted DA* [87] is a traditional computer vision problem consisting of adapting from a labeled source domain to a specific unlabeled target domain. We define *Untargeted DA* [35, 70] as adapting from a labeled source domain to a set of various unlabeled shifted domains. *Incremental DA* [83] consists in performing incremental steps of targeted
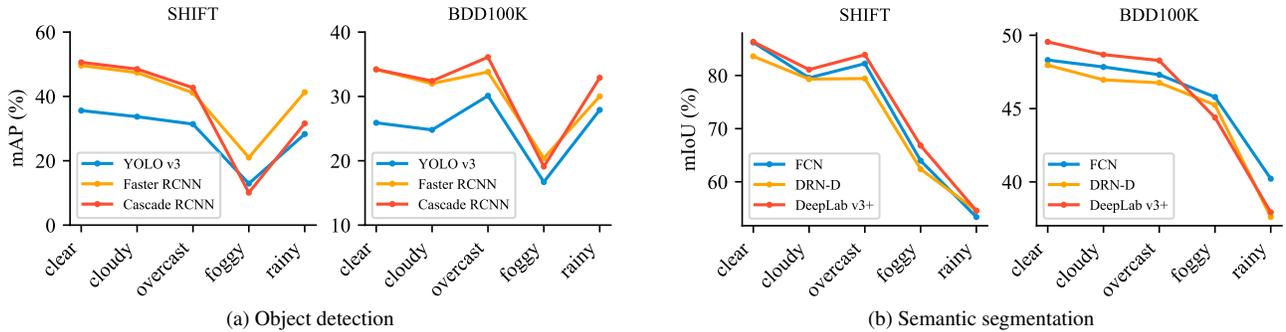
Figure 5. Performance degradation for different object detection (left) and semantic segmentation (right) methods under different weather conditions. Every model is trained under clear weather conditions and tested on other domains. SHIFT shows a similar trend as BDD100K.

| Task | Method | Metric | clear-daytime | partly cloudy | overcast | foggy | rainy | dawn/dusk | night |
|------|--------|--------|---------------|---------------|----------|-------|-------|-----------|-------|
| Semantic segmentation | DRN-D [91] | mIoU (%) ↑ | 83.6 | 79.3 | 79.4 | 62.4 | 54.6 | 60.8 | 42.8 |
| Instance segmentation | Mask R-CNN [19] | mAP (%) ↑ | 39.3 | 39.4 | 34.0 | 18.7 | 35.0 | 30.7 | 13.1 |
| Object detection | Faster R-CNN [6] | mAP (%) ↑ | 46.9 | 47.4 | 41.1 | 21.0 | 41.3 | 37.3 | 15.4 |
| MOT | QDTrack [47] | MOTA (%) ↑ | 56.2 | 53.4 | 46.2 | 25.0 | 41.9 | 44.7 | 16.5 |
| Mono. depth estimation | AdaBins-UNet [3] | SILog ↓ | 9.6 | 10.0 | 8.9 | 12.0 | 10.3 | 19.7 | 27.9 |
| Optical flow estimation | RAFT [75] | EPE (px) ↓ | 2.26 | 2.01 | 2.35 | 2.60 | 2.43 | 4.17 | 8.85 |

Table 2. Performance degradation on SHIFT of different methods for different perception tasks under discrete domain shifts. Training domain is underlined. The test domains are weather variations in daytime (partly cloudy, overcast, foggy, rainy) and time of day variations in clear weather (dawn/dusk, night). ↑ (↓): the higher (lower) the better.

| Scenario | Baseline | Targeted DA | Incremental DA |
|----------|----------|-------------|----------------|
| daytime → night | 42.8 | 45.3 | **47.3** |
| clear → foggy | **62.4** | 59.1 | 57.3 |
| clear → rainy | 54.6 | 61.0 | **64.9** |

Table 3. Comparison of different adaptation strategies for semantic segmentation under three directions of domain shift. The source domain is underlined. Incremental DA improves over Targeted DA, except for the case when Targeted DA underperforms the baseline. (Baseline = without DA)

DA from the source domain to the target domain passing through intermediate discretely-shifted domains. *Continuous TTA* [86] aims at adapting frame by frame to a sequence presenting a continuously shifted domain from source to target domain.

**Implementation details.** For the adaptation tasks, we focus on semantic segmentation and use ADVENT [85] for the Targeted and Untargeted DA. The segmentation backbone is DRN-D-54 [92]. Incremental DA is performed as a series of Targeted DA steps, while for Continuous TTA we extend TENT [86] to semantic segmentation and iteratively apply it on every incoming frame. Every model is trained in the clear-daytime domain and tested under different weather domains. While our dataset provides finer domain labels depending on the severity of the perturbation, we group different degrees of severity to match the environmental labels in BDD100K [90] in order to assess the compatibility of conclusions drawn from our dataset with real-world trends.

## 4.1. Discrete Shifts

As outlined in Sec. 3.3, our dataset provides incremental discrete shifts along natural environmental directions. We investigate properties of discrete shifts on the multitude of supported tasks and report findings on domain adaptation and uncertainty estimation performance.

**Impact of domain shift.** We find that many mainstream algorithms for different perception tasks suffer performance drops under domain shift (Tab. 2), where the severity increases with the distance from the source domain. In particular, we train all models in the clear-daytime domain and test under different weather conditions, showing the overall negative impact of domain shift on all the vision tasks supported by our dataset. Nevertheless, in some specific cases a model may even perform better on a shifted domain, *e.g.* instance segmentation on overcast. We leverage the incremental domain shifts provided in our dataset to investigate in Tab. 3 different discrete adaptation strategies for semantic segmentation, *i.e.* Incremental DA and Targeted DA. We find that incrementally adapting from source to target domain improves the generalization to the target domain compared to direct Targeted DA. However, clear → foggy represents a challenging scenario for which both the adaptation strategies worsen the baseline performance.

**Real-world compatibility.** To establish a reliable benchmark we must first confirm that trends witnessed in our simulation dataset are compatible with real-world observations. We use BDD100K [90] for comparison because it

6

| | | clear-daytime | cloudy | overcast | foggy | rainy | dawn/dusk | night | OOD avg. |
|---|---|---|---|---|---|---|---|---|---|
| **SHIFT** | Softmax | 3.3 | 32.6 | 14.2 | 48.8 | 64.3 | 43.7 | 64.7 | 45.2 |
| | MCDO | 1.2 | 13.1 | 7.6 | 20.8 | 10.0 | 27.2 | 39.6 | 19.7 |
| | Ensemble | 1.4 | 12.3 | 7.5 | 23.4 | 8.9 | 18.7 | 36.9 | 18.0 |
| **BDD** | Softmax | 9.6 | 23.2 | 9.9 | 9.7 | 7.7 | 10.6 | 48.6 | 18.4 |
| | MCDO | 12.3 | 22.0 | 7.8 | 13.0 | 11.4 | 13.1 | 41.4 | 18.1 |
| | Ensemble | 12.6 | 18.8 | 9.2 | 11.7 | 11.8 | 13.9 | 39.8 | 17.5 |

Table 4. Calibration (ECE, %) of uncertainty estimation methods under distributional shift for semantic segmentation. The lower, the better. Source domain is clear-daytime. We find that calibration worsens far from the source, both for SHIFT and BDD100K.



Figure 6. Comparison of different adaptation strategies for semantic segmentation on daytime → night shifts at varying amounts of available sequences. TTA is the most effective under limited amounts of data. When enough data becomes available, Incremental DA outperforms all other alternatives.



Figure 7. Performance on the target domain of TTA for different sequence lengths. Best learning rate on target domain is highlighted by black boxes. Both source and target performance are highly sensitive to the learning rates. Dashed lines = before TTA.

features the largest subset of our tasks available in a real-world dataset with discrete domain labels. We study the domain shift effect on two fundamental perception tasks, *i.e.* 2D object detection and semantic segmentation, and show compatible trends for different methods trained on SHIFT and BDD100K (Fig. 5). We evaluate the one-stage method YOLO v3 [58], as well as the two-stage methods Faster R-CNN [59] and Cascade R-CNN [6] for object detection. For semantic segmentation, we consider three different methods, FCN [37], DRN-D [91], and DeepLab v3+ [7]. Our experiments suggest that the performance of different methods for semantic segmentation and object detection degrades under different domain shifts. Moreover, we find that the ranking of methods and the relative degradation trend is compatible between SHIFT and the real-world dataset BDD100K, confirming the usefulness of SHIFT and its consistency with the real world.

