

# From HD Maps to No Maps

Wolfram Burgard

---

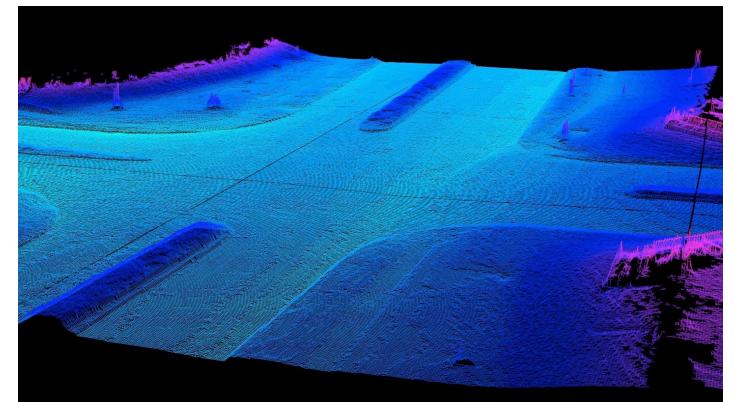
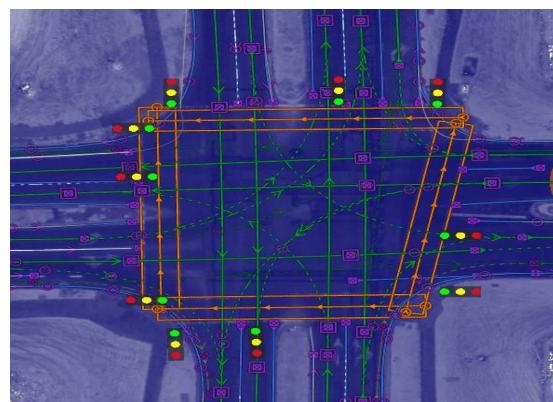
Joint work with: David Pannen, Martin Liebner, Jannik Zürn,  
Johan Vertens, Kshitij Sirohi, Juana Valeria Hurtado, Rohit  
Mohan, Daniel Büscher, Abhinav Valada ...



# Maps in Automated Driving and Robotics

Useful for

- Perception
- Tracking
- Localization
- Prediction
- Planning
- Control
- ...



# Challenges for HD Maps

- Expensive to acquire
- Assumptions about availability of features
- Change detection
- Expensive to update
- L5 barrier
- ...



## In this Talk

- Change detection through **crowdsourcing**
- **Learning** for on-the-fly estimation of HD map information
  - Semantics
  - Tracking
  - Topology

# Change Detection for HD Maps

2019 International Conference on Robotics and Automation (ICRA)  
Palais des congrès de Montréal, Montréal, Canada, May 20-24, 2019

## HD Map Change Detection with a Boosted Particle Filter

David Pannen<sup>1</sup>, Martin Liebner<sup>1</sup> and Wolfram Burgard<sup>2</sup>

**Abstract**—In this paper, we present a change detection algorithm that can run in real time as part of a backend-based stream processing pipeline. It can process the floating car data collected by series-production vehicles to detect changes in an automotive high definition digital (HD) map used for automated driving. The algorithm uses sensor readings matched with odometry, GNSS and landmark detections to localize the vehicle within the digital map. While all particles together represent the probability distribution for the vehicle's position at a given time, each individual particle also serves as a hypothesis about the vehicle's position. This is used to compute various metrics for how well the current sensor readings match the world model encoded in the HD map. The different metrics are evaluated by a number of weak classifiers that are used as input for a trained AdaBoost classifier. The achievable detection rate of a single vehicle is then compared to that of a simple crowd-based approach, where each vehicle votes on whether or not the current section of the road has changed.

### I. INTRODUCTION

To achieve the goal of safe, comfortable and efficient automated driving, high definition digital maps (HD maps) may be used to augment the vehicle's on-board sensors with prior beliefs about the environment and thus help to provide localization and sensor outlier detection. HD maps can be part of a fall-back strategy in case of sensor failure, but they also contribute to an increased automated driving comfort by enhancing the vehicle's foresight. Finally, they contain the baseline knowledge about drivable lanes, lane connectivity, and applicable traffic regulations that is necessary for automated driving.

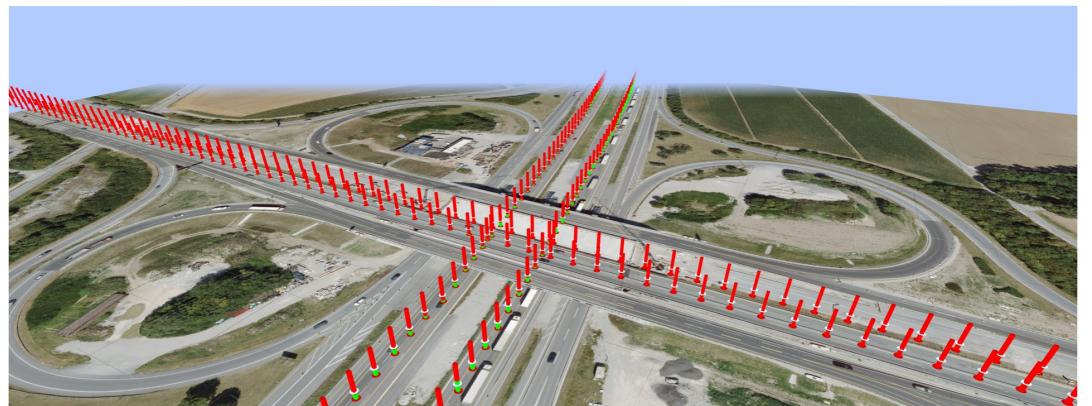
The current state of the art for creating HD maps uses mapping vehicles equipped with a suite of high accuracy sensors such as light detection and ranging (LiDAR) sensors or high accuracy global navigation satellite system (GNSS) sensors [1]. The problem with this approach is the high cost and resulting limited number of mapping vehicles. The mapping vehicle HERE, for example, at the end of 2015, used a fleet of just over 200 cars worldwide [1]. As those vehicles can cover only a small fraction of the road network per day, the frequency of visits for each road segment is low and changes cannot be incorporated into the map in time.

This poses a high risk for automated driving, as a faulty map could lead to dangerous situations or even accidents. Inconsistencies between the map and reality must therefore be detected as fast as possible. Due to the reasons stated above, this is not achievable with conventional mapping

<sup>1</sup>D. Pannen and M. Liebner are with BMW Group, 80788 Munich, Germany. [firstname.lastname@bmw.de](mailto:firstname.lastname@bmw.de)

<sup>2</sup>W. Burgard is with the Department of Computer Science, University of Freiburg, Georges-Köhler-Allee 080, 79110 Freiburg i. Br., Germany. [burgard@informatik.uni-freiburg.de](mailto:burgard@informatik.uni-freiburg.de)

## HD Map Change Detection with a Boosted Particle Filter. David Pannen, Martin Liebner and Wolfram Burgard, ICRA 2019



For evaluation purposes, we employ the developed algorithm to detect changes caused by motorway construction sites. This is primarily due to the fact that information about motorway construction sites is publicly available in Germany [7, 8] which facilitates obtaining ground-truth information. Also, motorway construction sites are of particular interest for the deployment of automated driving, as motorways are the primary target of most early automated driving systems. Since most of those systems cannot handle construction sites yet, especially without a valid map, changes in their location must be detected quickly in order to close affected sections for automated driving functions.

### A. Related Work

Detecting changes in a given map is an important problem with regards to life-long mapping and thus has been subject to extensive research in the robotics community. The approaches to solve this problem can be divided into two major categories: handling dynamic changes and handling semi-static changes. Dynamic changes refer to occlusions by fast moving objects such as pedestrians that can be detected directly from measurements as the time scale of their move-

# Creating HD Maps

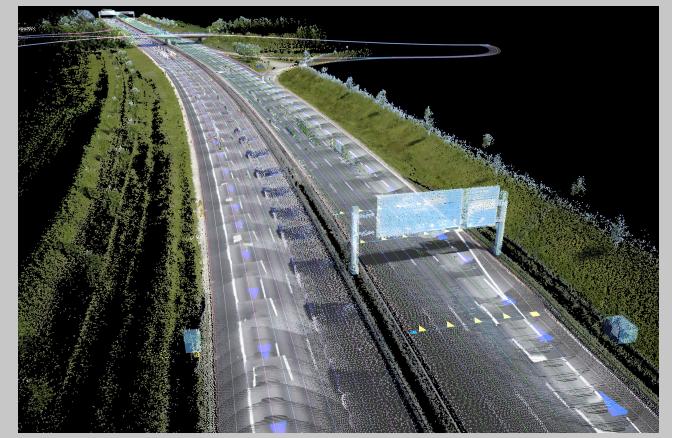
Automated driving promises increased

- Safety
- Comfort
- Efficiency



HD maps enable

- Sensor augmentation
- Extended foresight
- Fallback for sensor failures



Environmental changes frequently occur and make change detection necessary



State of the art mapping vehicles

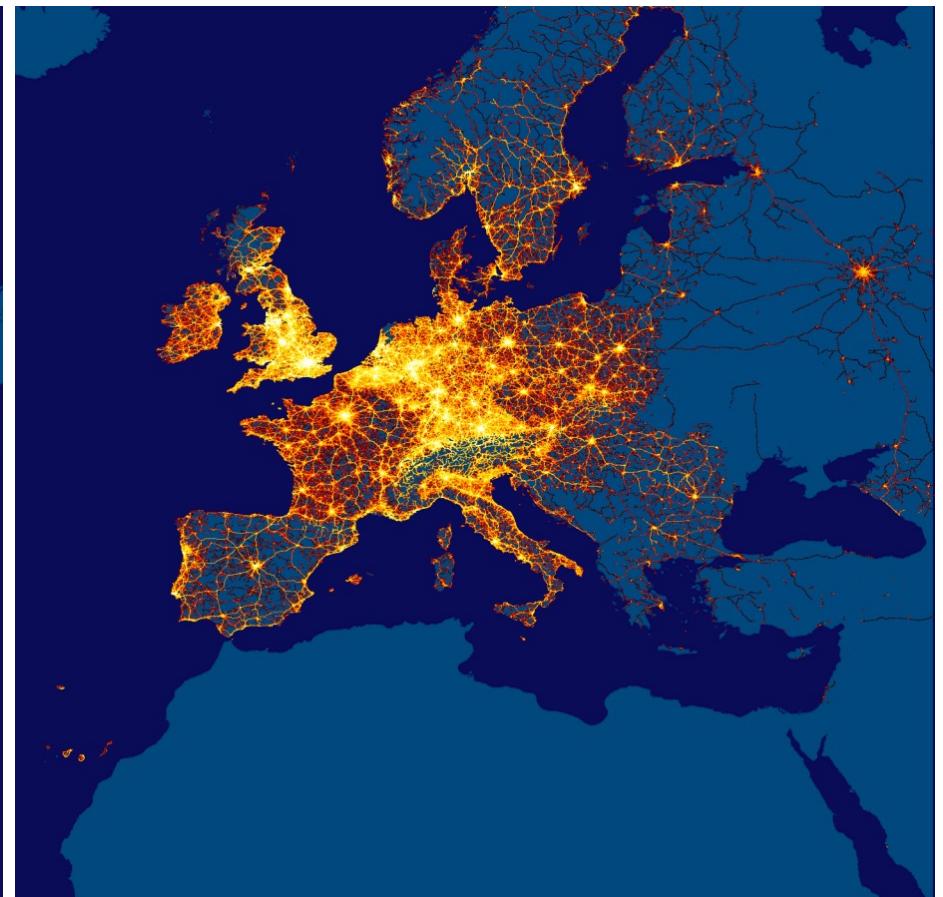
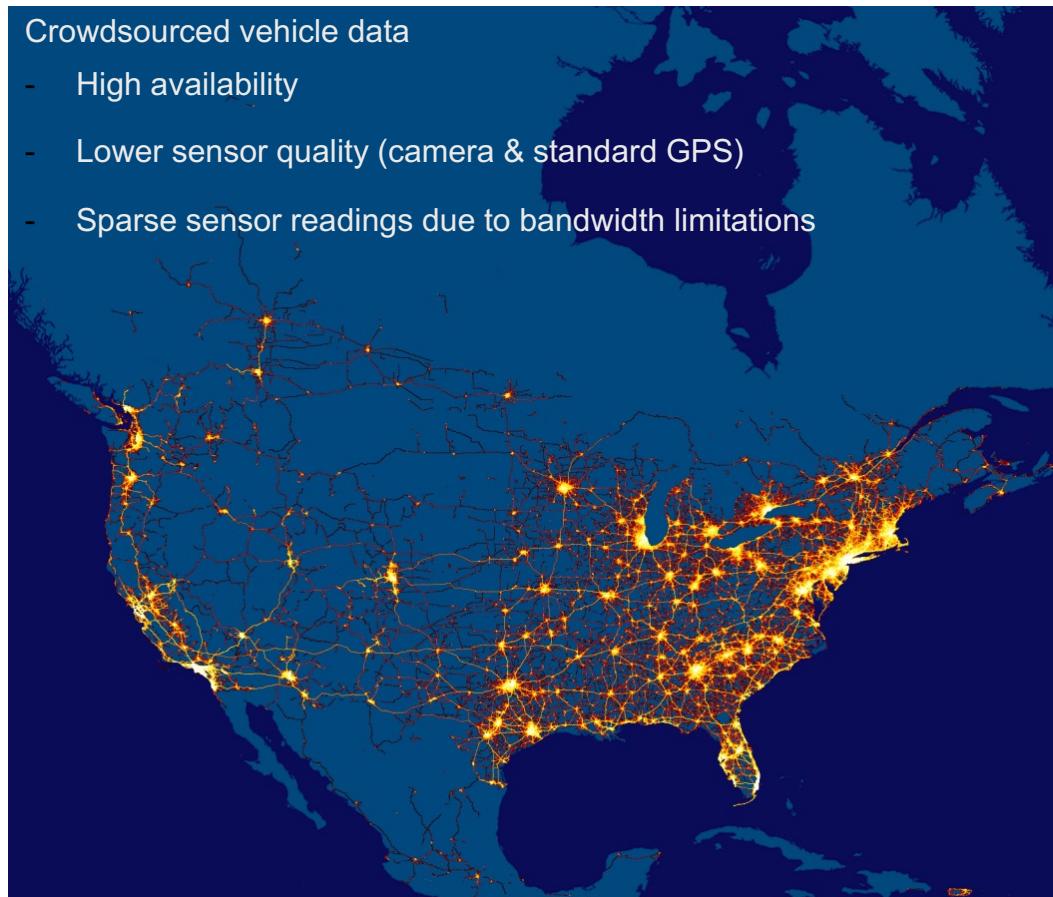
- Expensive sensors
- Small #vehicles
- Small frequency of visits



