

# FINAL PRESENTATION

# Instance Completion and Motion Estimation with Deep Shape Priors for Autonomous Driving

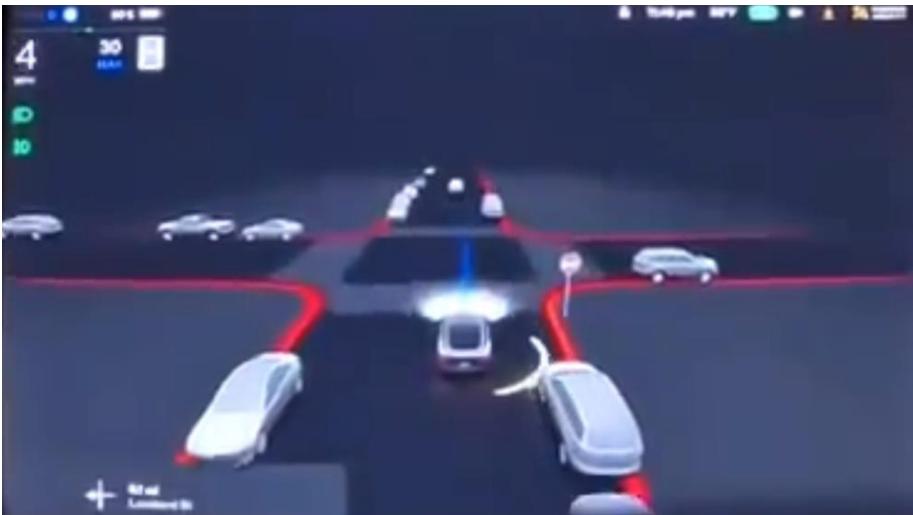
## Team members

Panyawat (Ohm) Rattana  
Shashank (Sai) Dammalapati