**Uncertainty estimation.** Autonomous driving systems must deal with life-threatening failure cases. To this end, uncertainty estimation represents a powerful tool to assess the reliability of a model's predictions. Following [17], we evaluate the Expected Calibration Error (ECE) to assess the calibration, *i.e.* correlation with model error, of uncertainty estimation methods under domain shift. In particular, we evaluate the Softmax Entropy baseline and traditional Bayesian techniques such as Monte-Carlo Dropout (MCDO) [13] and Deep Ensembles [32]. We observe that such uncertainty estimation methods are not well calibrated under domain shift, and that calibration worsens under incremental shifts on both SHIFT and BDD100K (Tab. 4). While some domains are more challenging in SHIFT than in BDD100K, the overall degradation of calibration observed on SHIFT is confirmed on BDD100K and the ranking of methods is preserved, further highlighting that conclusions drawn from our dataset transfer to the real world.

We hope that our dataset will help researchers providing solutions to the potentially life-threatening shortcomings of current DA and uncertainty estimation techniques.
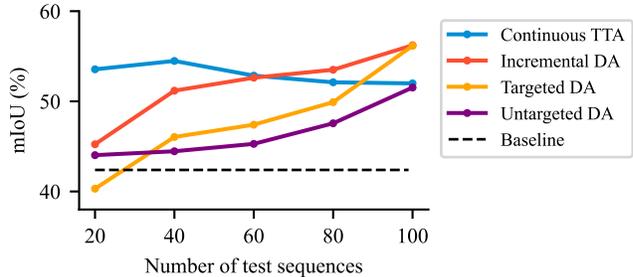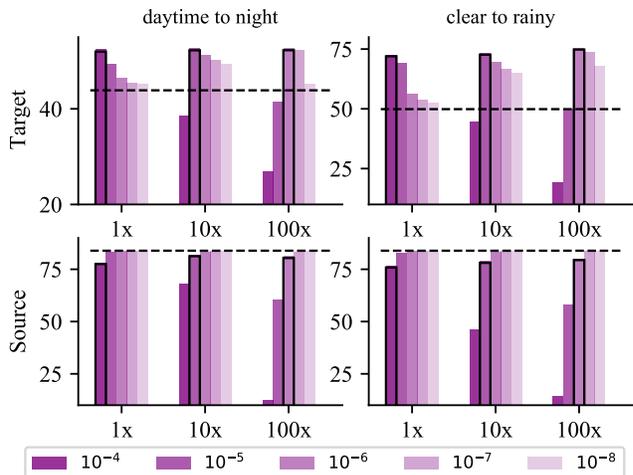
## 4.2. Continuous Shifts

A key feature of SHIFT is that of providing a set with continuous intra-sequence domain shifts, allowing to compare different adaptation strategies under continuous shifts and provide an in-depth analysis on TTA and its properties.

**Continual domain adaptation.** Fig. 6 compares four different adaptation strategies for semantic segmentation on an increasing number of sequences. Given a model pretrained on the source domain, *i.e.* clear-daytime, and the set of continuously shifting sequences from one domain to another, *i.e.* clear-daytime → night, we train the TTA algorithm on each frame of the incoming data stream. TTA is thus performed independently on each sequence. Final performance is averaged over all the sequences. For the other adaptation strategies, we divide the length of the sequence in 20 bins, consider each bin as a separate domain, and group corresponding bins from all the provided sequences. For
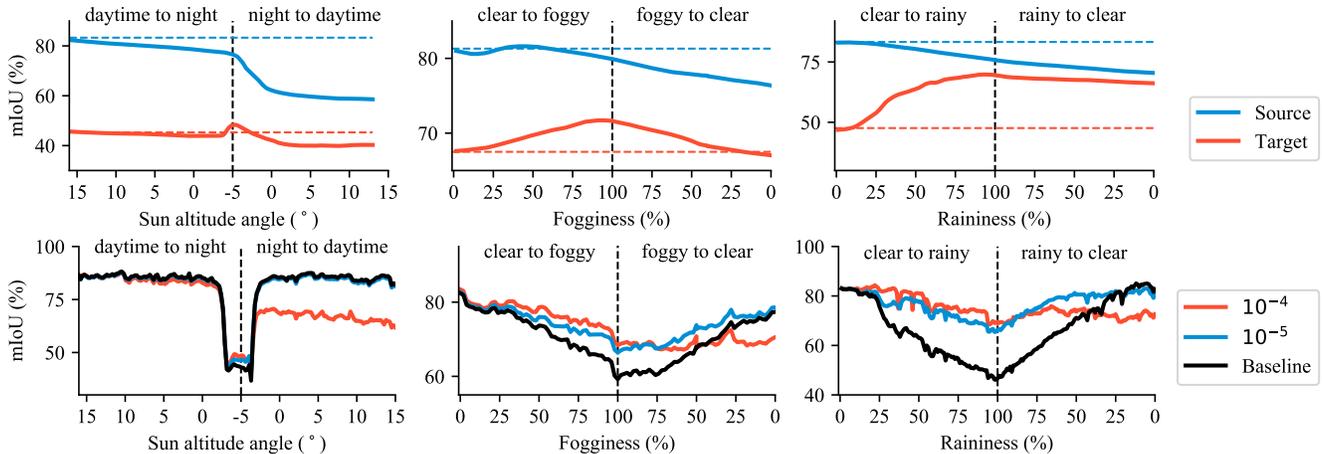
Figure 8. Performance of TTA for semantic segmentation under three types of domain shift: daytime → night, clear → foggy, clear → rainy. Each point corresponds to the performance of the model on the source (top-blue) / target (top-red) / current (bottom) domain finetuned up to that level of domain shift in the sequence. Horizontal lines in the bottom figure represent the original performance on source (blue) and target domain (red). After reaching the target domain, every sequence loops back to the original source domain. Catastrophic forgetting can be observed by the drop in source performance during TTA.

Targeted DA, we thus adapt directly to the last bin, corresponding to the night domain. Untargeted DA is instead applied on all the bins but the source one. Incremental DA is performed by incrementally adapting from one bin to the consecutive one until the end of the sequence is reached. In particular, we plot the average mIoU against the number of training sequences (Fig. 6). We find that TTA is extremely efficient under small target data availability compared to all other alternatives, and that Incremental DA is consistently more effective than both Targeted and Untargeted DA.

**Test-time adaptation.** As intra-sequence continuous shifts represent one of the main contributions of SHIFT, we further focus on TTA by using TENT [86] and evaluate the effect of the speed at which domain shift happens within a sequence on TTA performance (Fig. 7). This is made possible by the sets of sequences of different lengths (1x, 10x, 100x the basic sequence length).

Given a source and a target domain, *e.g.* daytime and night, each sequence starts from the source domain and reaches the target domain at mid-sequence length; then, it loops back to the original domain. We first observe that, depending on the domain shift speed, the learning rate can highly affect the outcome of the TTA (Fig. 7). Slower (faster) shifts will require lower (higher) learning rates. Moreover, after reaching the target domain at mid-sequence, the performance on the target domain has improved compared to its original value, while that on the source domain has dropped. According to Fig. 7 (1x), we find that the optimal learning rate in terms of adaptation to the target domain leads to the largest performance drop on the original source (Fig. 8, top). This problem, known as catastrophic forgetting [29] in the continual learn-

ing literature, has already been observed for class- and task-incremental learning.

To further investigate this issue, we loop back to the original domain after adapting to the target and find that, while the performance on the current target domains largely improves over the baseline (Fig. 8, bottom), the original source domain accuracy cannot be recovered (Fig. 8, top). While TTA has shown to be extremely effective to adapt on the go, a model adapted with TTA cannot be safely deployed on the original source domain. Showing that catastrophic forgetting also affects test-time adaptation further demonstrates the importance of providing continuously shifted sequences in driving datasets, and we hope that future research will attempt to mitigate this problem.

## 5. Conclusion

We introduce SHIFT, a multi-task driving dataset featuring the most important perception tasks under discrete and continuous domain shifts. Thanks to our dataset design, we demonstrate several new findings on different adaptation strategies and uncertainty estimation methods. Although simulation environments are still far from being a perfect representation of the real world, they allow inexpensive data collection and annotation. Moreover, we empirically demonstrate that conclusions drawn from our dataset hold in real-world datasets. To the best of our knowledge, SHIFT is the largest synthetic dataset for autonomous driving, providing the most inclusive set of annotations and conditions. We hope that providing the first dataset with realistic continuous domain shifts will contribute to shaping the data collection paradigm for real-world driving datasets and promote advances in test-time learning and adaptation.