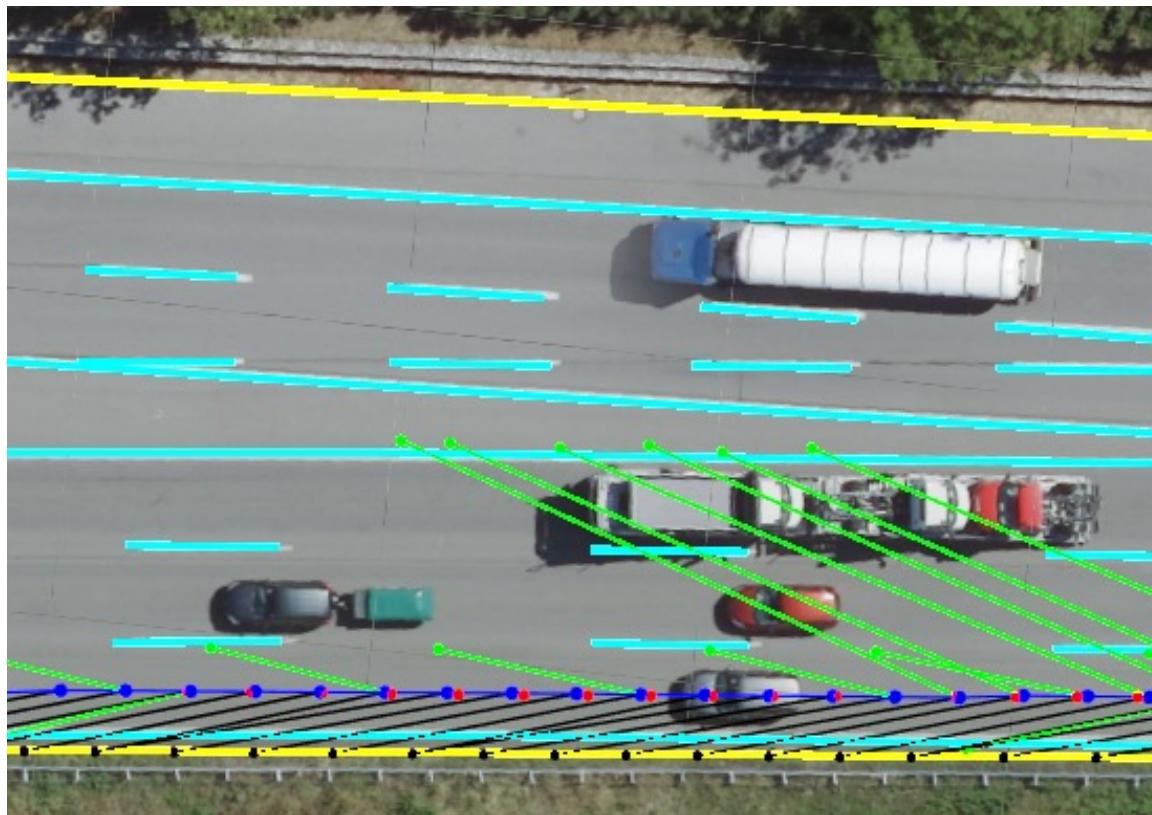
# With Crowdsourced Vehicle Data

Crowdsourced vehicle data

- High availability
- Lower sensor quality (camera & standard GPS)
- Sparse sensor readings due to bandwidth limitations



# HD Map Representation & Measurements



## Crowdsourced data

- Vehicle position
- Odometry measurement
- GNSS measurement
- Lane marking observation
- Road edge observation

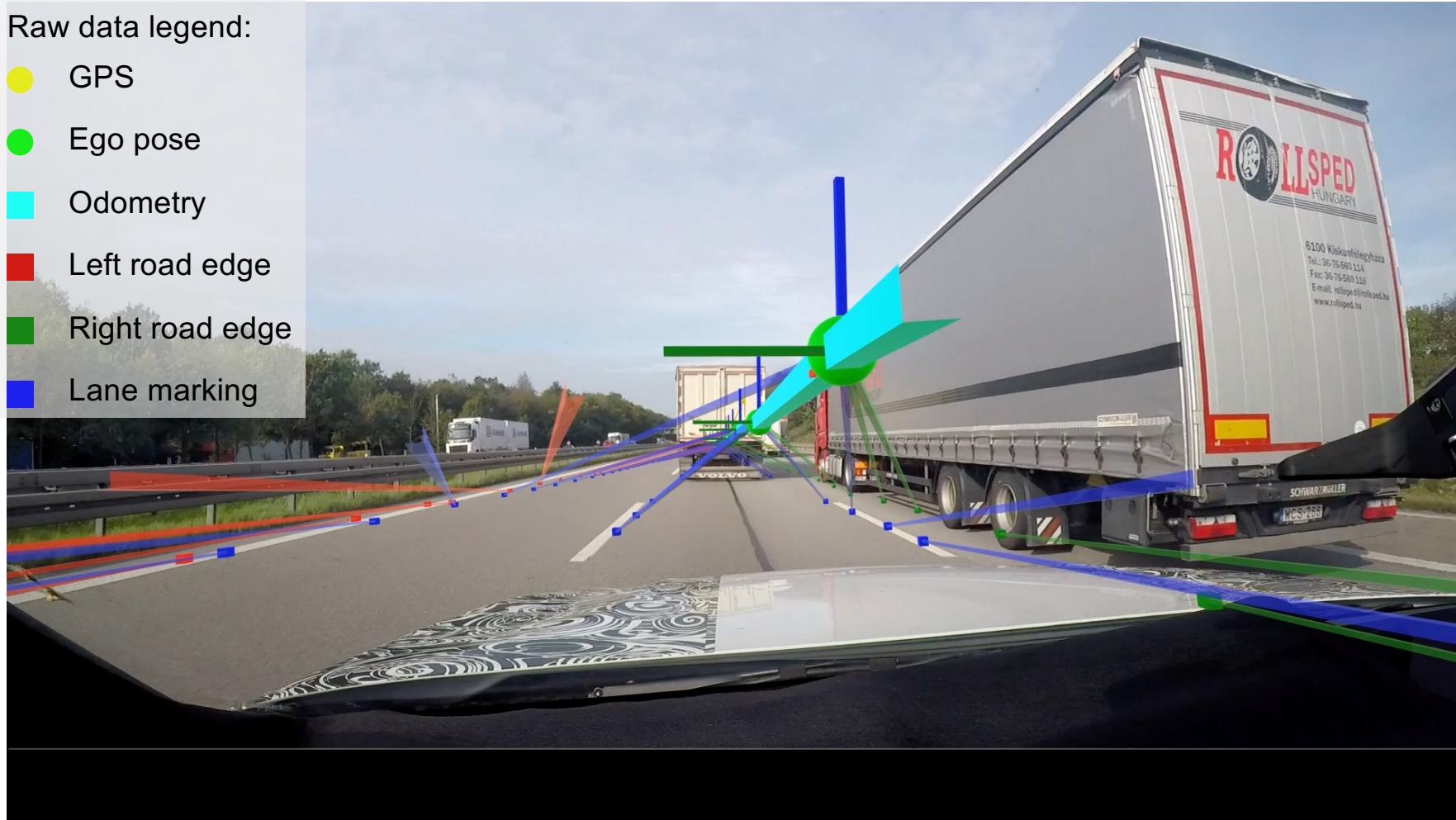
## HD map

- Lane marking
- Road Edge

# Features

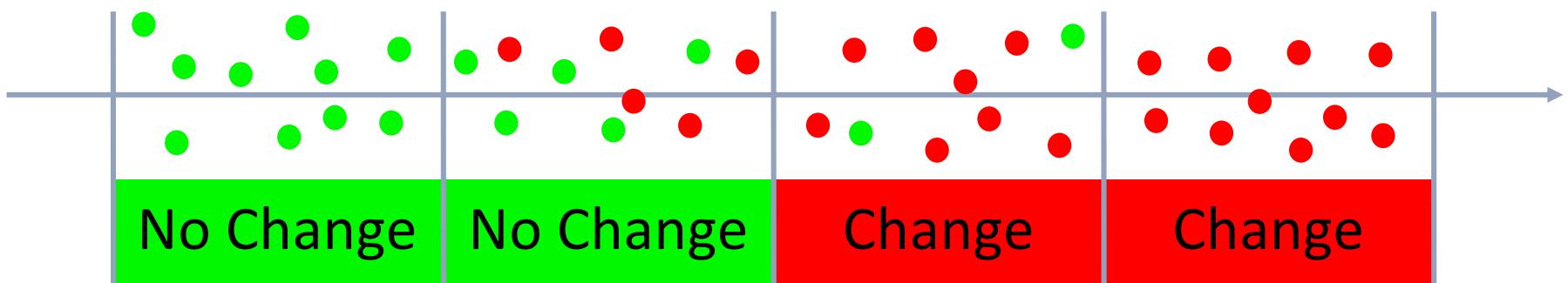
Raw data legend:

- GPS
- Ego pose
- Odometry
- Left road edge
- Right road edge
- Lane marking

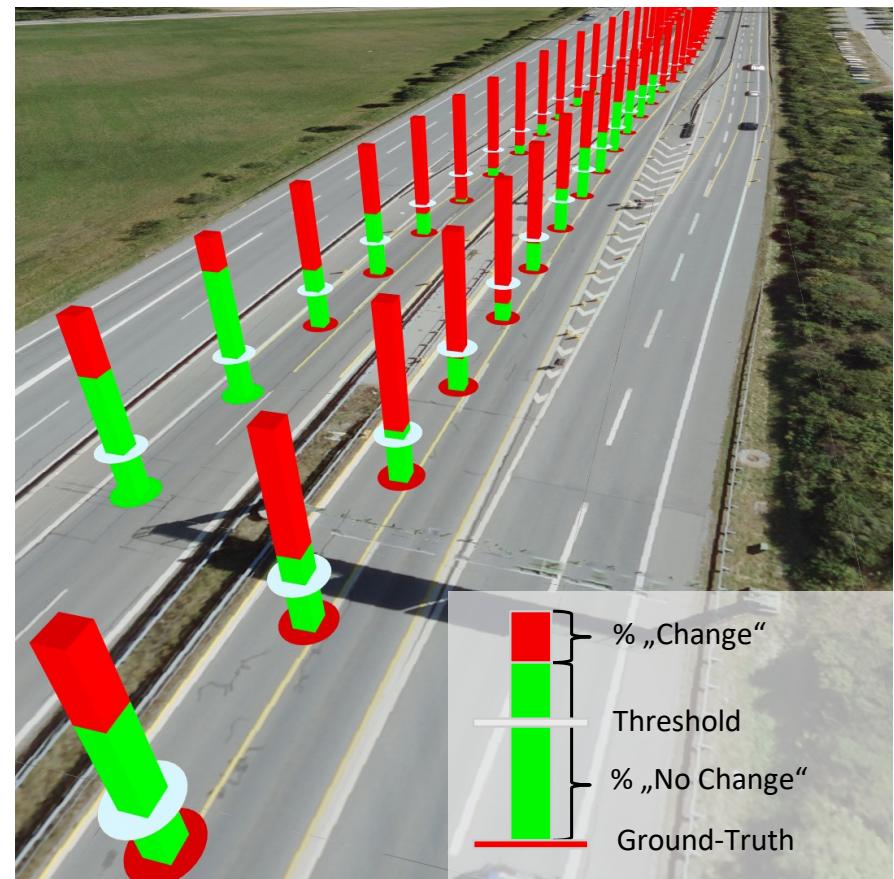
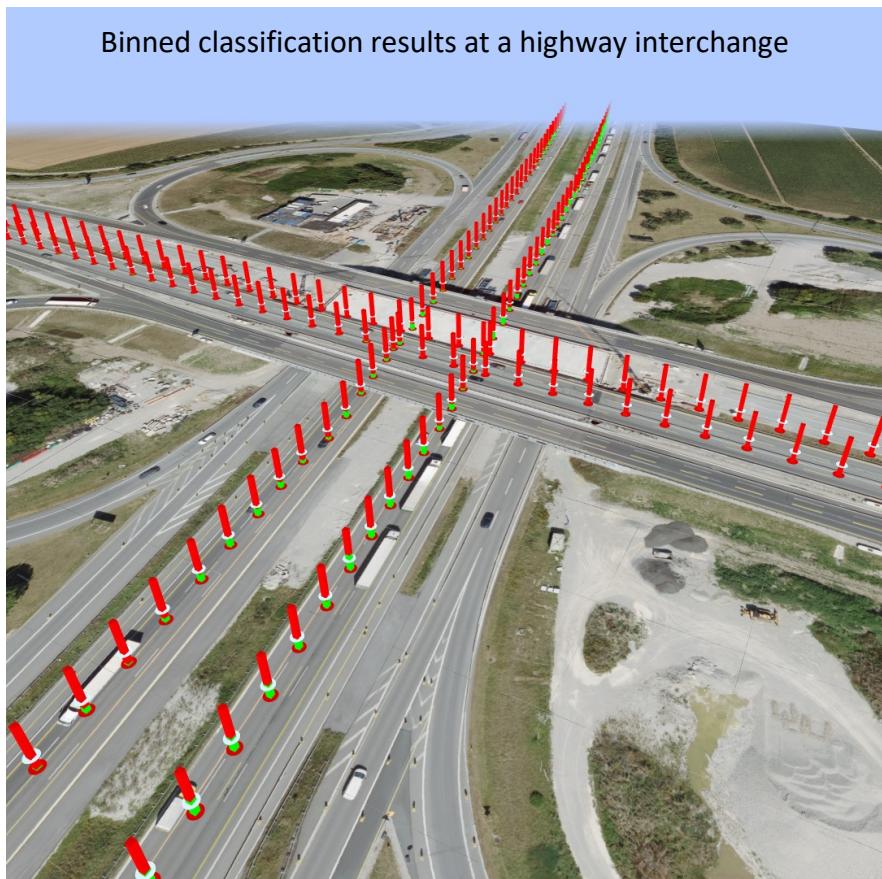