# References

[1] Scalabel: A scalable open-source web annotation tool. https://scalabel.ai/. Accessed: 2021-11-16. 4

[2] Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 79(1):151–175, 2010. 1

[3] Shariq Farooq Bhat, Ibraheem Alhashim, and Peter Wonka. Adabins: Depth estimation using adaptive bins. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4009–4018, 2021. 6, 19

[4] Gabriel J Brostow, Julien Fauqueur, and Roberto Cipolla. Semantic object classes in video: A high-definition ground truth database. *Pattern Recognition Letters*, 30(2):88–97, 2009. 2, 3

[5] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multi-modal dataset for autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11621–11631, 2020. 2, 3, 13

[6] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6154–6162, 2018. 6, 7, 19

[7] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017. 7

[8] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *ECCV*, 2018. 19

[9] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016. 2, 3, 13

[10] Dengxin Dai and Luc Van Gool. Dark model adaptation: Semantic image segmentation from daytime to nighttime. In *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*, pages 3819–3824. IEEE, 2018. 1, 2, 3

[11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 2

[12] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. CARLA: An open urban driving simulator. In *Proceedings of the 1st Annual Conference on Robot Learning*, pages 1–16, 2017. 2

[13] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059, 2016. 2, 4, 7

[14] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *International conference on machine learning*, pages 1180–1189. PMLR, 2015. 2, 3

[15] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *The journal of machine learning research*, 17(1):2096–2030, 2016. 3

[16] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013. 2, 3, 13, 14, 19

[17] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International Conference on Machine Learning*, pages 1321–1330. PMLR, 2017. 7

[18] Guy Hacohen and Daphna Weinshall. On the power of curriculum learning in training deep networks. In *International Conference on Machine Learning*, pages 2535–2544. PMLR, 2019. 2

[19] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 6, 16, 19

[20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 19

[21] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8340–8349, 2021. 3, 4

[22] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *Proceedings of the International Conference on Learning Representations*, 2019. 2, 3, 4

[23] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019. 4

[24] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016. 2

[25] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei Efros, and Trevor Darrell. Cycada: Cycle-consistent adversarial domain adaptation. In *International conference on machine learning*, pages 1989–1998. PMLR, 2018. 3

[26] Hou-Ning Hu, Yung-Hsu Yang, Tobias Fischer, Trevor Darrell, Fisher Yu, and Min Sun. Monocular quasi-dense 3d object tracking. *arXiv preprint arXiv:2103.07351*, 2021. 3

[27] Jiongchao Jin, Arezou Fatemi, Wallace Lira, Fenggen Yu, Biao Leng, Rui Ma, Ali Mahdavi-Amiri, and Hao Zhang. RaidaR: A rich annotated image dataset of rainy street scenes. *arXiv preprint arXiv:2104.04606*, 2021. 3

[28] Aditya Khosla, Tinghui Zhou, Tomasz Malisiewicz, Alexei A Efros, and Antonio Torralba. Undoing the damage of dataset bias. In *European Conference on Computer Vision*, pages 158–171. Springer, 2012. 2

[29] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017. 8

[30] Daniel Kondermann, Rahul Nair, Katrin Honauer, Karsten Krispin, Jonas Andrulis, Alexander Brock, Burkhard Gussefeld, Mohsen Rahimimoghaddam, Sabine Hofmann, Claus Brenner, et al. The hci benchmark suite: Stereo and flow ground truth with uncertainties for urban autonomous driving. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 19–28, 2016. 2, 3

[31] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. The cifar-10 dataset. *online: http://www. cs. toronto. edu/kriz/cifar. html*, 55:5, 2014. 2, 4

[32] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in neural information processing systems*, pages 6402–6413, 2017. 2, 4, 7

[33] Qicheng Lao, Xiang Jiang, Mohammad Havaei, and Yoshua Bengio. Continuous domain adaptation with variational domain-agnostic feature replay. *arXiv preprint arXiv:2003.04382*, 2020. 4

[34] Yann LeCun. The mnist database of handwritten digits. *http://yann. lecun. com/exdb/mnist/*, 1998. 4

[35] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Deeper, broader and artier domain generalization. In *Proceedings of the IEEE international conference on computer vision*, pages 5542–5550, 2017. 2, 5

[36] Jeremiah Zhe Liu, Zi Lin, Shreyas Padhy, Dustin Tran, Tania Bedrax-Weiss, and Balaji Lakshminarayanan. Simple and principled uncertainty estimation with deterministic deep learning via distance awareness. *Conference on Neural Information Processing Systems*, 2020. 4

[37] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015. 7, 19

[38] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael Jordan. Learning transferable features with deep adaptation networks. In *International conference on machine learning*, pages 97–105. PMLR, 2015. 2

[39] Antonio Loquercio, Mattia Segu, and Davide Scaramuzza. A general framework for uncertainty estimation in deep learning. *IEEE Robotics and Automation Letters*, 5(2):3153–3160, 2020. 2

[40] Andrey Malinin, Neil Band, German Chesnokov, Yarin Gal, Mark JF Gales, Alexey Noskov, Andrey Ploskonosov, Liudmila Prokhorenkova, Ivan Provilkov, Vatsal Raina, et al. Shifts: A dataset of real distributional shift across multiple large-scale tasks. *arXiv preprint arXiv:2107.07455*, 2021. 3

[41] Jiageng Mao, Minzhe Niu, Chenhan Jiang, Xiaodan Liang, Yamin Li, Chaoqiang Ye, Wei Zhang, Zhenguo Li, Jie Yu, Chunjing Xu, et al. One million scenes for autonomous driving: Once dataset. 2021. 3

[42] Nikolaus Mayer, Eddy Ilg, Philip Hausser, Philipp Fischer, Daniel Cremers, Alexey Dosovitskiy, and Thomas Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4040–4048, 2016. 19

[43] Claudio Michaelis, Benjamin Mitzkus, Robert Geirhos, Evgenia Rusak, Oliver Bringmann, Alexander S Ecker, Matthias Bethge, and Wieland Brendel. Benchmarking robustness in object detection: Autonomous driving when winter is coming. *arXiv preprint arXiv:1907.07484*, 2019. 1, 3

[44] Krikamol Muandet, David Balduzzi, and Bernhard Schölkopf. Domain generalization via invariant feature representation. In *International Conference on Machine Learning*, pages 10–18. PMLR, 2013. 2

[45] Chaithanya Kumar Mummadi, Robin Hutmacher, Kilian Rambach, Evgeny Levinkov, Thomas Brox, and Jan Hendrik Metzen. Test-time adaptation to distribution shift by confidence maximization and input transformation. *arXiv preprint arXiv:2106.14999*, 2021. 4

[46] Gerhard Neuhold, Tobias Ollmann, Samuel Rota Bulo, and Peter Kontschieder. The mapillary vistas dataset for semantic understanding of street scenes. In *Proceedings of the IEEE international conference on computer vision*, pages 4990–4999, 2017. 2

[47] Jiangmiao Pang, Linlu Qiu, Xia Li, Haofeng Chen, Qi Li, Trevor Darrell, and Fisher Yu. Quasi-dense similarity learning for multiple object tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 164–173, 2021. 6

[48] Abhishek Patil, Srikanth Malla, Haiming Gang, and Yi-Ting Chen. The h3d dataset for full-surround 3d multi-object detection and tracking in crowded urban scenes. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 9552–9557. IEEE, 2019. 2, 3

[49] Quang-Hieu Pham, Pierre Sevestre, Ramanpreet Singh Pahwa, Huijing Zhan, Chun Ho Pang, Yuda Chen, Armin Mustafa, Vijay Chandrasekhar, and Jie Lin. A* 3d dataset: Towards autonomous driving in challenging environments. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2267–2273. IEEE, 2020. 3

[50] Quang-Hieu Pham, Pierre Sevestre, Ramanpreet Singh Pahwa, Huijing Zhan, Chun Ho Pang, Yuda Chen, Armin Mustafa, Vijay Chandrasekhar, and Jie Lin. A* 3d dataset: Towards autonomous driving in challenging environments. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2267–2273. IEEE, 2020. 2

[51] Matthew Pitropov, Danson Evan Garcia, Jason Rebello, Michael Smart, Carlos Wang, Krzysztof Czarnecki, and Steven Waslander. Canadian adverse driving conditions dataset. *The International Journal of Robotics Research*, 40(4-5):681–690, 2021. 2, 3

[52] Matteo Poggi, Alessio Tonioni, Fabio Tosi, Stefano Mattoccia, and Luigi Di Stefano. Continual adaptation for deep

stereo. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. 2

[53] Janis Postels, Hermann Blum, Yannick Strümpler, Cesar Cadena, Roland Siegwart, Luc Van Gool, and Federico Tombari. The hidden uncertainty in a neural networks activations. *arXiv preprint arXiv:2012.03082*, 2020. 2, 4

[54] Janis Postels, Francesco Ferroni, Huseyin Coskun, Nassir Navab, and Federico Tombari. Sampling-free epistemic uncertainty estimation using approximated variance propagation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2931–2940, 2019. 2

[55] Janis Postels, Mattia Segu, Tao Sun, Luc Van Gool, Fisher Yu, and Federico Tombari. On the practicality of deterministic epistemic uncertainty. *arXiv preprint arXiv:2107.00649*, 2021. 2, 4

[56] Joaquin Quiñonero-Candela, Masashi Sugiyama, Anton Schwaighofer, and N Lawrence. Covariate shift and local learning by distribution matching, 2008. 3

[57] Joseph Redmon and Ali Farhadi. Yolo9000: better, faster, stronger. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7263–7271, 2017. 19