# Particle Filters to Detect Innovation

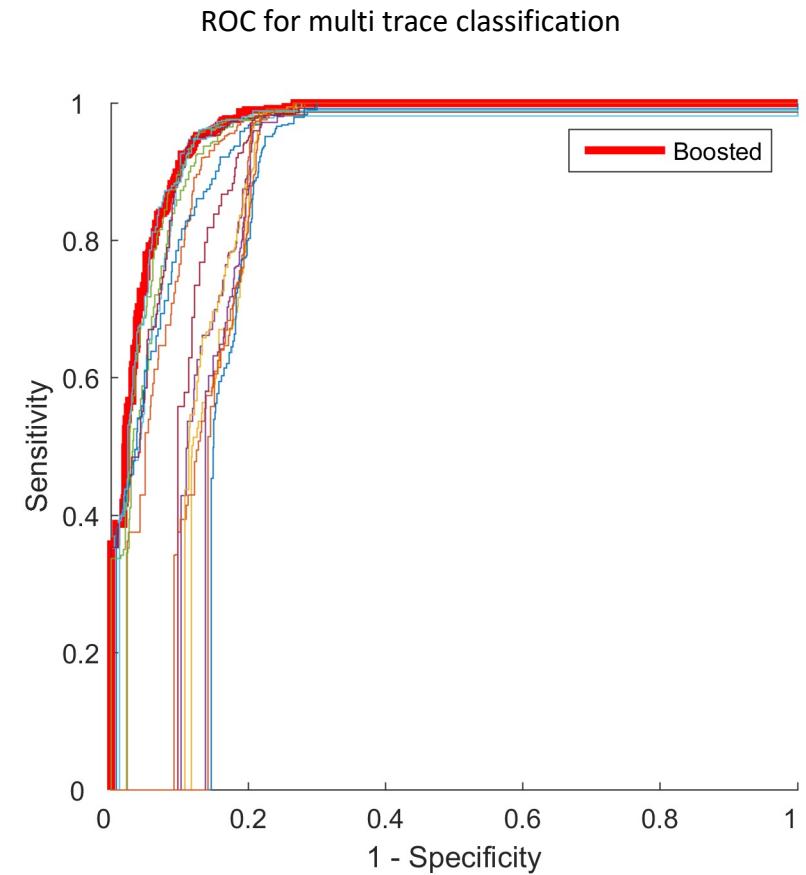
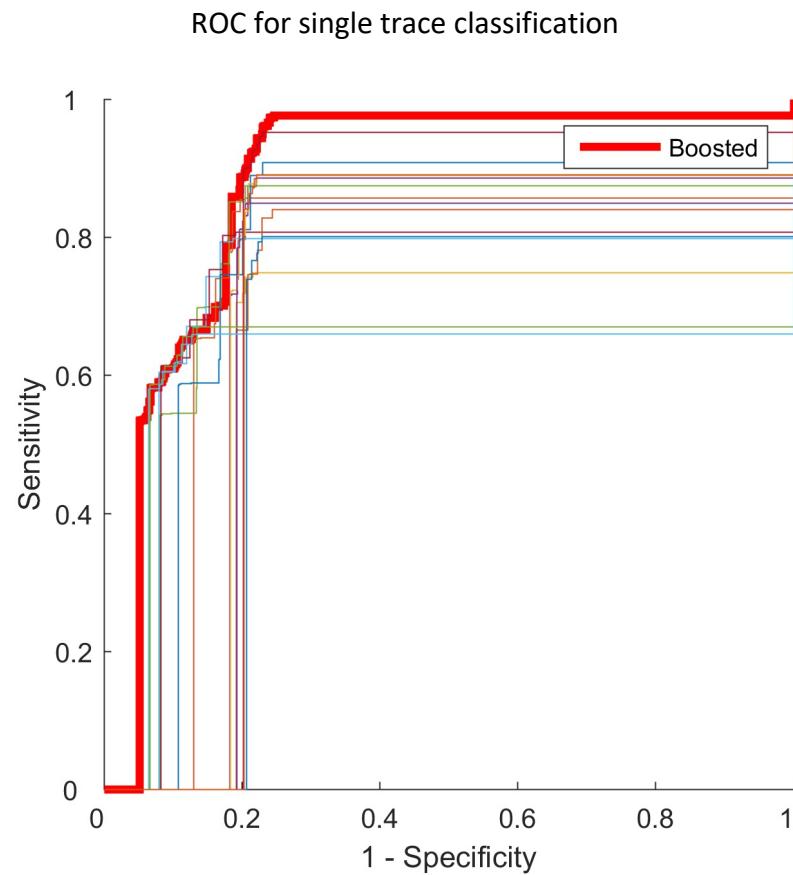
- Odometry
- GNSS
- Lane markings
- Road edges
- Adaboost for classification
- Accumulation of classification results in bins along the road



# Qualitative Results

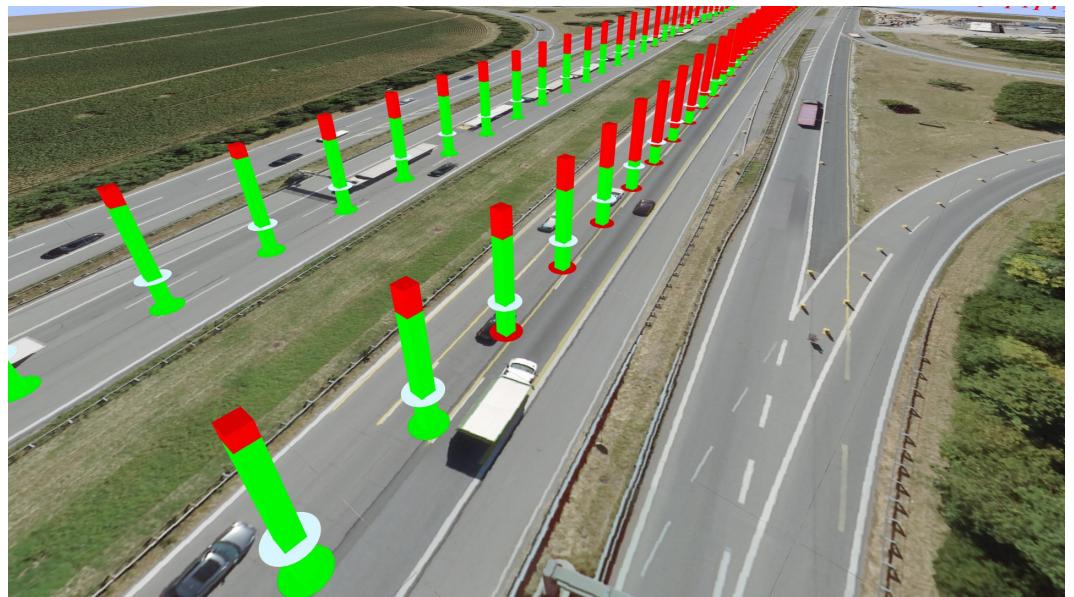


# Results



# Contribution

- New HD map change detection approach based on particle filter localization and boosting
- Using comparably sparse and noisy sensor data
- Crowdsourcing significantly improves single traversal classification with a sensitivity of 98.7% and specificity of 81.2%
- Evaluation showed successful detection of real-world changes in construction sites



# Challenges for Learning Maps on the Fly

- Semantics
- Data association
  - Tracking
  - Traffic light associations
  - ...
- Topology
  - Road structure
  - Intersection topology
  - ...

# Learning Semantics & Dynamics

## EfficientLPS: Efficient LiDAR Panoptic Segmentation

Kshitij Sirohi<sup>\*1</sup>, Rohit Mohan<sup>\*1</sup>, Daniel Büscher<sup>1</sup>, Wolfram Burgard<sup>1,2</sup>, and Abhinav Valada<sup>1</sup>

**Abstract**—Panoptic segmentation of point clouds is a crucial task that enables autonomous vehicles to comprehend their vicinity using their highly accurate and reliable LiDAR sensors. Existing top-down approaches tackle this problem by either combining independent task-specific networks or translating methods from the image domain ignoring the intricacies of LiDAR data and the resulting challenges. In this paper, we present the novel top-down Efficient LiDAR Panoptic Segmentation (EfficientLPS) architecture that addresses multiple challenges in segmenting LiDAR point clouds including distance-dependent sparsity, severe occlusions, large scale-variations, and re-projection errors. EfficientLPS comprises of a novel shared backbone that encodes with strengthened geometric transformation modeling capacity and aggregates semantically rich range-aware multi-scale features. It incorporates new scale-invariant semantic and instance segmentation heads along with the panoptic fusion module which is supervised by our proposed panoptic periphery loss function. Additionally, we formulate a regularized panoptic labeling framework to further improve the performance of EfficientLPS by training on unlabeled data. We benchmark our proposed method on two large-scale LiDAR datasets: nuScenes, for which we also provide ground truth annotations, and SemanticKITTI. Notably, EfficientLPS sets the new state-of-the-art on both these datasets.

**Index Terms**—Scene Understanding, Semantic Segmentation, Instance Segmentation, Panoptic Segmentation.

### I. INTRODUCTION

AUTONOMOUS vehicles are required to operate in challenging urban environments that consist of a wide variety of agents and objects, making comprehensive perception a critical task for robust and safe navigation. Typically, perception tasks are focused on independently reasoning about the semantics of the environment and recognition of object instances. Recently, panoptic segmentation [1] which unifies semantic and instance segmentation has emerged as a popular scene understanding problem that aims to provide a holistic solution. Panoptic segmentation simultaneously segments the scene into ‘stuff’ classes that comprise of background objects or amorphous regions such as road, vegetation, and buildings, as well as ‘thing’ classes that represent distinct foreground objects such as cars, cyclists, and pedestrians. Panoptic segmentation has been extensively studied in the image domain [1]–[4], facilitated by the ordered structure of images being supported by well-researched convolutional networks. However, only a handful of methods have been proposed for panoptic segmentation of LiDAR point clouds [5], [6]. LiDARs have become an indispensable sensor for autonomous vehicles due

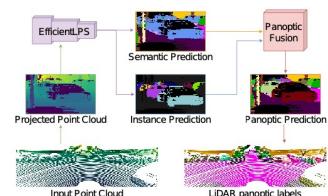


Fig. 1: Overview of the top-down EfficientLPS architecture that consists of a shared backbone to learn spatially-aware features from the projected point cloud and additional heads to learn semantic and instance specific features which are fused in the panoptic fusion module. The network explicitly utilizes the range information in the backbone, semantic head and fusion module to mitigate the problems due to the projection and distance-dependent sparsity of LiDAR point clouds.

to their illumination independence and geometric description of the scene, making scene understanding using LiDAR point clouds an essential capability. However, the typical unordered, sparse, and irregular structure of point clouds pose several unique challenges.

To this end deep learning methods that rely on grid based convolutions to address these challenges typically follow two different directions. They either project the point cloud into the 3D voxel space and employ 3D convolutions on them [7], [8], or they project the point cloud into the 2D space [6], [9], [10] and employ the well-researched 2D Convolutional Neural Networks (CNNs). While voxel-based method achieve high accuracy, they are computationally more expensive and require substantial memory to store the voxelized point clouds. The 2D projection based methods on the other hand, yield a more denser representation and require comparatively lesser computational resources, but they suffer from information loss during projection, blurry CNN outputs, and incorrect label assignment to the occluded points during re-projection. Therefore, there is a need to bridge this gap with a method that has the advantages of fast and memory-efficient 2D convolutions while mitigating the problems due to the projection.

In this work, we present the novel Efficient LiDAR Panoptic Segmentation (EfficientLPS) architecture that effectively addresses the aforementioned challenges by employing a 2D CNN for the task while explicitly utilizing the unique 3D information provided by point clouds. EfficientLPS consists of a shared backbone comprising our novel Proximity Convolution

**EfficientLPS: Efficient LiDAR Panoptic Segmentation.** Kshitij Sirohi, Rohit Mohan, Daniel Büscher, Wolfram Burgard, Abhinav Valada

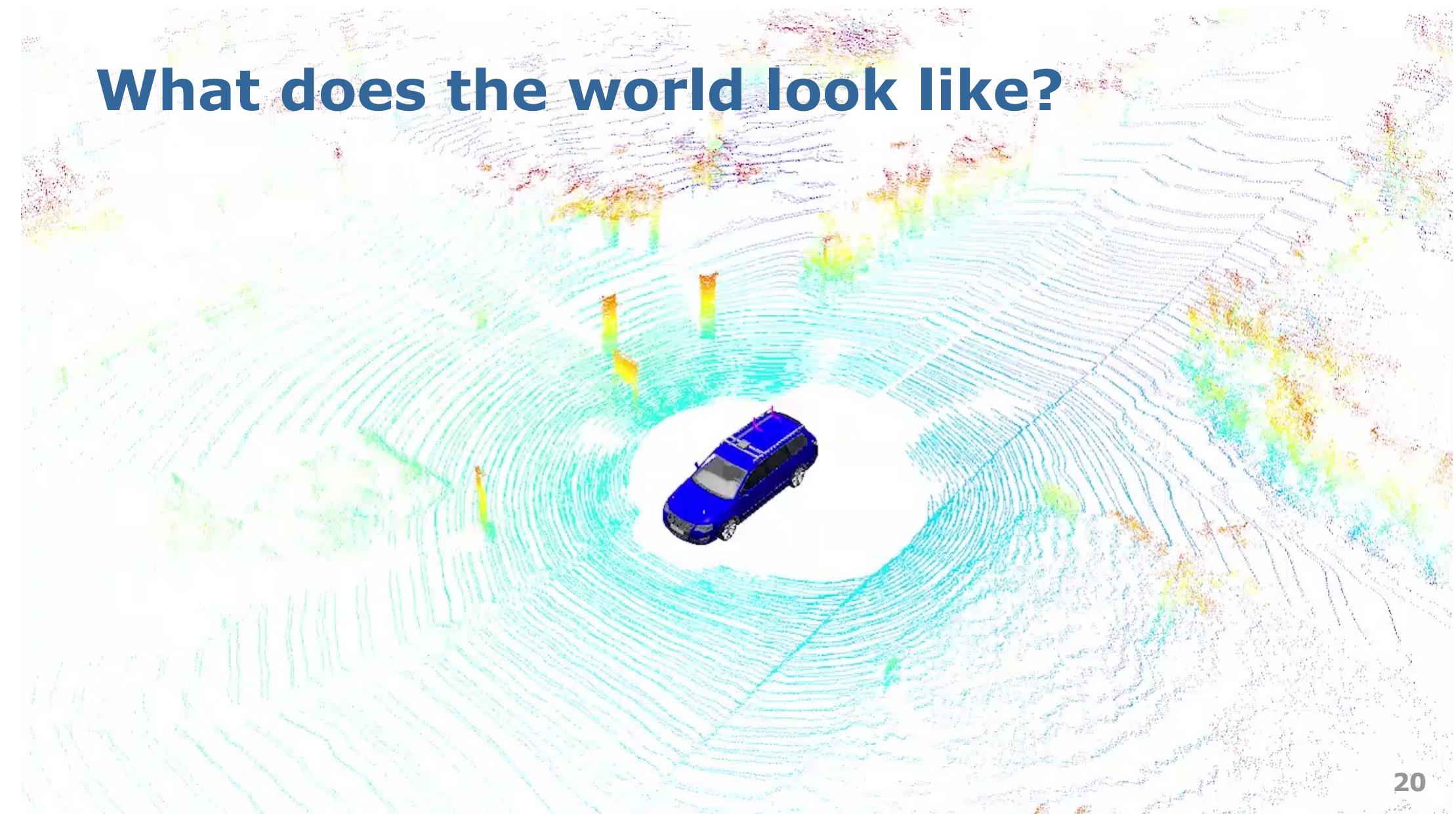


<sup>\*</sup>These authors contributed equally.

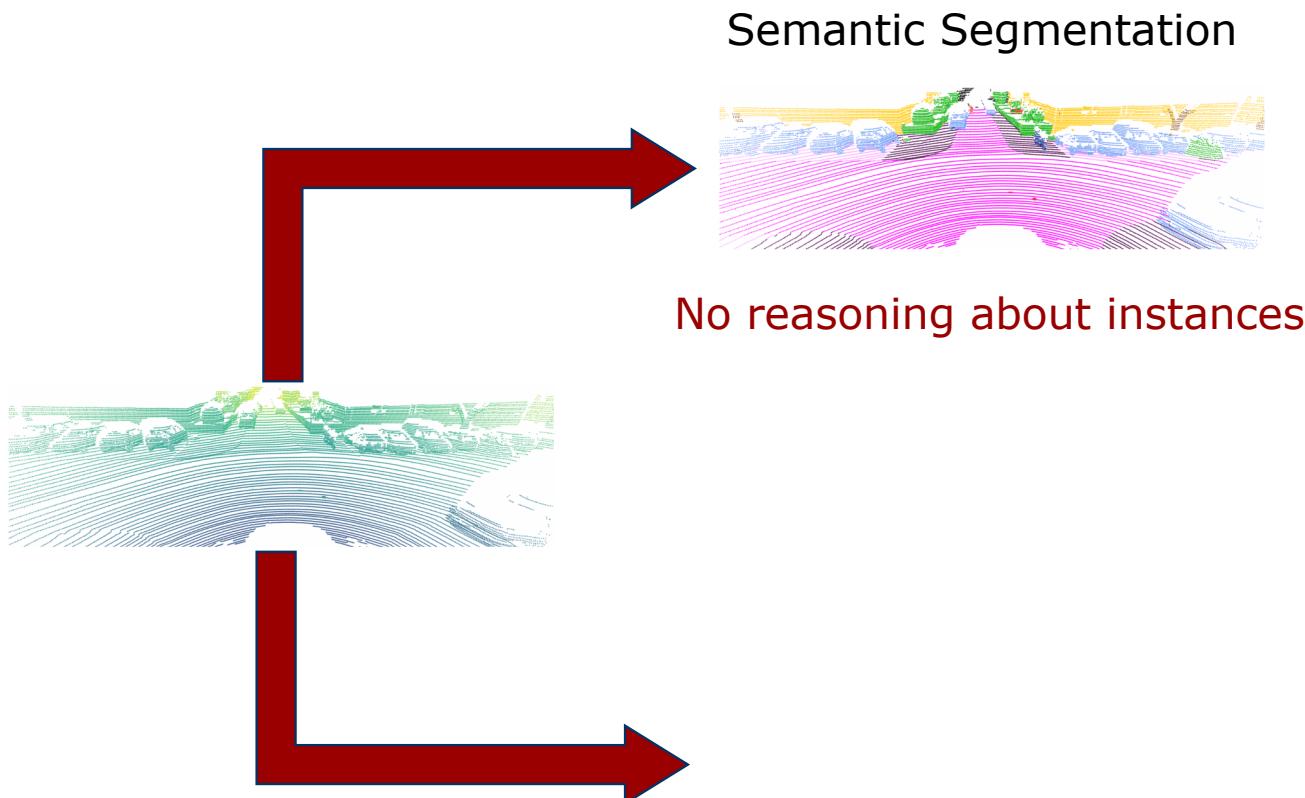
<sup>1</sup>Department of Computer Science, University of Freiburg, Germany

<sup>2</sup>Toyota Research Institute, Los Altos, USA.