[58] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018. 7, 19

[59] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28:91–99, 2015. 7, 18, 19

[60] Stephan R. Richter, Zeeshan Hayder, and Vladlen Koltun. Playing for benchmarks. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 2232–2241, 2017. 3, 13, 19

[61] Stephan R Richter, Vibhav Vineet, Stefan Roth, and Vladlen Koltun. Playing for data: Ground truth from computer games. In *European conference on computer vision*, pages 102–118. Springer, 2016. 3

[62] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 16, 19

[63] German Ros, Laura Sellart, Joanna Materzynska, David Vazquez, and Antonio M Lopez. The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3234–3243, 2016. 2, 3

[64] Lukas Ruff, Jacob R Kauffmann, Robert A Vandermeulen, Grégoire Montavon, Wojciech Samek, Marius Kloft, Thomas G Dietterich, and Klaus-Robert Müller. A unifying review of deep and shallow anomaly detection. *Proceedings of the IEEE*, 2021. 2

[65] Kate Saenko, Brian Kulis, Mario Fritz, and Trevor Darrell. Adapting visual category models to new domains. In *European conference on computer vision*, pages 213–226. Springer, 2010. 2

[66] Kuniaki Saito, Donghyun Kim, Stan Sclaroff, Trevor Darrell, and Kate Saenko. Semi-supervised domain adaptation via minimax entropy. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8050–8058, 2019. 3

[67] Christos Sakaridis, Dengxin Dai, Simon Hecker, and Luc Van Gool. Model adaptation with synthetic and real data for semantic dense foggy scene understanding. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 687–704, 2018. 3

[68] Christos Sakaridis, Dengxin Dai, and Luc Van Gool. Semantic foggy scene understanding with synthetic data. *International Journal of Computer Vision*, 126(9):973–992, 2018. 2, 3

[69] Christos Sakaridis, Dengxin Dai, and Luc Van Gool. ACDC: The adverse conditions dataset with correspondences for semantic driving scene understanding. *arXiv preprint arXiv:2104.13395*, 2021. 2, 3

[70] Mattia Segu, Alessio Tonioni, and Federico Tombari. Batch normalization embeddings for deep domain generalization. *arXiv preprint arXiv:2011.12672*, 2020. 5

[71] Baochen Sun, Jiashi Feng, and Kate Saenko. Correlation alignment for unsupervised domain adaptation. In *Domain Adaptation in Computer Vision Applications*, pages 153–171. Springer, 2017. 3

[72] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2446–2454, 2020. 2, 3, 13

[73] Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei Efros, and Moritz Hardt. Test-time training with self-supervision for generalization under distribution shifts. In *International Conference on Machine Learning*, pages 9229–9248. PMLR, 2020. 2, 4

[74] Jean-Philippe Tarel, Nicholas Hautiere, Laurent Caraffa, Aurélien Cord, Houssam Halmaoui, and Dominique Gruyer. Vision enhancement in homogeneous and heterogeneous fog. *IEEE Intelligent Transportation Systems Magazine*, 4(2):6–20, 2012. 3

[75] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *European conference on computer vision*, pages 402–419. Springer, 2020. 6, 19

[76] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages 23–30. IEEE, 2017. 2

[77] Alessio Tonioni, Oscar Rahnama, Thomas Joy, Luigi Di Stefano, Thalaiyasingam Ajanthan, and Philip HS Torr. Learning to adapt for stereo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9661–9670, 2019. 2

[78] Alessio Tonioni, Fabio Tosi, Matteo Poggi, Stefano Mattoccia, and Luigi Di Stefano. Real-time self-adaptive deep stereo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 195–204, 2019. 2

[79] Frederick Tung, Jianhui Chen, Lili Meng, and James J Little. The raincouver scene parsing benchmark for self-driving in adverse weather and at night. *IEEE Robotics and Automation Letters*, 2(4):2188–2193, 2017. 2, 3

[80] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7167–7176, 2017. 2

[81] Girish Varma, Anbumani Subramanian, Anoop Namboodiri, Manmohan Chandraker, and CV Jawahar. IDD: A dataset for exploring problems of autonomous navigation in unconstrained environments. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1743–1751. IEEE, 2019. 3

[82] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. Deep hashing network for unsupervised domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5018–5027, 2017. 2

[83] Riccardo Volpi, Diane Larlus, and Grégory Rogez. Continual adaptation of visual representations via domain randomization and meta-learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4443–4453, 2021. 2, 4, 5

[84] Riccardo Volpi, Hongseok Namkoong, Ozan Sener, John Duchi, Vittorio Murino, and Silvio Savarese. Generalizing to unseen domains via adversarial data augmentation. *arXiv preprint arXiv:1805.12018*, 2018. 2

[85] Tuan-Hung Vu, Himalaya Jain, Maxime Bucher, Matthieu Cord, and Patrick Pérez. Advent: Adversarial entropy minimization for domain adaptation in semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2517–2526, 2019. 3, 5, 6

[86] Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully test-time adaptation by entropy minimization. *Int. Conf. Learn. Represent. (ICLR)*, 2021. 2, 4, 5, 6, 8

[87] Mei Wang and Weihong Deng. Deep visual domain adaptation: A survey. *Neurocomputing*, 312:135–153, 2018. 2, 5

[88] Xinshuo Weng, Yunze Man, Dazhi Cheng, Jinhyung Park, Matthew O'Toole, and Kris Kitani. All-In-One Drive: A Large-Scale Comprehensive Perception Dataset with High-Density Long-Range Point Clouds. *arXiv*, 2020. 2, 3

[89] Jingkang Yang, Kaiyang Zhou, Yixuan Li, and Ziwei Liu. Generalized out-of-distribution detection: A survey. *arXiv preprint arXiv:2110.11334*, 2021. 2

[90] Fisher Yu, Haofeng Chen, Xin Wang, Wenqi Xian, Yingying Chen, Fangchen Liu, Vashisht Madhavan, and Trevor Darrell. Bdd100k: A diverse driving dataset for heterogeneous multitask learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 2, 3, 6

[91] Fisher Yu, Vladlen Koltun, and Thomas Funkhouser. Dilated residual networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 472–480, 2017. 6, 7, 16, 18, 19

[92] Fisher Yu, Vladlen Koltun, and Thomas Funkhouser. Dilated residual networks. In *Computer Vision and Pattern Recognition (CVPR)*, 2017. 6

# Appendix

We provide additional details on our dataset in Sec. A. In particular, we report the sensor layout (Sec. A.1), annotation details (Sec. A.2), extensive information on dataset generation (Sec. A.3) and dataset statistics (Sec. A.4).

Moreover, we conduct additional experiments in Sec. B. We provide baselines on multitask learning under continual domain shift (Sec. B.1), and conduct ablation studies on joint training with real-world data (Sec. B.2) and the optimal dataset size for each task (Sec. B.3). Further, we propose a qualitative comparison between properties of SHIFT and the VIPER dataset [60] (Sec. B.4), and ablate on the model failures on the rainy and foggy domains (Sec. B.5).

Implementation details for each experiment conducted in this work are reported in Sec. C for full reproducibility.

## A. Dataset Details

The detailed user guide and additional information can be found at https://www.vis.xyz/shift.

### A.1. Reference systems and sensor layout

The dataset has three levels of reference systems: *world*, *vehicle*, and *camera*. The world system represents the absolute position of objects. The vehicle system is used for storing all 3D annotations. The camera systems are the reference systems used for each individual camera.

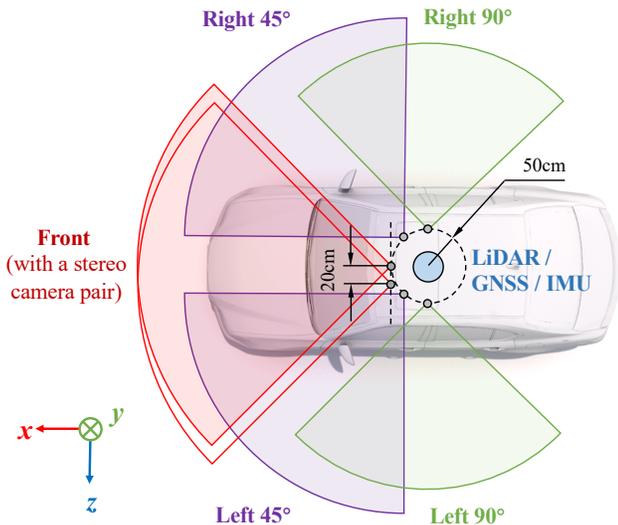Tab. 5 summarizes the supported sensors. We set up the



Figure 9. The vehicle system and the sensor layout. Except for stereo cameras, all the cameras are located on a circle centered at the vehicle reference system's origin (blue dot). LiDAR and motion sensors are located at the origin. Axes directions of the vehicle system are shown at the bottom left corner. *Best viewed in color.*

| Sensor | Data type | Position |
|---|---|---|
| RGB camera | 24-bit RGB | $5 \times$ RGB cameras (front, left / right $45°$, left / right $90°$). |
| Stereo camera | 24-bit RGB | Additional RGB camera offsetting 20cm toward left from the center. |
| Depth camera | 24-bit Gray | Same as front view RGB camera. |
| Optical flow | 32-bit UV | Same as front view RGB camera. |
| GNSS / IMU | Vector | Center of vehicle. |

Table 5. The data type and position settings of the sensors.

vehicle system following KITTI's convention and the right-hand rule. Specifically, the origin is located at the center of the ego-vehicle (marked as the blue dot in Fig. 9). Its $x$, $y$, and $z$ axes point in the right, down, and front directions, respectively (Fig. 9, bottom left). All the sensors are located on a circle centered at the vehicle reference system's origin, except for the stereo cameras that are placed on the left to the front camera, with a horizontal displacement of 20cm. All the cameras have a field-of-view (FoV) of $90°$. The 128-channel LiDAR sensor has a vertical FoV range of $[-10°, +10°]$ and a scan rate of 1.12M points per second.

Annotations stored in the vehicle system can be easily converted into the camera systems. Here, the front cameras and LiDAR sensor have camera systems identical to the vehicle system, so no conversion is needed for them. For other cameras, a vehicle-to-camera matrix (*i.e.* intrinsic and extrinsic parameters) is provided to transform the annotations so that they fit each camera.

### A.2. Annotation details

We present detailed specifications for the annotation set provided in SHIFT.

**Object detection** is a fundamental localization task for scene understanding and a basis for numerous downstream driving tasks, including multiple object tracking (MOT) and object re-identification (ReID). We provide 2D/3D bounding box annotations and object identities for six categories of traffic participants, *i.e.* car, truck, bus, bicycle, motorcycle, and pedestrian, together with the visibility attributes 'occluded' and 'truncated'. Moreover, for each box, we provide fine-grained object classes (*e.g.* vehicle model type).

While previous datasets only provide 7 DoF (*i.e.* only yaw angle) 3D boxes [5, 16, 72], we provide 9 DoF annotations and use the Euler angle system (*i.e.*, yaw, roll, pitch) to represent the orientation for bounding boxes in 3D space.

**Image segmentation** is a fundamental pixel-level perception task. For each frame, we provide panoptic (*i.e.* instance and semantic) segmentation labels on the 23 classes of the Cityscapes [9] annotation scheme. Together with 2D bounding boxes, segmentation labels can be used in multi-object tracking and segmentation (MOTS) and multi-object

panoptic tracking (MOTP) tasks.

**Depth estimation** is an essential step to extend the 2D perception tasks into the 3D setting. We provide the depth labels aligned with the front-view RGB camera to enable image- and video-based monocular and stereo depth estimation. Depth resolution is 1mm.

**Optical flow estimation** is an essential task for driving algorithms involving motion. However, existing large-scale datasets typically do not provide optical flow annotations due to the high labeling cost. Representing the relative motion between each pixel in a pair of images, optical flow can be instrumental in object tracking and ego-motion tasks. We provide the optical flow labels in the UV map format, also used in KITTI [16].

### A.3. Data generation pipeline

We introduce the pipeline that used to generate the discrete and continuous domain shifts.

**Disrete shift.** As discussed in Sec. 3.3, we set up an efficient sampling pipeline that can cover a diverse combination of conditions. To determine the environmental parameters of each sequence, we use a technique similar to random search. In Tab. 6, we define 4 categories of domain shifts, *e.g.*, time of day, weather, vehicle density, and pedestrian density. For the $i$-th category ($1 \leq i \leq 4$), we define a set of *candidate* domains, $\mathcal{H}_i = \{h_i^{(1)}, \cdots, h_i^{(n_i)}\}$, where each candidate $h_i^{(j)}$ corresponds to a certain group of environmental parameters, defined in the Tab. 6. Note that the parameter can be a fixed value or a set of values. For the set of values, we again uniformly sample one value out of the set.

Our sampling method for the discrete domain shifts can be summarized as following. A sequence is generated with a fixed parameter vector $\boldsymbol{\theta} = \theta_1 \cup \cdots \cup \theta_m$, where each $\theta_i$ is sampled uniformly across all candidates in the $i$-th category, *i.e.*,

$$h_i^{(j)} \sim \text{Uniform}(\mathcal{H}_i), \quad \forall i = \{1, 2, 3, 4\} \quad (1)$$

$$\theta_i \sim h_i^{(j)} \quad (2)$$

This pipeline guarantees the uniform marginal distribution of candidates conditioned on any category. Using this pipeline, we can easily add data without breaking the distribution of domains. Moreover, any subset of sequences of SHIFT has the same distribution, allowing a fair experiment on the impact of data amount.

**Continuous shifts.** For sequences with continuous domain shift, the change of parameters happens on one specific domain category $c$, while others are kept unchanged, *i.e.* the frame at time $t$ is generated with the parameter vector

$$\boldsymbol{\theta}(t) = \theta_1 \cup \cdots \cup f_c(t) \cup \cdots \cup \theta_4 , \quad (3)$$

where $f_c(t)$ is obtained by linear interpolation of the states listed in the 'Environmental parameters' column in Tab. 6. Specifically, $f_c(t)$ is obtained by interpolating the points

$$(t, \theta) = [(0, \theta_{c,\text{begin}}), (0.2, \theta_{c,\text{intermediate}}), (1, \theta_{c,\text{end}})] , \quad (4)$$

where $t \in [0, 1]$ represents the degree of continuous shift from the minimum to the maximum parameter allowed for a given domain category.

Our dataset provides 300, 120, 30 continuous shift sequences in 1x, 10x, 100x length respectively, where the base length (1x) is 200 frames. Since continuous domain shift is generated by interpolating between the state of the initial frame and the state of the final frame, the domain shift speed is inversely proportional to the sequence length in frames. Furthermore, we provide an additional set of 150 sequences of base length presenting domain shifts simultaneously happening along multiple domain shift directions within the same sequence.

**Domain labeling details.** The degree of shift for each domain category is quantified by a numerical value called *severity*. For weather conditions, we use percentage values to indicate the degree of severity, where $0\%$ corresponds to clear weather conditions and $100\%$ represents the most extreme condition allowed by the CARLA simulator for a given weather direction, *e.g.* cloudiness, precipitation, fog density, or fog distance. We describe the time of day using the Sun's altitude angle to disentangle the lighting condition with the sunrise/sunset time. For the object densities, we use the number of objects per frame as the severity (Tab. 6).

### A.4. Dataset statistics

SHIFT is diverse in bounding box scale. Fig. 10 (left) plots the object density measured by boxes per frame and shows coverage from 0 to 30 boxes/frame for SHIFT. We compare the distribution with the BDD100K's MOT set. Due to the sparsity of the vehicle/pedestrians density domains, our dataset has on average a higher density of frames counting less bounding boxes than BDD100K, but the crowded frames ($\geq$ 20 boxes/frame) show similar trends. Moreover, Fig. 10 (right) shows the distribution of bounding box sizes, defined as $\sqrt{wh}$ where $w$ and $h$ are the width and height of a box. SHIFT covers diverse box sizes ranging from 10 to 650 pixels. We also observe that our dataset has 41.2% bounding boxes smaller than 15 pixels while BDD100K has 30.9%, showing that our dataset provides challenging conditions for small object detection and tracking.

## B. Additional Experiments

To further highlight the usefulness of SHIFT, we conduct experiments on multitask learning (Sec. B.1) and joint

14

| Category $\mathcal{H}_i$ | Candidate dom. $h_i^{(j)}$ | | BDD100K eq. | Environmental parameters | Degrees of shift |
|---|---|---|---|---|---|
| Time of day | noon<br>morning / afternoon | } daytime | Sun altitude angle = {90, 75, 60, 45, 30}<br>Sun altitude angle = {15, 10, 5} | altitude angle $\in [-5, 90]$ |
| | dawn / dusk<br>sunrise / sunset | } dawn / dusk | Sun altitude angle = {4, 3, 2}<br>Sun altitude angle = {1, 0, -1} | |
| | night<br>dark night | } night | Sun altitude angle = {-2, -3}<br>Sun altitude angle = {-4, -5} | |
| Weather | clear<br>slight cloudy | } clear | cloudiness = {0, 5}<br>cloudiness = {10, 15} | cloudiness $\in [0, 100]$ |
| | partly cloudy<br>overcast | partly cloudy<br>overcast | cloudiness = {25, 50, 70}<br>cloudiness = 100 | |
| | small rain<br>mid rain<br>heavy rain | rainy | cloudiness = 70; precipitation = 20; deposit = 60; fog den. = 3<br>cloudiness = 80; precipitation = 50; deposit = 80; fog den. = 3<br>cloudiness = 100; precipitation = 100; deposit = 100; fog den. = 7 | precipitation $\in [0, 100]$ |
| | small fog<br>heavy fog | foggy | cloudiness = 60; fog density = 30; fog distance = 15<br>cloudiness = 80; fog density = 90; fog distance = 20 | fog density $\in [0, 100]$ |
| Vehicle density | sparse<br>moderate<br>crowded | -<br>-<br>- | num of vehicle = 50<br>num of vehicle = 100<br>num of vehicle = 250 | vehicle per map,<br>vehicle per frame |
| Pedestrian density | sparse<br>moderate<br>crowded | -<br>-<br>- | num of pedestrians = 100<br>num of pedestrians = 200<br>num of pedestrians = 400 | pedestrian per map,<br>pedestrian per frame |