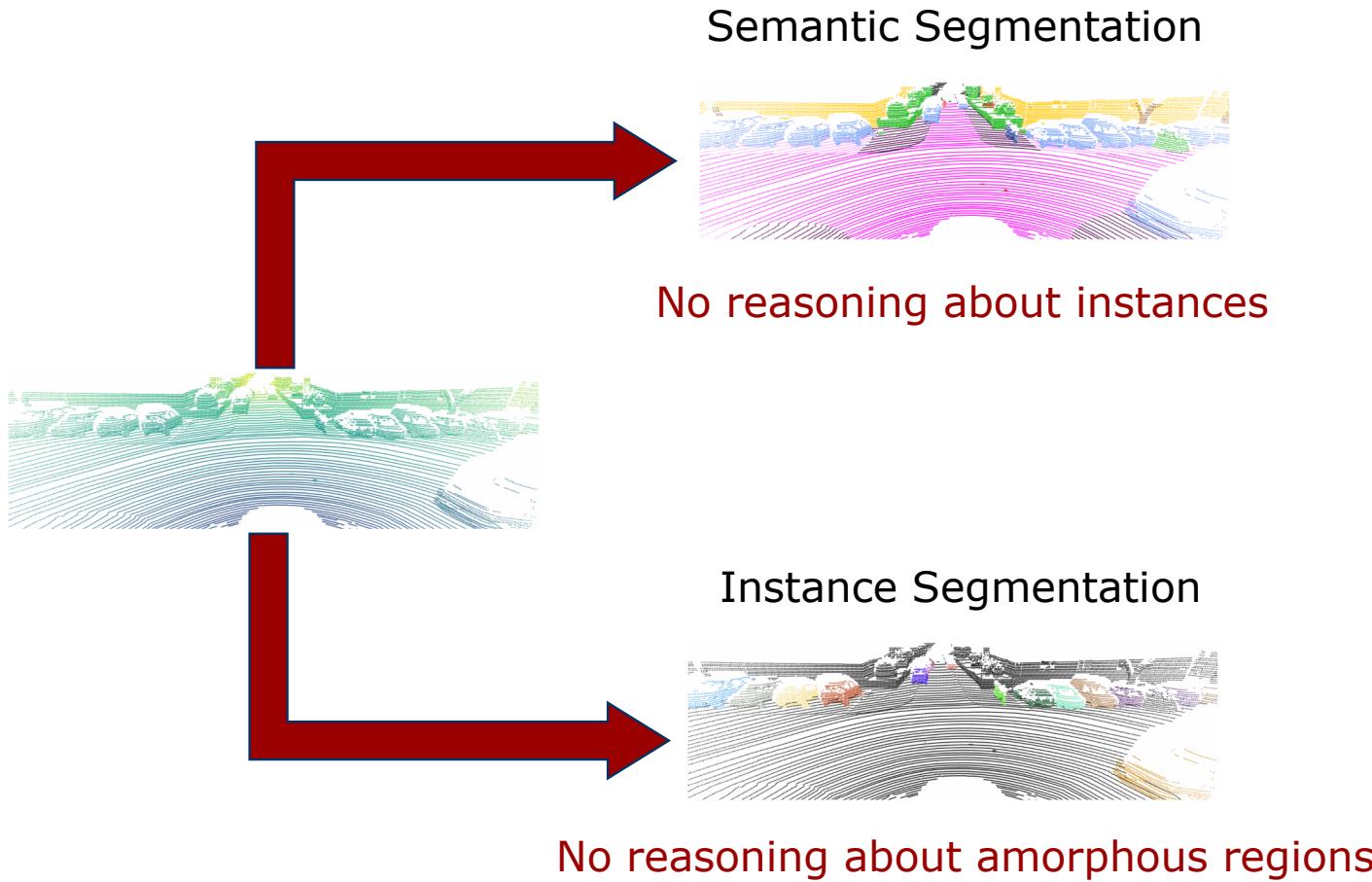
# What does the world look like?



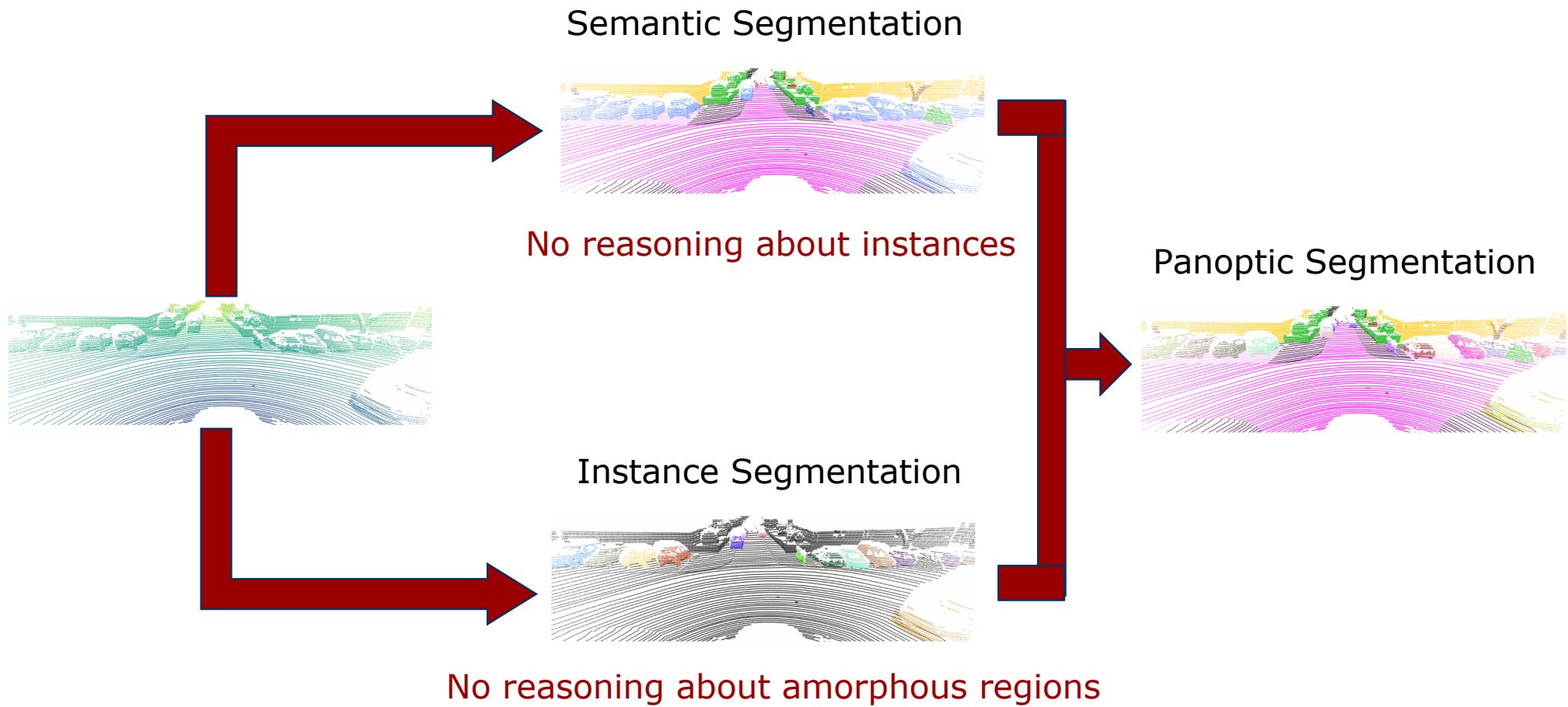
# Semantic Tasks



# Semantic Tasks

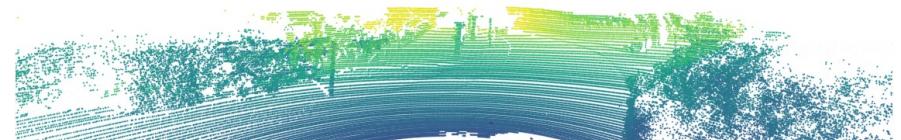


# Semantic Tasks



# Challenges in using LiDAR data

- LiDARs provide an accurate illumination independent geometric description
- Unordered and irregular structure of point clouds
- Distance dependent sparsity
- Memory and time intensive 3D convolutions
- Information loss during point cloud projection into a 2D space



# LiDAR Scan Projection



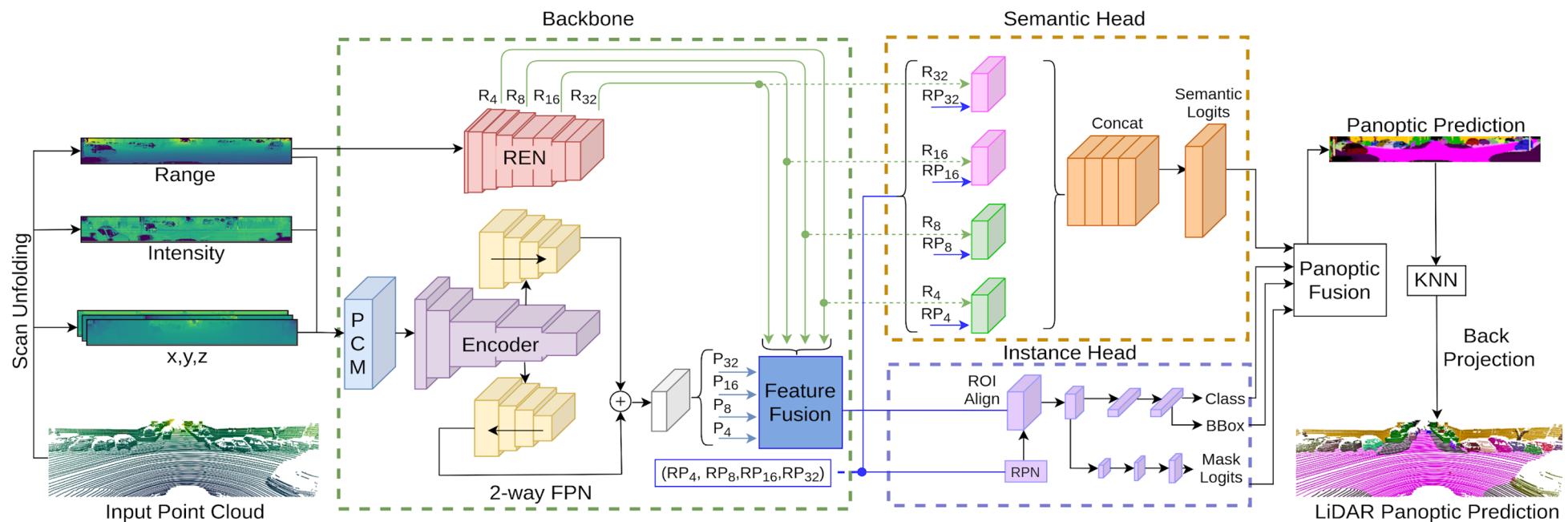
Spherical Projection



Scan Unfolding Projection

- Spherical projection (Milioto et al., 2019) suffers from systematic point occlusions
- Scan unfolding projection (Triess et al., 2020) reduces systematic occlusion artifacts
- Improves performance

# EfficientLPS Architecture



- Scan unfolding projection
- Backbone: PCM + Encoder + REN + 2-way FPN
- Semantic Head, Instance Head, Panoptic Fusion Module
- Reprojection into 3D using kNNs

# Proximity Convolution Module

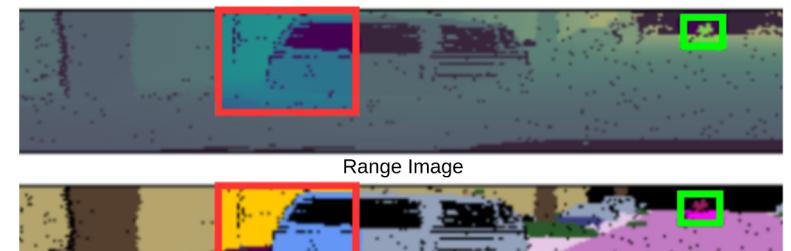
- Exploits range information to augment spatial sampling locations
- Enhances the geometric modeling capabilities of convolutions
- Effectively captures contextual information from sparse points



Standard Convolution 3x3



3x3 Convolution on range filtered neighbors

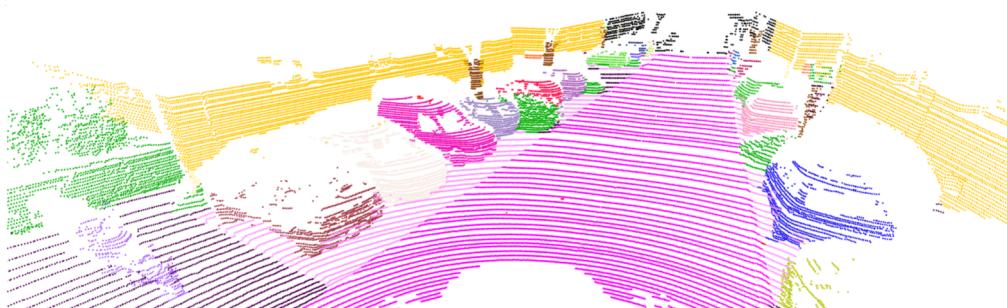


Labeled Segments

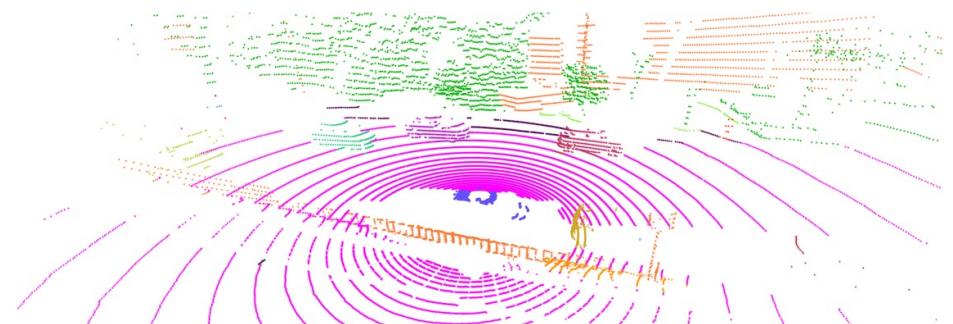
Standard Convolution

Proximity Convolution

# Evaluation Datasets



SemanticKITTI



nuScenes

# Benchmarking Results

- Ranked #1 on SemanticKITTI and nuScenes datasets for panoptic segmentation
- Novel regularized pseudo labeling framework further improves the performance