Table 6. Definitions of the domain category and candidate domains, used for discrete domain shifts. Each category has a group of candidate domains. For each candidate domain, we show its equivalent domain label in BDD100K and the environmental parameters for simulation.

| Continuous shift type | Environmental parameters | | |
|---|---|---|---|
| | **Beginning state** ($t = 0$) | **Intermediate state** ($t = 0.2$) | **End state** ($t = 1$) |
| Time of day | Sun altitude angle = 90 | - | Sun altitude angle = -5 |
| Cloudiness | cloudiness = 0 | - | cloudiness = 100 |
| Raininess | cloudiness = 0, precipitation = 0,<br>deposit = 0, fog density = 0 | cloudiness = 80, precipitation = 50,<br>deposit = 80, fog density = 3 | cloudiness = 100, precipitation = 100,<br>deposit = 100, fog density = 7 |
| Fogginess | cloudiness = 0, fog density = 0,<br>fog distance = 0 | cloudiness = 60, fog density = 30,<br>fog distance = 15 | cloudiness = 80, fog density = 90,<br>fog distance = 20 |

Table 7. Definitions of parameters used for continuous domain shifts. The parameters are updated for every frames during driving. The value of parameters are determined by linear interpolation between the state of beginning, intermediate (if applicable) and end.
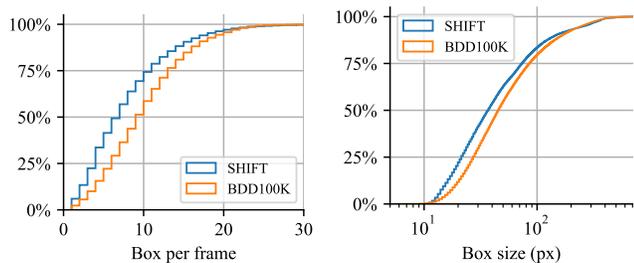


Figure 10. Cumulative distributions of the bounding box per frame (left) and bounding box size measured in $\sqrt{wh}$ (right). We only count the objects in the front camera view. SHIFT covers various object densities and a wide range of object scale.

training with real-world data (Sec. B.2). We also investigate the optimal dataset size and sampling rate (Sec. B.3).

## B.1. Multitask learning

In this experiment, we study whether different perception tasks mutually benefit or interfere with each other when jointly learned with a shared feature extractor. The wide variety of tasks supported in SHIFT unlocks new opportunities to investigate different combinations of perception tasks. Special attention is also paid to the robustness of multitask models under incrementally shifted domains.

Specifically, we consider four different perception tasks: semantic segmentation, instance segmentation, monocular depth estimation, and optical flow estimation. Each task requires the model to learn a distinct encoding function: semantic segmentation requires intermediate activations to encode pixel-level information, instance segmentation requires instance-level information, depth estimation requires contextual information and object priors that allow to con-

| Task | Train | Metric | Source | OOD | | | | | | OOD Avg. | $\Delta_{\text{Source}}$ | $\Delta_{\text{OOD}}$ | $\Delta_{\text{S}\to\text{O}}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | clear-daytime | cloudy | overcast | foggy | rain | dawn/dusk | night | | | | |
| Semantic segmentation (S) | S | mIoU (%) ↑ | 69.1 | 40.6 | 40.6 | 21.5 | 19.6 | 18.1 | 8.9 | 24.9 | - | - | -64.0% |
| | S + D | | **75.2** | 53.8 | 52.6 | 24.3 | 26.6 | 24.0 | 9.9 | 31.9 | **8.9%** | 28.1% | -57.6% |
| | S + F | | 69.4 | 51.8 | 54.7 | 26.4 | 22.4 | 22.7 | 9.8 | 31.3 | 0.4% | 25.8% | -54.9% |
| | S + D + F | | 71.8 | 50.0 | 51.9 | 23.5 | 24.0 | 22.1 | 9.5 | 30.2 | 3.8% | 21.2% | -58.0% |
| | S + I | | 74.8 | **63.9** | **68.1** | **41.0** | **36.8** | **37.3** | **23.6** | **45.1** | 8.2% | **81.3%** | **-39.7%** |
| | S + D + I | | 75.0 | 62.4 | 65.1 | 37.4 | 35.4 | 35.3 | 20.5 | 42.7 | 8.6% | 71.6% | -43.1% |
| | S + F + I | | 72.5 | 58.3 | 59.6 | 35.8 | 27.2 | 28.8 | 14.4 | 37.3 | 4.9% | 50.1% | -48.5% |
| | S + D + F + I | | 74.7 | 60.6 | 59.6 | 37.1 | 32.7 | 33.1 | 19.3 | 40.4 | 8.1% | 62.3% | -45.9% |
| Depth estimation (D) | D | SILog ↓ | 17.8 | 28.3 | 23.1 | 81.9 | 46.3 | 54.6 | 63.2 | 49.6 | - | - | -64.1% |
| | D + S | | 16.9 | 25.2 | 22.4 | 65.7 | 43.0 | 49.4 | 57.6 | 43.9 | 5.6% | 13.0% | -61.6% |
| | D + F | | 19.3 | 25.3 | 20.4 | 66.6 | 45.3 | 50.3 | 54.4 | 43.7 | -7.8% | 13.4% | -55.8% |
| | D + S + F | | 19.6 | 26.9 | 24.7 | 67.8 | 45.1 | 52.1 | 56.6 | 45.5 | -9.2% | 8.9% | -56.9% |
| | D + I | | 17.3 | 21.0 | 16.8 | 66.4 | 35.1 | 42.4 | 48.4 | 38.3 | 2.9% | 29.3% | -54.8% |
| | D + S + I | | **16.0** | **19.5** | **15.4** | 61.1 | **31.2** | **38.5** | **42.7** | **34.7** | **11.0%** | **42.8%** | -53.8% |
| | D + F + I | | 17.8 | 21.4 | 17.9 | **47.9** | 36.4 | 39.9 | 46.3 | 35.0 | 0.1% | 41.8% | **-49.1%** |
| | D + S + F + I | | 17.6 | 21.9 | 18.3 | 53.2 | 37.1 | 42.7 | 50.5 | 37.3 | 1.0% | 33.0% | -52.7% |
| Optical flow estimation (F) | F | EPE (px) ↓ | **6.0** | **6.7** | **6.4** | 9.0 | 9.7 | 9.1 | 11.0 | 8.6 | - | - | -30.8% |
| | F + S | | 7.8 | 8.3 | 8.5 | 10.4 | 12.0 | 10.9 | 12.6 | 10.4 | -23.1% | -17.4% | -25.7% |
| | F + D | | 6.0 | 6.9 | 6.4 | 9.4 | 10.4 | 9.6 | 11.8 | 9.1 | -0.2% | -5.0% | -34.2% |
| | F + D + S | | 6.1 | 8.5 | 8.3 | 10.6 | 12.1 | 11.0 | 13.0 | 10.6 | -2.1% | -18.5% | -42.4% |
| | F + I | | 9.8 | 9.6 | 9.6 | 10.4 | 11.3 | 10.6 | 12.1 | 10.6 | -38.9% | -18.5% | **-7.7%** |
| | F + S + I | | 7.7 | 8.2 | 7.9 | 9.7 | 10.6 | 9.8 | 11.8 | 9.7 | -22.7% | -10.8% | -20.2% |
| | F + D + I | | 8.0 | 8.3 | 8.2 | 9.9 | 10.6 | 10.0 | 11.9 | 9.8 | -25.5% | -12.1% | -18.4% |
| | F + D + S + I | | 8.1 | 8.4 | 8.4 | 10.1 | 11.0 | 10.2 | 12.1 | 10.0 | -26.4% | -14.0% | -19.2% |
| Instance segmentation (I) | I | mAP (%) ↑, vehicles | 63.9 | 57.4 | 65.7 | 21.9 | 31.2 | 22.7 | 6.6 | 34.2 | - | - | 46.4% |
| | I + S | | 64.9 | 59.1 | 66.2 | 26.4 | **34.4** | 27.1 | **14.3** | **37.9** | 1.5% | **10.7%** | 41.6% |
| | I + S + D | | 65.0 | 57.9 | 64.9 | 25.9 | 32.6 | 26.1 | 10.9 | 36.4 | 1.6% | 6.3% | 44.0% |
| | I + S + F | | 62.3 | 57.1 | 64.2 | 21.6 | 31.6 | 23.6 | 7.7 | 34.3 | -2.5% | 0.1% | 45.0% |
| | I + D | | **65.9** | **59.3** | 66.9 | **26.8** | 32.2 | **26.3** | 11.4 | 37.2 | **3.1%** | 8.5% | 43.6% |
| | I + D + F | | 65.8 | 50.2 | **67.0** | 21.5 | 31.0 | 22.9 | 6.7 | 33.2 | 2.9% | -3.0% | 49.5% |
| | I + F | | 63.1 | 56.9 | 65.1 | 20.4 | 28.9 | 21.7 | 5.0 | 33.0 | -1.3% | -3.6% | 47.7% |
| | I + S + D + F | | 64.8 | 57.9 | 65.3 | 22.8 | 31.5 | 23.1 | 8.3 | 34.8 | 1.4% | 1.6% | 46.3% |