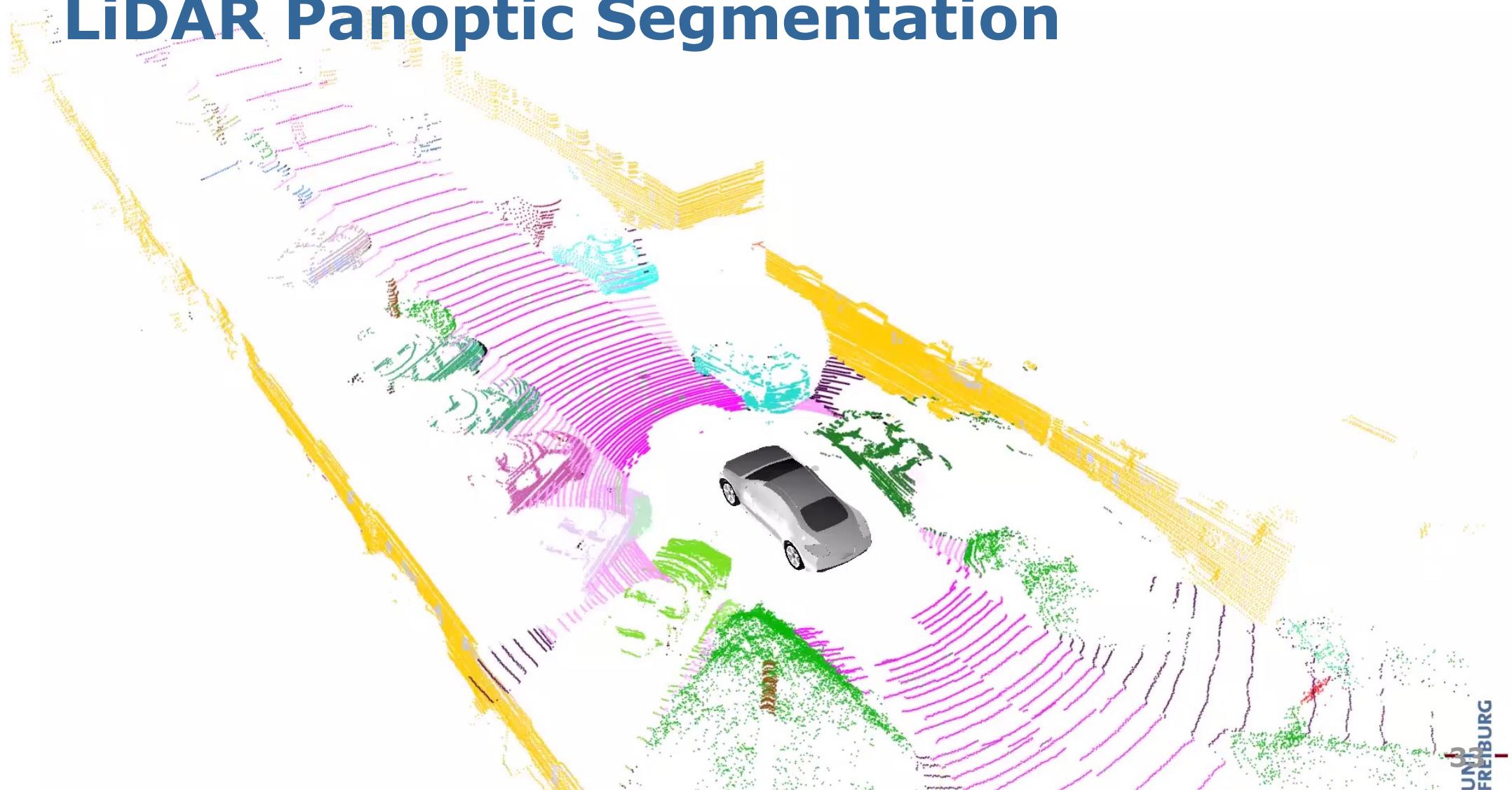
SemanticKITTI Leaderboard

Method	PQ
RangeNet++ [23] + PointPillars [43]	37.1
KPConv [21] + PointPillars [43]	44.5
LPSAD [6]	38.0
PanopticTrackNet [10]	43.1
Panoster [5]	52.7
EfficientLPS (ours)	<b>57.4</b>

nuScenes Validation Set

Method	PQ
RangeNet++ [23] + Mask R-CNN [38]	45.2
PanopticTrackNet [10]	50.0
KPConv [21] + Mask R-CNN [38]	50.1
EfficientLPS (ours)	<b>59.2</b>

# LiDAR Panoptic Segmentation



# How Do We Go Beyond?

Panoptic Segmentation



Frame-by frame, inconsistent over time

# How Do We Go Beyond?

Panoptic Segmentation



Frame-by frame, inconsistent over time

Multi-Object Tracking



Bbox tracking performance is saturating

# How Do We Go Beyond?

Panoptic Segmentation



Multi-Object Tracking



Multi-Object Panoptic Tracking



Instance IDs are  
temporally associated

# Multi-Object Panoptic Tracking (MOPT)

## MOPT: Multi-Object Panoptic Tracking

Juana Valeria Hurtado Rohit Mohan Wolfram Burgard Abhinav Valada  
University of Freiburg



Figure 1: MOPT output overlaid on the image LiDAR input from the Virtual KITTI 2 (top row) and SemanticKITTI (bottom row) datasets. MOPT unifies semantic segmentation, instance segmentation and multi-object tracking to yield segmentation of stuff classes and segmentation of thing classes with temporally consistent instance IDs. Observe that the tracked thing instances retain their color-code temporally in subsequent timesteps.

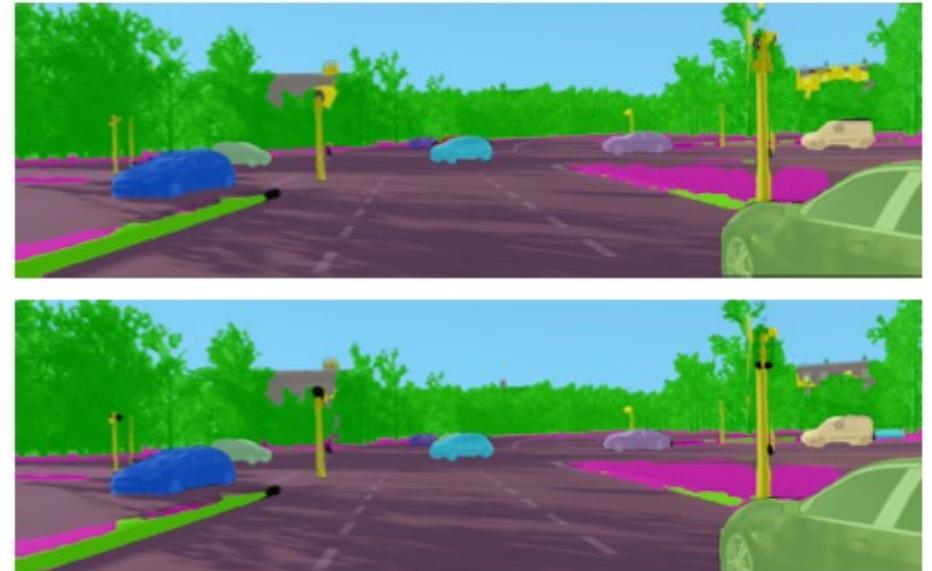
### Abstract

Comprehensive understanding of dynamic scenes is a critical prerequisite for intelligent robots to autonomously operate in their environment. Research in this domain, which encompasses diverse perception problems, has primarily been focused on addressing specific tasks individually rather than modeling the ability to understand dynamic scenes holistically. In this paper, we introduce a novel perception task denoted as multi-object panoptic tracking (MOPT), which unifies the conventionally disjoint tasks of semantic segmentation, instance segmentation, and multi-object tracking. MOPT allows for exploiting pixel-level semantic information of thing and stuff classes, temporal coherence, and pixel-level associations over time, for the mutual benefit of each of the individual sub-problems. To facilitate quantitative evaluations of MOPT in a unified manner, we propose the soft panoptic tracking quality (sPTQ) metric. As a first step towards addressing this task, we propose the novel PanopticTrackNet architecture that builds upon the state-of-the-art top-down panoptic segmentation network EfficientPS by adding a new tracking head to simultaneously learn all sub-tasks in an end-to-end manner. Additionally, we present several strong baselines that combine predictions from state-of-the-art panoptic segmentation and multi-object tracking models for comparison. We present extensive quantitative and qualitative evaluations of both vision-based and LiDAR-based MOPT that demonstrate encouraging results.

### 1. Introduction

Comprehensive scene understanding is a critical challenge that requires tackling multiple tasks simultaneously

MOPT: Multi-Object Panoptic Tracking.  
Juana Valeria Hurtado, Rohit Mohan,  
Wolfram Burgard, Abhinav Valada

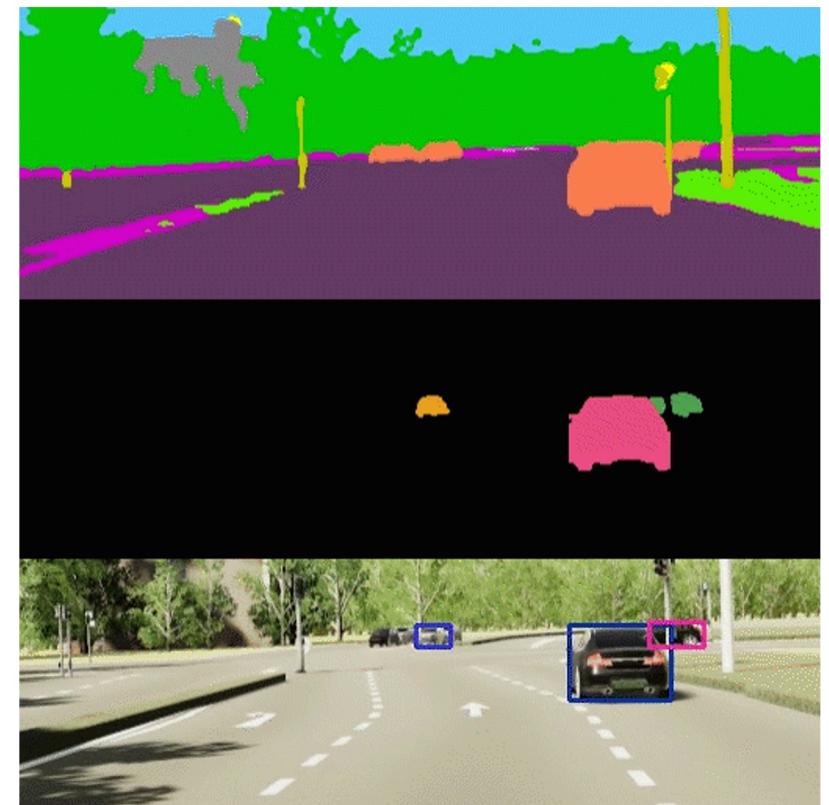


to detect, localize, and identify the scene elements as well as to understand the occurring context, dynamics, and relationships. These fundamental scene comprehension tasks are crucial enablers of several diverse applications [21, 26, 22, 25, 30] including autonomous driving, robot navigation, augmented reality and remote sensing. Typically, these problems have been addressed by solving distinct perception tasks individually, i.e., image\pointcloud recognition, object detection and classification, semantic segmentation, instance segmentation, and tracking. The state of the art in these tasks have been significantly advanced since the advent of deep learning approaches, however their performance is no longer increasing at the same groundbreaking pace [18]. Moreover, as most of these tasks are required to be performed simultaneously in real-world applications, the scalability of employing several individual models is becoming a limiting factor. In order to mitigate this emerging problem, recent works [8, 27, 28, 17] have made efforts to exploit common characteristics of some of these tasks by jointly modeling them in a coherent manner.

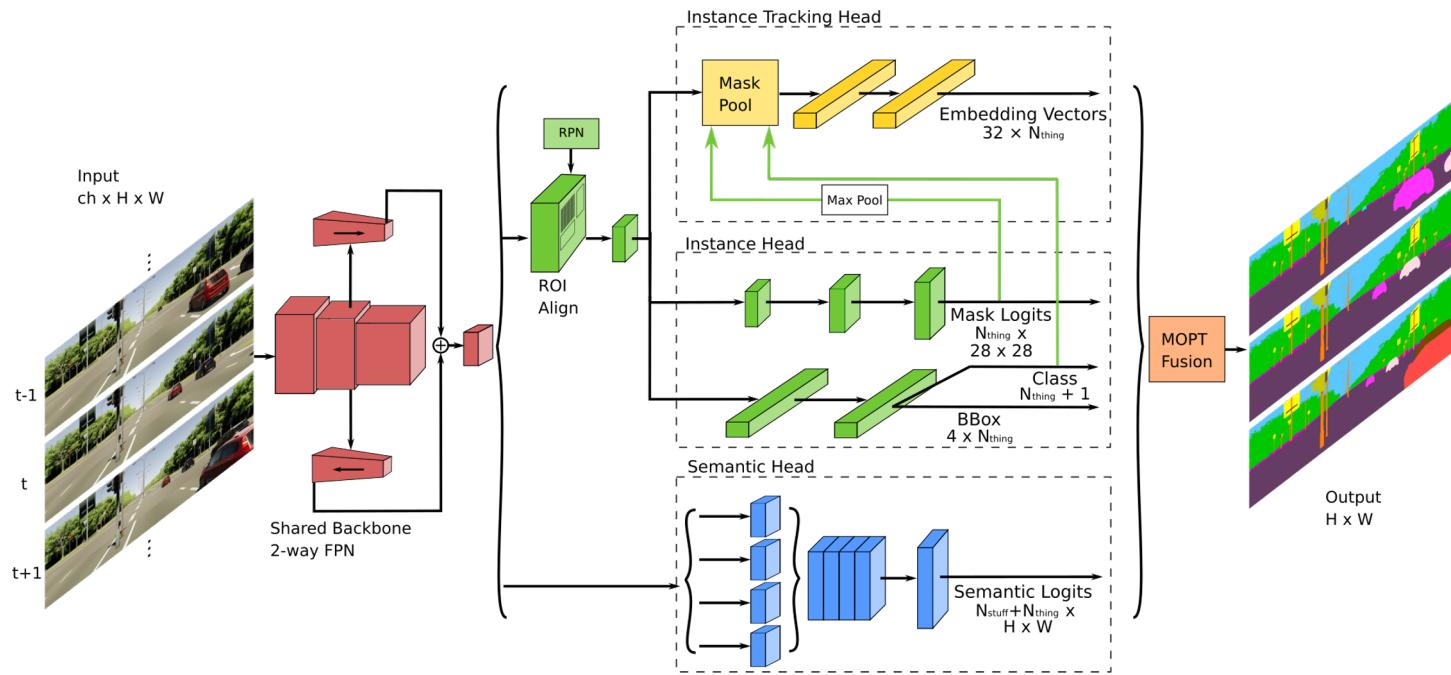
Two such complementary tasks have gained a substantial amount of interest in the last few years due to the availability of public datasets [1, 3] and widely adopted benchmarks [28, 1]. The first of which is panoptic segmentation [8] that unifies semantic segmentation of stuff classes which consist of amorphous regions and instance segmentation of thing classes which consist of countable objects. While the second task is Multi-Object Tracking and Segmentation (MOTS) [28] which extends multi-object tracking to the pixel level by unifying instance segmentation of thing classes. Since the introduction of these tasks, considerable advances have been made in both panoptic segmentation [7, 31, 15, 14] and MOTS [16, 12, 11, 29] which