Table 8. Multitask learning performances. We evaluate 15 combinations of 4 perception tasks: semantic segmentation (S), monocular depth estimation (D), optical flow estimation (F), and instance segmentation (I). The combinations of S + I, S + D, and S + D + I significantly improve on both tasks' source and OOD performance in their respective tasks. ↑ (↓): the higher (lower) the better.

vert 2D images to 3D cues, and optical flow requires to encode a function of two images that embodies information on motion perception.

**Multitask model.** To compose a unified multitask model, we use the segmentation model DRN-D-54 [91] as feature extractor and combine it with the heads required for other tasks. The DRN-D-54 model has 8 sequential residual blocks with dilated convolutions and transposed convolutions at the end to generate segmentation results. Here, all the modules of DRN-D-54 are used for semantic segmentation. For instance segmentation, we rely on the Feature Pyramid Network (FPN) [96], Region Proposal Network (RPN), and ROIAlign modules identical to those introduce in Mask R-CNN [19]. FPN uses the 2nd to 5th blocks' outputs of the DRN-D network. For the optical flow and depth estimation, we adapt the decoders similar to FlowNet [94] and U-Net [62]. The decoder has 5 sequential blocks, where each block has one up-sampling layer, followed by a short-cut connection from the feature extractor's corresponding block, and a series of convolution layers. Together with the feature extractor, we obtain an encoder-decoder structure commonly used in dense prediction tasks.

**Experiment setup.** We traverse all 15 combinations for the 4 tasks mentioned above. Our multitask model is trained with 5,000 frames sampled from the clear-daytime domain in SHIFT and evaluated under different discrete domain shifts. To fit the multitask model into the GPU memory, we reduce the image size to $640 \times 400$ pixels. Please note that the performance will be slightly affected by the size-reduced images and thus, it is not directly comparable to our baseline experiments in 2 of the main paper. All combinations are trained for 100 epochs, when convergence is reached for all tasks.

**Experimental results** are summarized in Tab. 8. Every model is trained on the clear-daytime domain and tested on
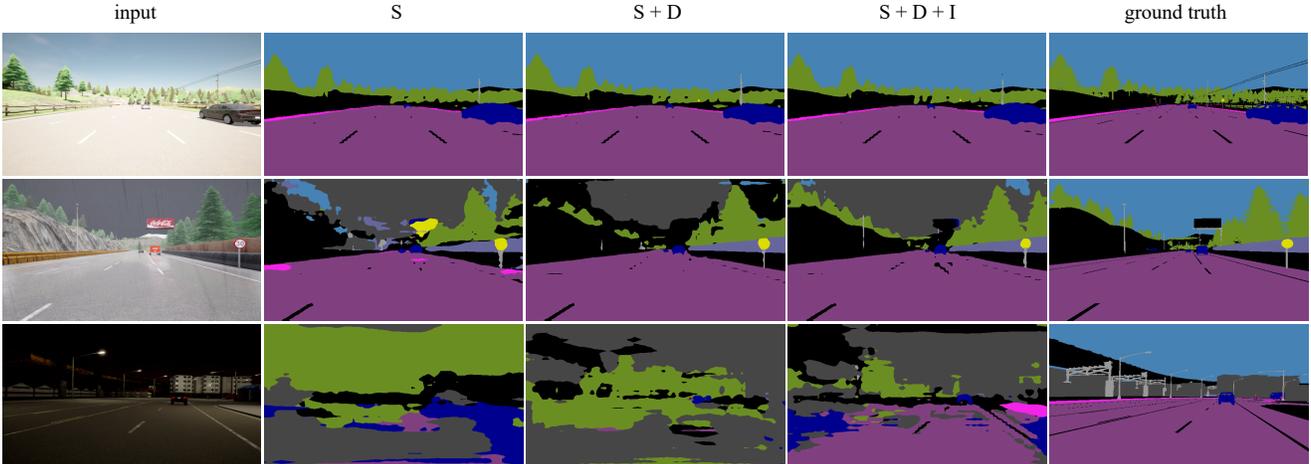
Figure 11. Qualitative results on semantic segmentation for models trained on the clear-daytime domain. Each column represents a model trained on a different task combination: semantic segmentation (S), semantic segmentation + depth estimation (S + D), semantic segmentation + depth estimation + instance segmentation (S + D + I). The three rows show the results on respectively the clear-daytime, rainy, and night domain. The combinations S + D and S + D + I improve the performance against domain shifts.

different types of shifted domains, indicated with OOD in the Table. We report the average performance on the out-of-distribution domains as OOD avg. The columns $\Delta_{\text{Source}}$ and $\Delta_{\text{OOD}}$ report for different multitask models the relative Source / OOD avg. performance change on a given task with respect to the performance of a single-task model trained on that specific task. The column $\Delta_{S \to O.}$ reports for different multitask models the relative change from Source to OOD avg. performance on a given task. Below are our observations.

***Multitask learning improves robustness.*** We observe that specific combinations of tasks largely improve the single-task model performance on the source domain. For instance, the combination of semantic segmentation (S) + depth (D) + instance segmentation (I) boosts the source domain performance by 8.6% / 11.0% / 1.6% on the respective tasks. Similar improvements are observed for other combinations, including S + I and S + D. We visualize the results of these combinations in Fig. 11. This is possibly due to the intertwined nature of such tasks. In particular, depth and semantics both need to learn contextual features from neighboring pixels, and both instance and semantics segmentation need to segment parts of the image.

Further, multitask learning often significantly increases the generalization of a model to domain shifts. For example, the combination of S + D + I improves the OOD performance in the respective tasks by 71.6% / 42.8% / 6.3%. The improvements are substantially greater than the improvements on the source domain, suggesting that the increase in model's robustness is not attributable to the increase in the overall model's performance as seen on the source domain. We argue that this is potentially due to the model learning

more general features that are shareable across tasks and, consequently, also more general under domain shifts. For example, the addition of instance segmentation typically causes the greatest robustness improvements. This might be due to the complex nature of the instance segmentation task, which requires to encode features capable of both detecting and segmenting objects in an image.

***Instance segmentation can only be improved mildly.*** Instance segmentation is only improved at most by 10.7% on OOD performance by other tasks. As previously mentioned, we hypothesize that instance segmentation already learns more general features due to its nature. Thus, the addition of other tasks provides only mild improvements. On the other hand, however, when combined with other tasks, instance segmentation largely boosts their robustness, *e.g.* S + I and D + I.

***Optical flow is heavily affected by other tasks.*** Unlike the previous tasks that benefit from multitask learning, optical flow shows a different behavior. Although optical flow can improve other tasks' robustness (*e.g.* S + F and D + F), the optical flow itself is negatively affected by the addition of other tasks. When jointly trained with other tasks, its performance drops by a large margin, ranging from -0.2% to -38.9%. A possible explanation is that the optical flow task, which takes a pair of frames as input, learns a different encoding function than other non-temporal tasks requiring only one frame. To learn a feature extractor shared across the two different types of inputs, the model shows to sacrifice its effectiveness on the task requiring two images. This suggests that combining different tasks is not trivial; instead, it requires extensive evaluation and comparison. SHIFT provides a playground to develop novel mul-

titask learning techniques and to investigate and solve the multiple challenges presented by such an interesting problem.

***Domain shift is only partially mitigated.*** While the model's robustness can be improved by multitask learning, the domain shifts provided in SHIFT still pose a tremendous threat to the robustness under domain shift. For all the evaluated tasks, the minimum average OOD performance drop with respect to the corresponding source performance ($\Delta_{S \to O.}$) amounts to $\sim 40\%$. Under extreme conditions, *e.g.* foggy and night, the performances are degraded even more than 60%, which indicates real-life risks if autonomous vehicles heavily rely on such models.

By introducing SHIFT, which supports multi-domain and multitask studies in a single dataset, we hope to foster future research on multitask domain adaptation algorithms to counteract these domain gaps effectively. Moreover, we hope that the continuous domain shifts provided in our dataset will shed new light on this challenging problem.