# MOPT: Multi-Object Panoptic Tracking

- Encourage holistic modeling of dynamic scenes
- Unifies semantic segmentation, instance segmentation, multi-object tracking
- Exploit temporal consistency for learning better semantics
- Efficient and scalable model for dynamic scene understanding

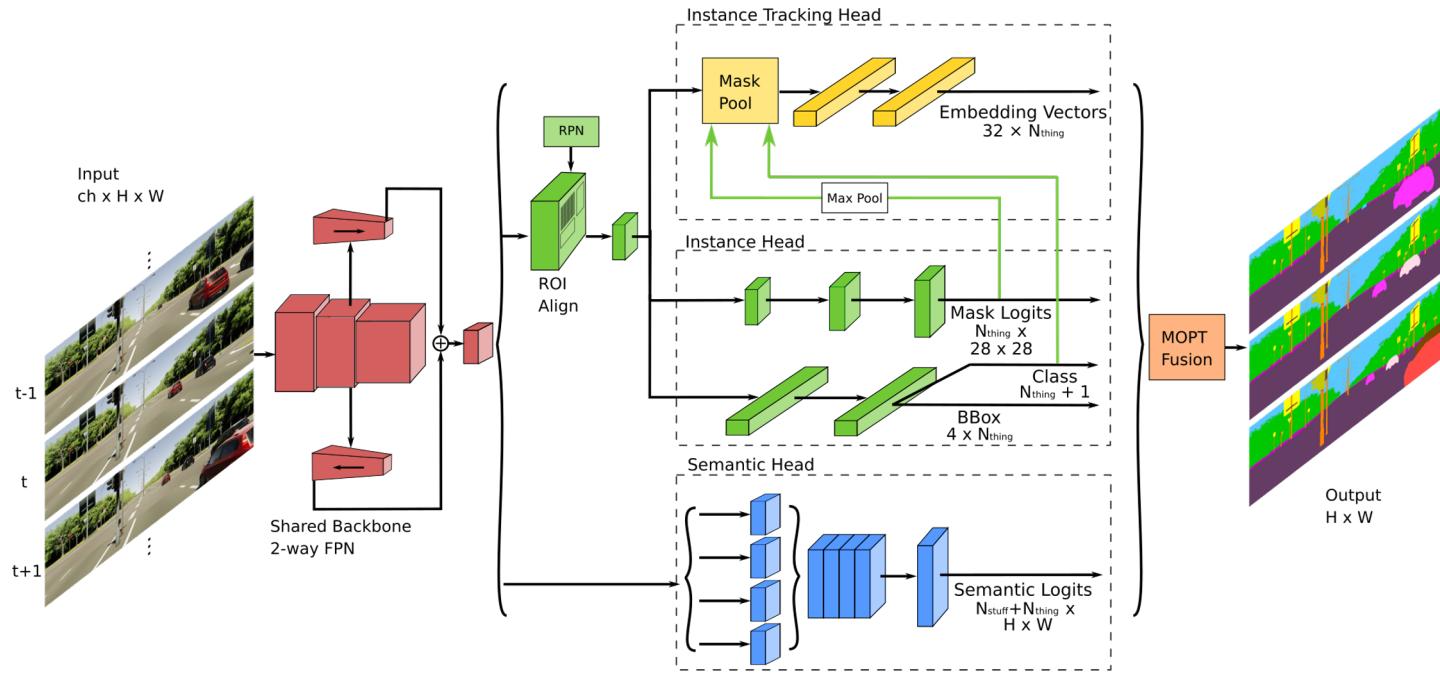


# PanopticTrackNet Architecture



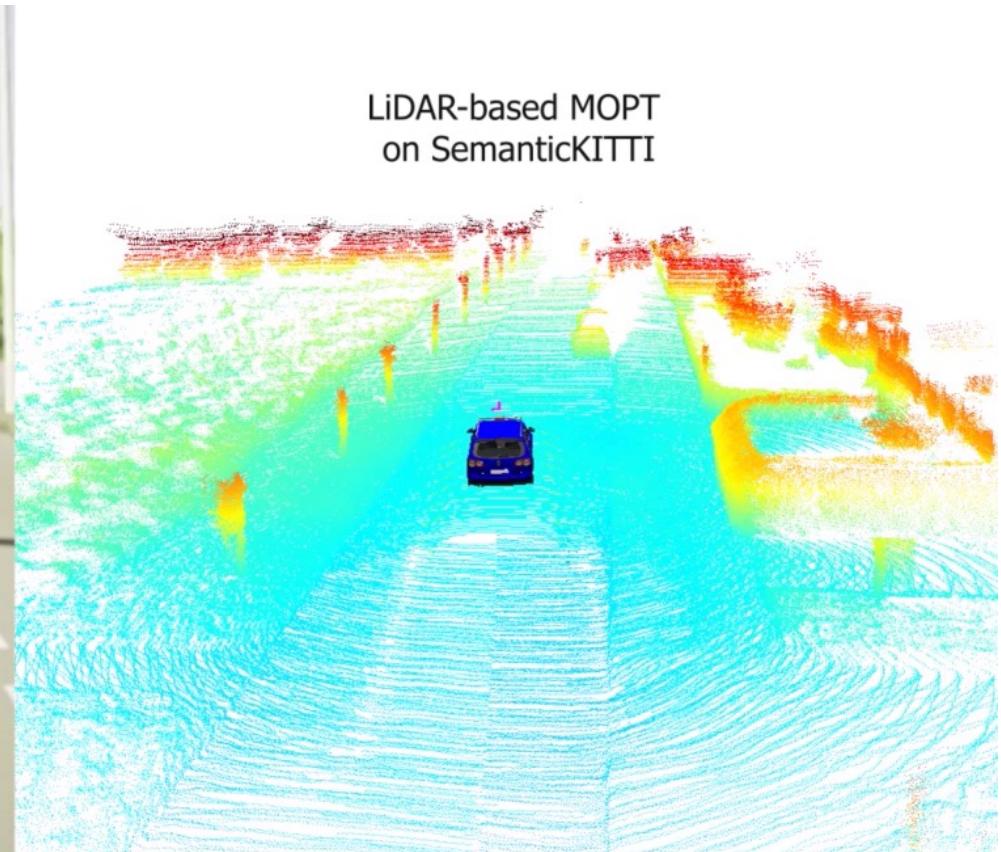
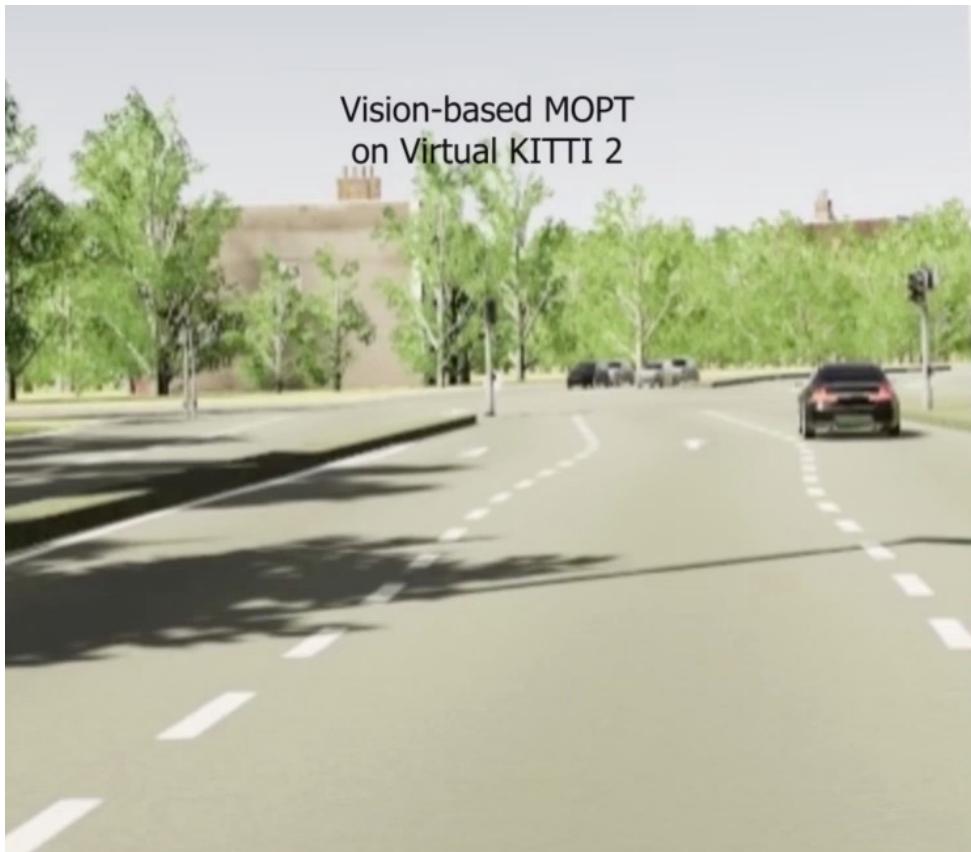
- Built upon our EfficientPS architecture
- Three task specific heads for semantic segmentation, instance segmentation, and multi-object tracking

# PanopticTrackNet Architecture



We adaptively fuse task-specific outputs in our fusion module to yield the panoptic output with temporally tracked instances

# Evaluation Datasets



# Results: Vision-Based MOPT

Network	sPTQ (%)	PTQ (%)	Params. (M)	FLOPs (B)	Time (ms)
Seamless [ 1 ] + Track R-CNN [ 2 ]	45.66	45.28	79.91	273.96	115
EfficientPS [ 3 ] + Track R-CNN [ 2 ]	46.68	46.3	69.91	232.80	115
EfficientPS [ 3 ] + MaskTrack R-CNN [ 4 ]	46.17	45.99	120.57	224.43	117
PanopticTrackNet (ours)	<b>47.27</b>	<b>46.67</b>	<b>45.08</b>	<b>167.40</b>	<b>114</b>

# Results: LiDAR-Based MOPT

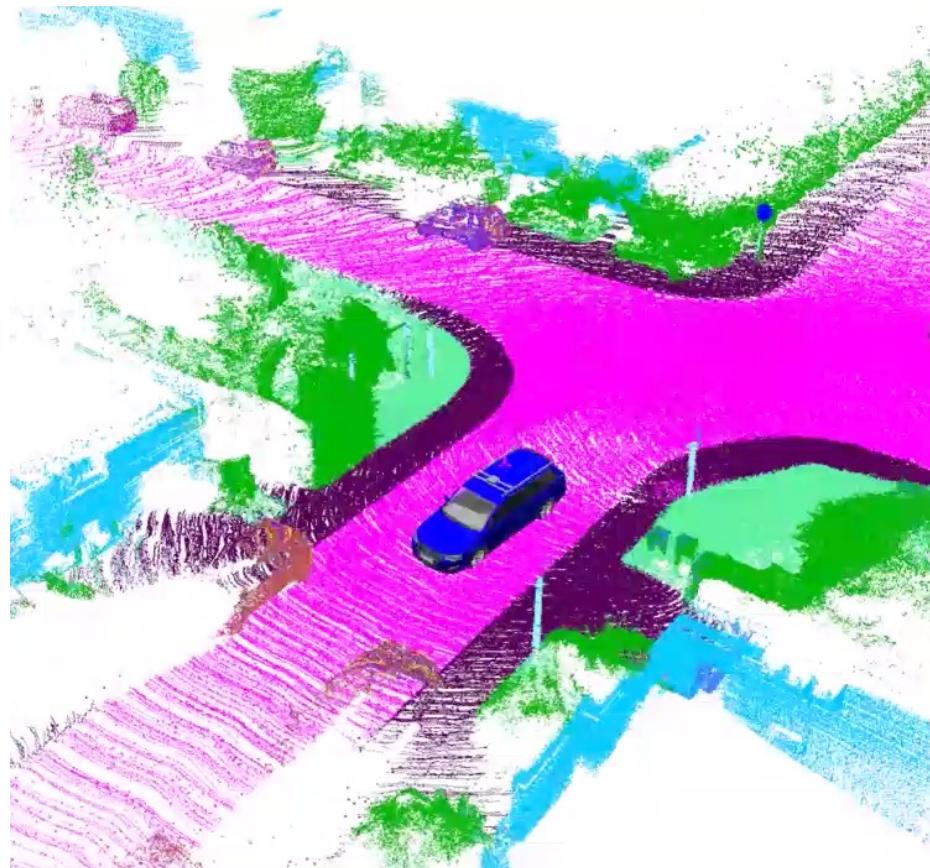
Network	sPTQ (%)	PTQ (%)	Params. (M)	FLOPs (B)	Time (ms)
RangeNet++ [ 1 ] + PointPillars [2] + Track R-CNN [ 3 ]	42.22	41.94	110.74	695.51	409
KPConv [ 4 ] + PointPillars [2] + Track R-CNN [ 3 ]	46.04	45.50	89.96	438.34	514
EfficientPS [ 5 ] + Track R-CNN [ 3 ]	44.50	43.96	79.02	379.73	148
EfficientPS [ 5 ] + MaskTrack R-CNN [ 6 ]	44.03	43.72	120.64	445.90	151
PanopticTrackNet (ours)	<b>48.23</b>	<b>47.89</b>	<b>45.13</b>	<b>300.81</b>	<b>146</b>

# Qualitative Results: Vision-Based MOPT

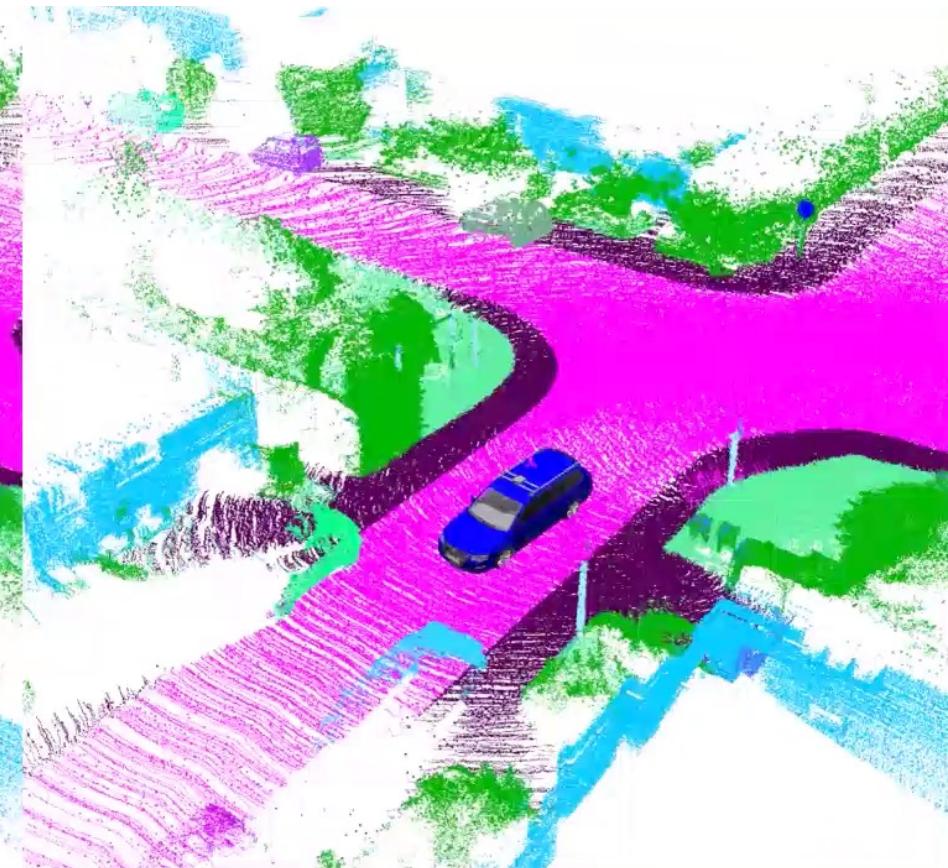


# Results: LiDAR Panoptic Tracking

Panoptic Segmentation



MOPT



# Results: LiDAR Panoptic Tracking

# Robust Semantic Segmentation by Domain Adaption

HeatNet: Bridging the Day-Night Domain Gap in Semantic Segmentation with Thermal Images

Johan Vertens\*, Jannik Zürn\*, and Wolfram Burgard

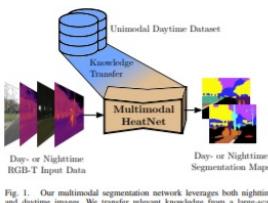
**Abstract**—The majority of learning-based semantic segmentation methods are optimized for daytime scenarios and favorable lighting conditions. Real-world driving scenarios, however, often involve environmental conditions such as nighttime illumination or glare which remains a challenge for existing approaches. In this work, we propose a multimodal semantic segmentation model that can be applied during daytime and nighttime. To this end, we use RGB images to leverage thermal images to make our network significantly more robust. We avoid the expensive annotation of nighttime images by leveraging an existing daytime RGB-dataset and propose a teacher-student training approach that transfers the dataset's knowledge to the nighttime domain. We also propose a domain adaptation method to align the learned feature spaces across the domains and propose a novel two-stage training scheme. Furthermore, due to a lack of thermal data for autonomous driving, we need datasets containing over 20,000 time-synchronized and aligned RGB-thermal image pairs. In this context, we also present a novel target-less calibration method that allows for automatic robust extrinsic and intrinsic thermal camera calibration. Among others, we empirically demonstrate and show state-of-the-art results for nighttime semantic segmentation.

## I. INTRODUCTION

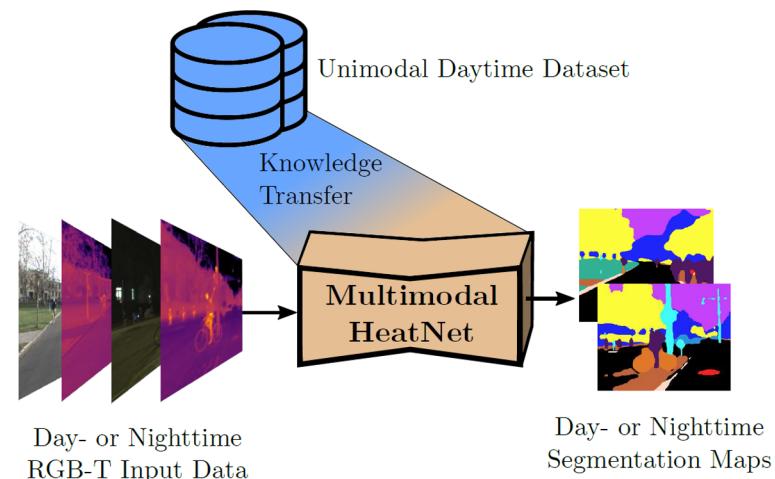
Robust and accurate semantic segmentation of urban scenes is one of the enabling technologies for autonomous driving in complex and cluttered driving scenarios. Recent years have shown great progress in RGB image segmentation for autonomous driving [36], [5], which were predominantly demonstrated in favorable daytime illumination conditions. On the other hand, semantic segmentation of nighttime benchmark datasets [5], [18], these models tend to generalize poorly to adverse weather conditions and low illumination levels present at nighttime. This constraint becomes especially apparent in rural areas where artificial lighting is weak or scarce. In autonomous driving, to ensure safety and situation awareness, robust perception in these conditions is a vital prerequisite.

Transfer learning and domain adaptation approaches aim at narrowing the domain gap between a source domain, where supervised learning from labelled data is possible, to a target domain, where labelled data is either sparse or not available. Such approaches, as demonstrated in [28] or [35], allow to adapt a given segmentation model to a different domain. These approaches, however, do not leverage a complementary modality such as thermal infrared images that can contain more relevant information to solve a given task.

\*These authors contributed equally. All authors are with the University of Freiburg, Germany. Wolfram Burgard is also with the Toyota Research Institute, Los Altos, USA. Corresponding author: vertens@informatik.uni-freiburg.de



**HeatNet: Bridging the Day-Night Domain Gap in Semantic Segmentation with Thermal Images. Johan Vertens, Jannik Zürn and Wolfram Burgard, IROS 2020**



# Motivation



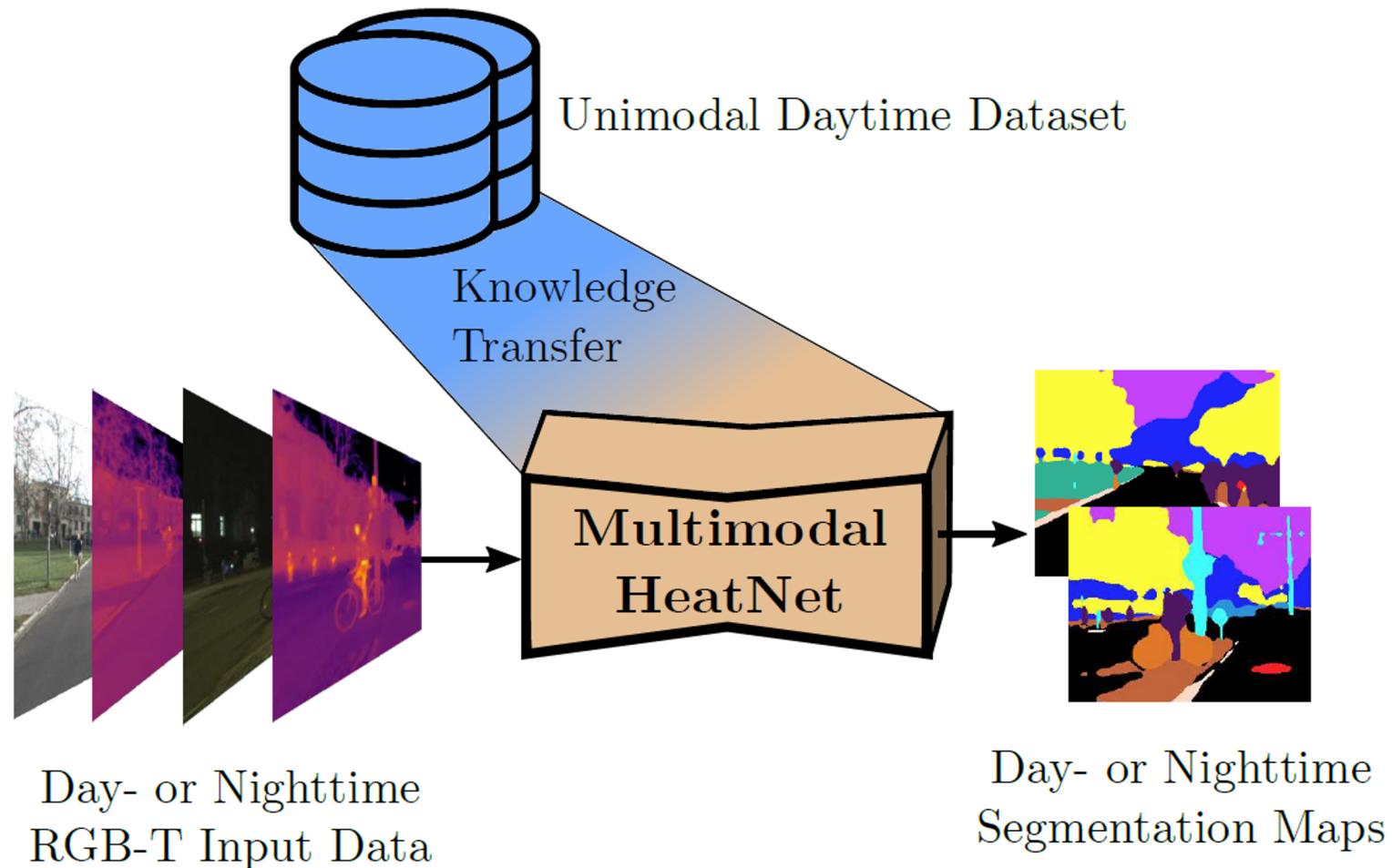
Semantic segmentation prediction during nighttime using conventional CNN trained on publicly available datasets

# Motivation

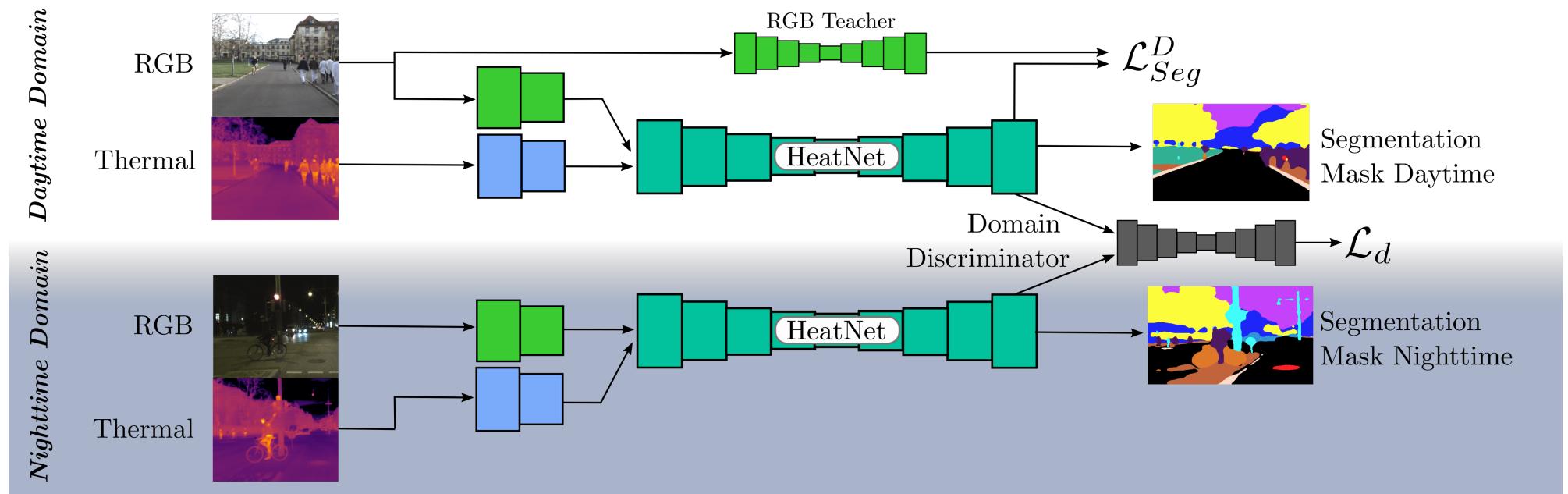


Thermal infrared images exhibit small domain gap between daytime and nighttime

# Approach



# Approach

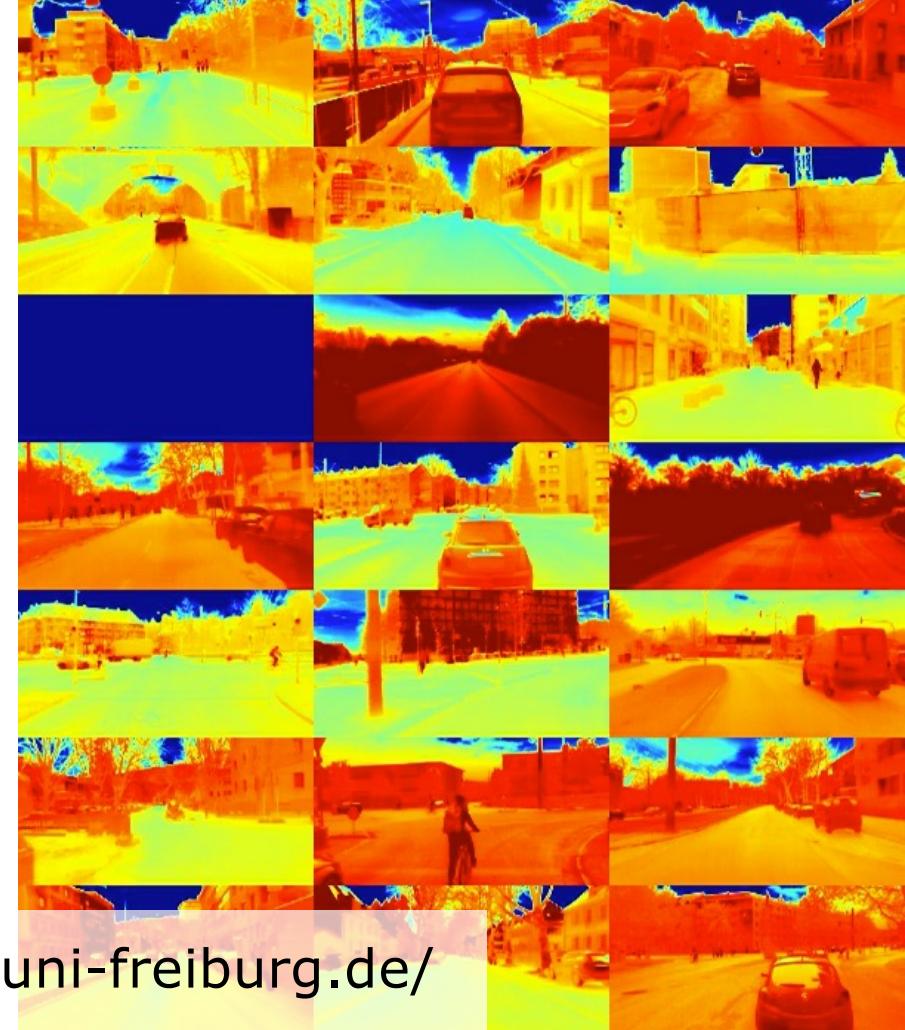


$$\mathcal{L}_{p_1} = \mathcal{L}_s^D + \lambda [0 - C(S_N)]^2, \quad \mathcal{L}_{p_2} = \frac{1}{HW} \sum_{h,w} \begin{cases} [0 - C(S_X)]^2, & \text{if } X = D \\ [1 - C(S_X)]^2, & \text{if } X = N \end{cases}$$