## B.2. Joint training with real-world data

We investigate whether the domain variations in our dataset in combination with a specific domain of real-world data can make a model more robust to domain shift compared to a model only trained on the real-world data. Specifically, we jointly train the model with the source domain data (*i.e.*, clear daytime) from BDD100K and all domain variations from ours. The model is then evaluated on other domains of BDD100K. We employ the Faster R-CNN [59] as the model for object detection and DRN-D-54 [91] for semantic segmentation. The models are learned with the same amount of data from BDD100K but with different amounts of data from SHIFT.

**Object detection** results are shown in Tab. 9. We observe that the joint training provides a relative improvement of the source domain and OOD performance amounting to 2.52% and 3.40%, respectively.

**Semantic segmentation** has similar trends. As shown in Tab. 10, source domain mIoU improves from 46.04% to 51.20%, with a relative improvement of 10.34%. Moreover, out-of-domain mIoU rises by a relative 5.30% from 37.37% to 39.76%.

These results suggest that, if a model is trained on a limited real-world domain, jointly training with the variety of domains provided by our dataset will improve the robustness of the model to real-world shifts.

## B.3. Dataset size

To understand the impact of dataset size and optimize the design of the dataset, we conduct ablation studies on: (1) sampling rate and (2) amount of sequences. Every model is trained on clear-daytime sequences.

| Training set | Source domain | | OOD avg. | |
|---|---|---|---|---|
| | AP | AP$_{75}$ | AP | AP$_{75}$ |
| BDD100K | 0.318 | 0.312 | 0.265 | 0.251 |
| BDD100K + 2k frames | 0.320 | 0.327 | 0.267 | 0.267 |
| BDD100K + 5k frames | **0.326** | **0.334** | **0.274** | **0.271** |
| BDD100K + 10k frames | 0.325 | 0.329 | 0.254 | 0.238 |

Table 9. Joint training for object detection. Generalization ability is improved with a proper amount of data.

| Training set | Source domain | OOD avg. |
|---|---|---|
| BDD100K | 46.04 | 37.37 |
| BDD100K + 6k frames | 47.11 | 38.56 |
| BDD100K + 12k frames | **51.20** | **39.76** |
| BDD100K + 24k frames | 51.09 | 39.23 |

Table 10. Joint training for semantic segmentation. We report the mIoU. Generalization ability is improved with a proper amount of data.

| Frame rate (Hz) | 0.1 | 0.2 | 0.5 | 1 | 5 | 10 |
|---|---|---|---|---|---|---|
| # Frames ($\times$1k) | 7.5 | 15 | 37.5 | 75 | 375 | 750 |
| Seg. (mIoU, %) | 62.6 | 62.9 | **63.1** | 63.0 | 62.9 | - |
| Det. (mAP, %) | 40.6 | 43.1 | 45.8 | 46.8 | **48.4** | - |
| MOT (MOTA, %) | 25.6 | 34.7 | 45.2 | 49.3 | 54.1 | **54.9** |

Table 11. Performance of different tasks at increasing sampling rates. Training and testing on the same 1500 sequences from all domains.

| Training sequence | 350 | 750 | 1500 | 2000 | 3000 |
|---|---|---|---|---|---|
| Seg. (mIoU, %) | 59.4 | 61.4 | 63.0 | 62.6 | **63.1** |
| Det. (mAP, %) | 41.2 | 45.1 | 46.8 | 48.0 | **50.1** |

Table 12. Performance of different tasks at increasing sequences number. Training and testing on the data of 1Hz from all domains.

**Frame rate.** To avoid the model learning from redundant information, we study what is the optimal sampling rate to achieve the best performance on a given task. Here, we test the semantic segmentation, object detection, and multiple obeject tracking performance on a set of images sampled at different frame rates from a fixed set of 2000 sequences. We notice that performance of different tasks starts to saturate at different sampling rates (Tab. 11). For image-based tasks, such as segmentation and detection, we argue that the information provided by adjacent frames can be redundant, and increasing the sampling rate over a certain threshold have insignificant benefits on the resulting model performance. However, for video-based tasks, like multi-object tracking, the inter-frame information is crucial. A lower frame rate leads to lose a considerable amount of information, thus severely reducing the model performance (Tab. 11, third row).

Our dataset is collected at a fixed frame rate of 10Hz, which is necessary to support a wide range of perception

tasks. However, according to the experiments on the sampling rate, we also provide a subset sampled at 1Hz for image-based perception tasks.

**Amount of sequences** is another factor affecting the performance. Here, we test semantic segmentation and object detection performance on a varying number of sequences sampled at 1Hz. In Tab. 12, we find that the performance continuously increases up to 3000 sequences. However, the performance gain is diminishing the more sequences we add. This is potentially due to the limited environmental variation in the simulator. To balance between size and learning performance, we set the total number of sequences to 3000 for the discrete set. Together with our sampling pipeline (Sec. A.3), the current size of SHIFT guarantees that for each BDD100K's domain label, we have more than 500 corresponding sequences for training and testing.

### B.4. Comparison with VIPER

As a synthetic dataset, VIPER [60] also presents sequences from discrete domain shifts. Here, we compare the segmentation performance under domain shifts in VIPER, SHIFT and BDD100K (Tab. 13). We find that the adverse conditions presented in VIPER provide a less relevant threat to model generalization, highlighting how SHIFT mimics more closely real-world trends.

| Dataset | daytime ($M_0$) | sunset | night | rain | $\frac{\max \Delta M}{M_0}$ |
|---|---|---|---|---|---|
| VIPER | 59.3 | 57.6 | 55.1 | 53.0 | -10.6% |
| SHIFT (*ours*) | 83.6 | 60.4 | 42.8 | 54.6 | -48.8% |
| BDD100K | 47.9 | - | 20.6 | 37.6 | -57.0% |

Table 13. Out-of-distribution performance on different datasets of a segmentation model (DRN-D) trained on the daytime domain. The last column represents the maximal relative performance drop w.r.t. source.

### B.5. Error analysis for foggy and rainy domains

As noticeable in Fig. 5, detection and segmentation models show a slightly different behavior under different types of domain shift. While it is worth noticing that segmentation and detection have different label sets, we here analyze the differences in performance on the two domains presenting the largest discrepancy across the two tasks, *i.e.* foggy and rainy. For example, we find that the most drastically affected class for segmentation in the rainy domain is 'sky' (-69% mIoU w.r.t. clear-daytime), with 25% of the corresponding pixels misclassified as 'building', as opposed to only 2% under foggy conditions. For object detection, we find that most of the errors come from missed detections. The shifted domains lower the classification confidence below the pre-selected threshold, with foggy posing a greater challenge (car AP drops by 74% on foggy vs 40% on rainy).

## C. Implementation Details

In this section, we describe the implementation details and metrics for each task in Tab. 2 and Fig. 5.

**Object detection.** We compare Faster R-CNN [59], Cascade R-CNN [6], and YOLO v3 [58]. The backbone network for the first two methods is ResNet-50 [20], while YOLO v3 uses DarkNet [57] as its backbone. We use the mean Average Precision (mAP) as the metric for 2D bounding boxes. We train the models on 50k frames of data, following the "1x" schedule provided in the `mmdetection` library [93].

**Semantic segmentation.** We also compare three models for semantic segmentation, DeepLab v3+ [8], Fully Convolutional Network (FCN) [37], and DRN-D-54 [91]. All three models use the ResNet-50 [20] as the backbone. We train the models with 20k frames of data until they converge (approximately 150 epochs). We use the mean IoU (mIoU) metric for all evaluations on semantic segmentation.

**Instance segmentation.** We use Mask R-CNN [19] with ResNet-50 backbone and follow the same training routine as Faster R-CNN [59]. A segmentation mAP metric is used for evaluation.

**Depth estimation.** We use AdaBins [3] for the depth estimation experiments. It uses a U-Net-like [62] backbone structure and predicts depth with adaptive bins. The model is trained using its official implementation. We follow KITTI's benchmark on depth estimation [16]. Specifically, we use the Scale-invariant Logarithm (SILog) metric evaluated on the central crop of the image (*i.e.* Eigen Crop).

**Optical flow estimation.** We use RAFT [75] for optical flow estimation. The model is fine-tuned from pre-trained weights on the Things Dataset [42], with 10k frames of our data. The End-point Error (EPE) metric is used for evaluation.

## References

[93] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tianheng Cheng, Qijie Zhao, Buyu Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. MMDetection: Open mmlab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019. 19

[94] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. Flownet: Learning optical flow with convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2758–2766, 2015. 16

[95] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 6, 16, 19

[96] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017. 16

[97] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28:91–99, 2015. 7, 18, 19

[98] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 16, 19

[99] Fisher Yu, Vladlen Koltun, and Thomas Funkhouser. Dilated residual networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 472–480, 2017. 6, 7, 16, 18, 19