# Dataset

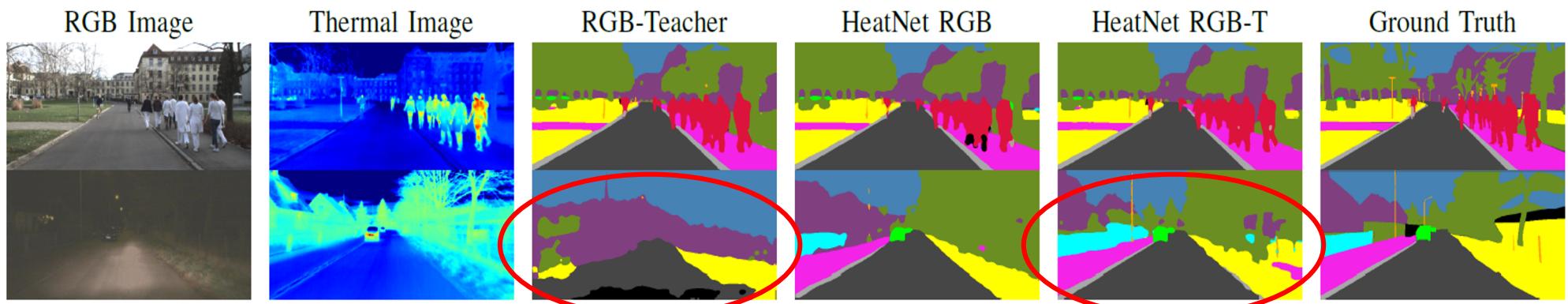


# Dataset



<http://thermal.cs.uni-freiburg.de/>

# Experimental Results

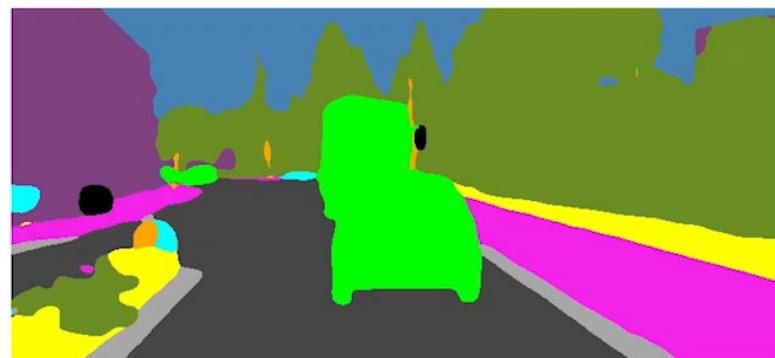
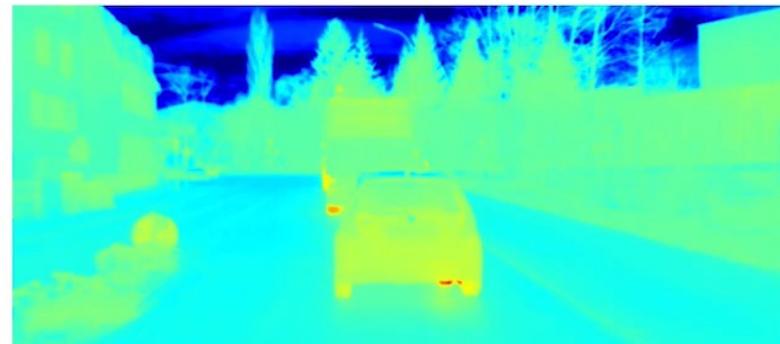


# Experimental Results

RGB



Thermal



HeatNet

# Topology (Lane Graph) Estimation

This work has been submitted to the IEEE for possible publication. Copyright may be transferred without notice, after which this version may no longer be accessible.

## Lane Graph Estimation for Scene Understanding in Urban Driving

Jannik Zürn\*, Johan Vertens\*, and Wolfram Burgard

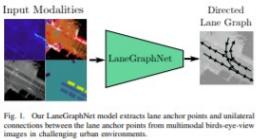


Fig. 1. Our LaneGraphNet model extracts lane anchor points and unidirectional connections between the lane anchor points from multimodal bird's-eye-view images in challenging urban environments.

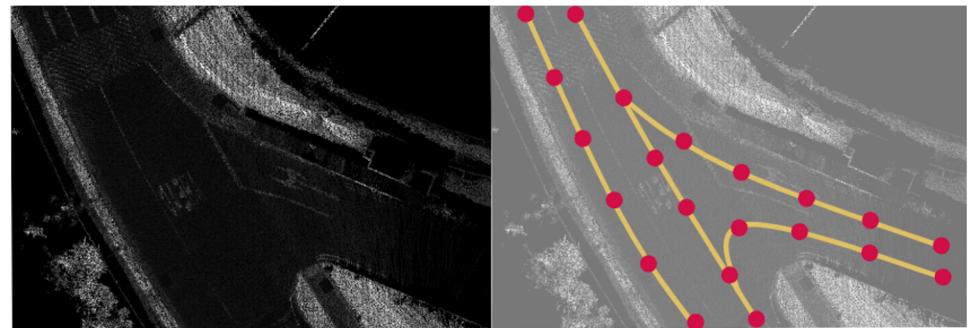
**Abstract**— Lane-level scene annotations provide invaluable data for training planning and complex environments such as urban areas and cities. However, obtaining such data is time-consuming and expensive since lane annotations have to be annotated manually by humans and are as such hard to scale to large areas. In this work, we propose a novel approach for lane geometry estimation from bird's-eye-view images. We formulate the problem of lane shape and lane connectivity estimation as a graph estimation problem where lane anchor points are nodes and lane segments are graph edges. We train a graph estimation model on multimodal bird's-eye-view data processed from the popular NuScenes dataset and its mapping expansion pack. We furthermore estimate the direction of the lane segments and lane boundaries from the graph model which results in a directed lane graph. We illustrate the performance of our LaneGraphNet model on the challenging NuScenes dataset and evaluate it with a quantitative evaluation. Our model shows promising performance for most evaluated urban scenes and can serve as a step towards automated generation of HD lane annotations for autonomous driving.

## INTRODUCTION

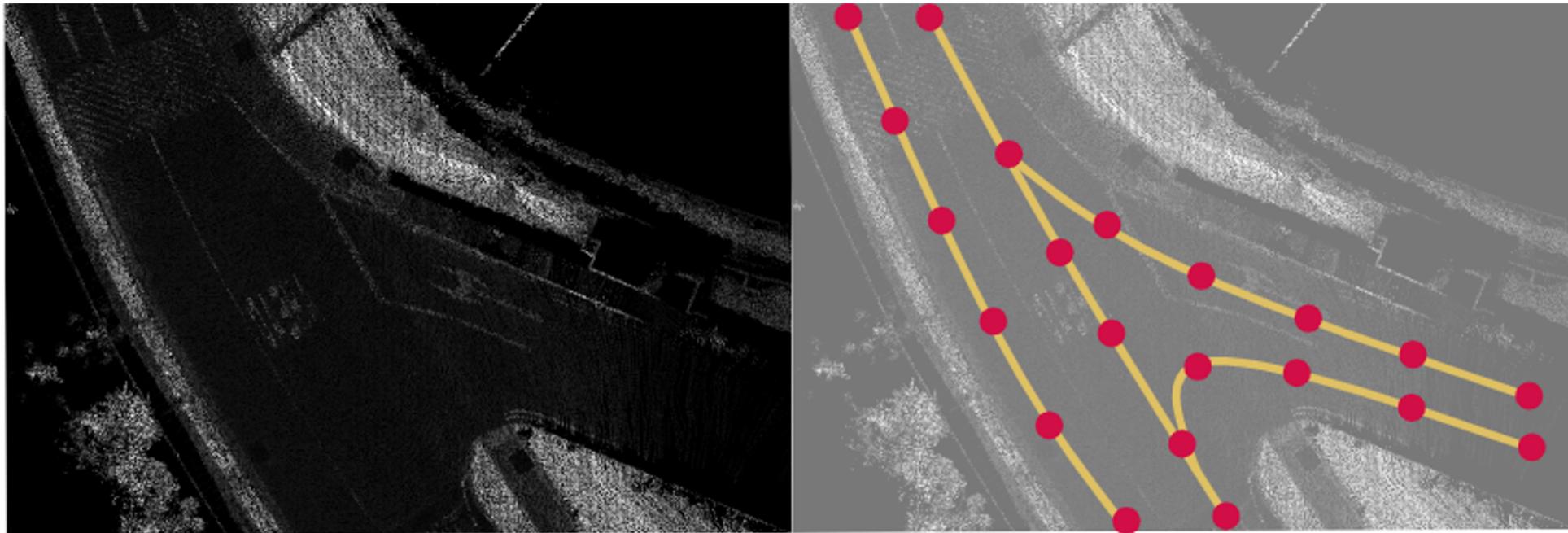
HD maps play an important role for autonomous navigation [18], [10], in particular in semi-structured and dynamic environments such as urban areas. In this work, we are interested in learning to infer geometric annotations for lane centerlines which are integral to HD maps and which allow autonomous vehicles to determine the position and orientation to access the exact location of lane anchor points and lane directions enables autonomous vehicles to determine legal trajectories from the current position to goal points, both in map-based autonomous driving and in mapless autonomous driving. Despite their utility for autonomous driving, lane-level annotations are expensive and time-consuming to obtain and are often only available to be obtained due to construction works or other changes to the road and its surroundings [18]. Due to the time-consuming manual annotation process, obtaining large-scale lane-level annotations remains a challenge to this day and poses a bottleneck for widespread deployment of fully autonomous driving. Therefore, automatic inference of lane centerline annotations from vehicle sensor readings is an important step towards achieving autonomous driving at scale.

\*These authors contributed equally. All authors are with the University of Freiburg, Germany. Wolfram Burgard is also with the Toyota Research Institute. De Anan Cai is the corresponding author. [zuerw@informatik.uni-freiburg.de](mailto:zuerw@informatik.uni-freiburg.de)

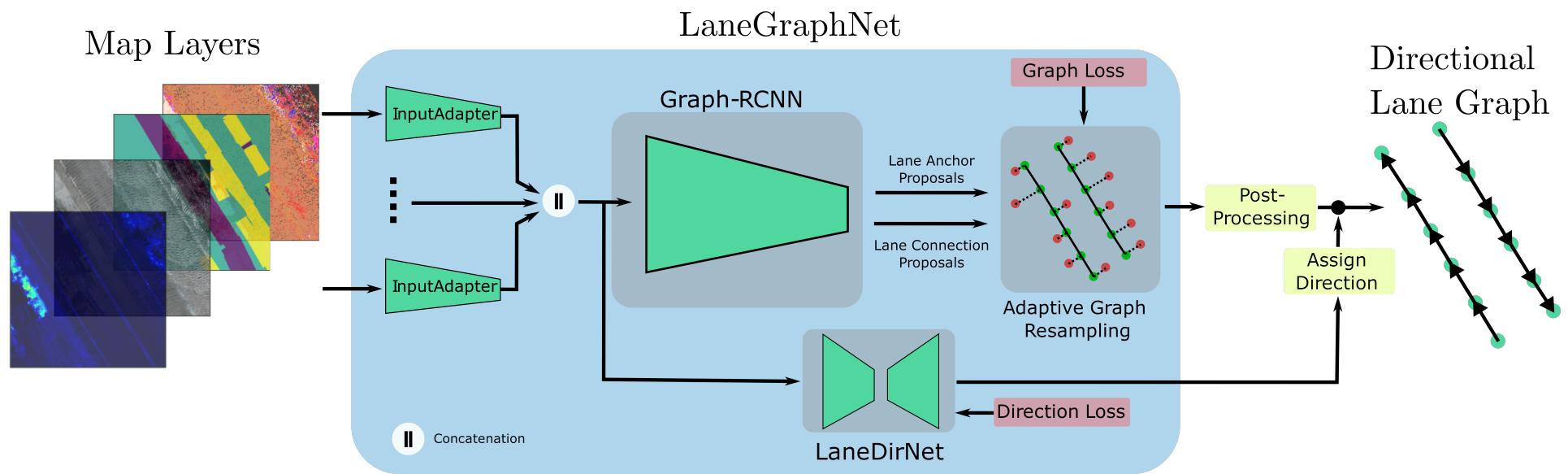
## Lane Graph Estimation for Scene Understanding in Urban Driving. Jannik Zürn, Johan Vertens and Wolfram Burgard, Submitted to RA-L



# Can we Learn HD-maps from BEV Sensor Data?

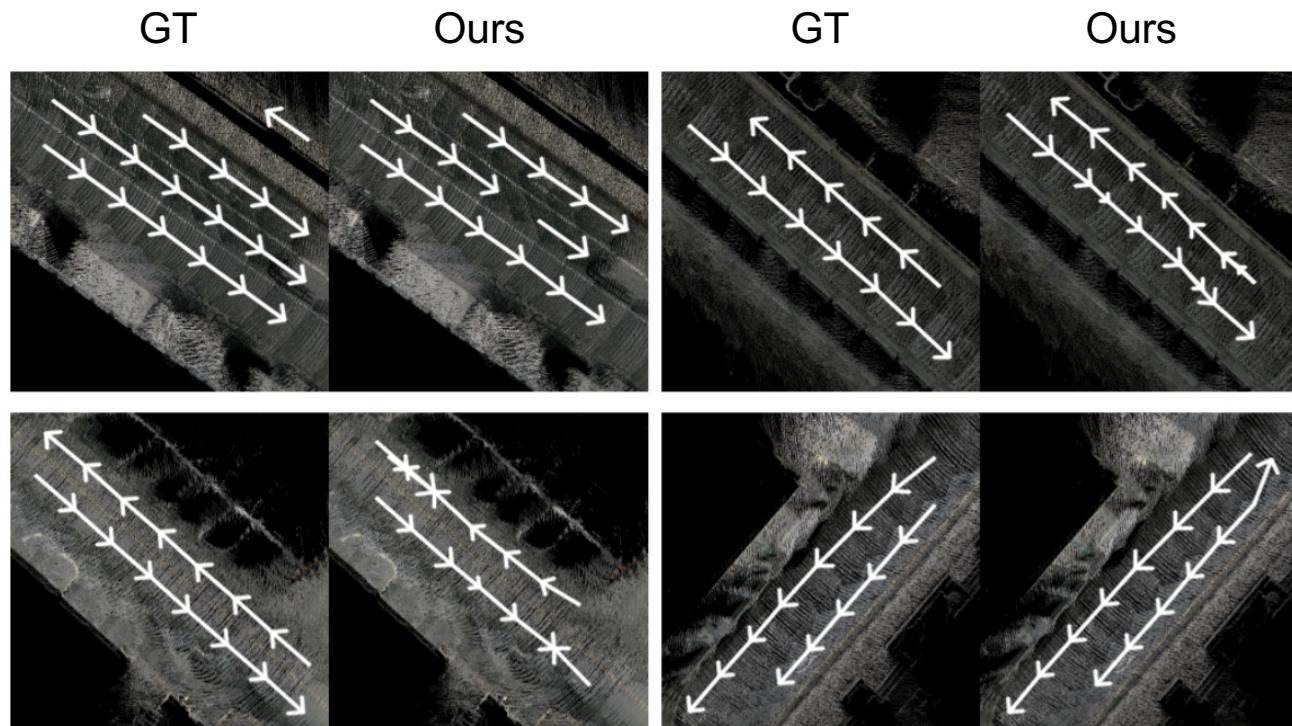
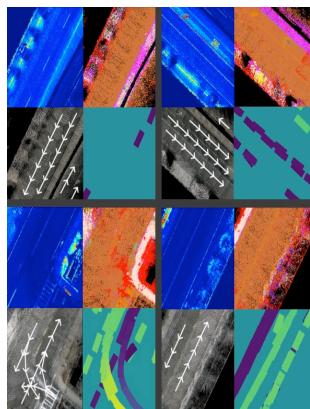


# Approach



# Approach, Qualitative Results

BEV-Modalities:  
Lidar, RGB, Semantics,  
Observed Car Poses



# Summary

- HD maps are useful but limiting
- If we commit to HD maps we need to have proper approaches for change detection
- If we want to get rid of HD maps we need substantially extend the perception capabilities
- The key technology for this will be **Machine Learning**