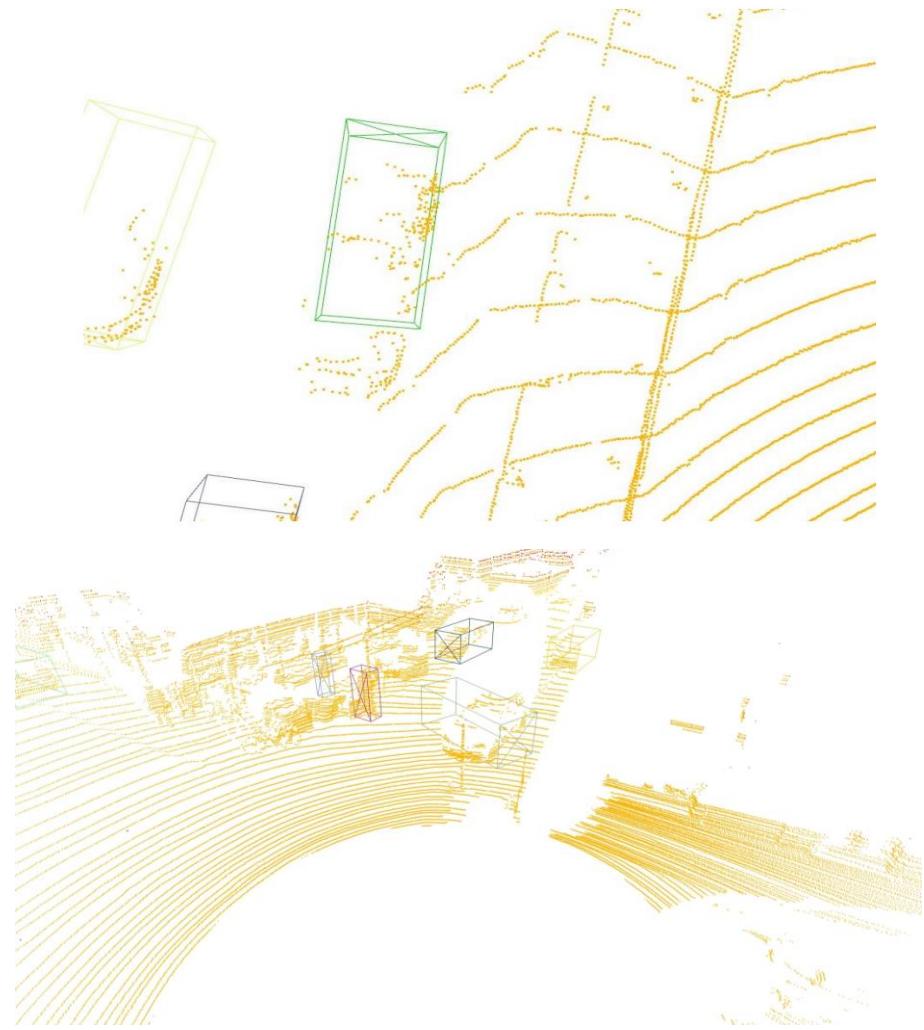
## Supervisors

Xingguang (Starry) Zhong  
Yue Pan

# Motivation



Driving Scene

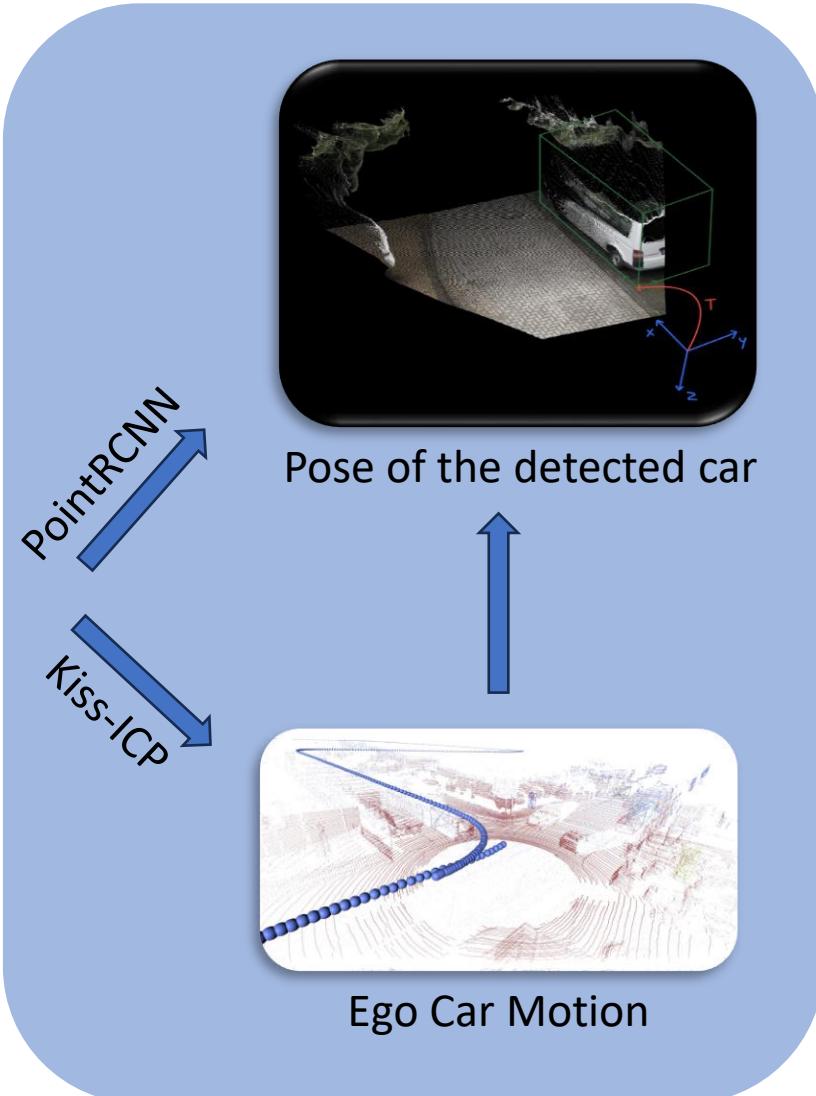


Incorrect Bounding Box Predictions

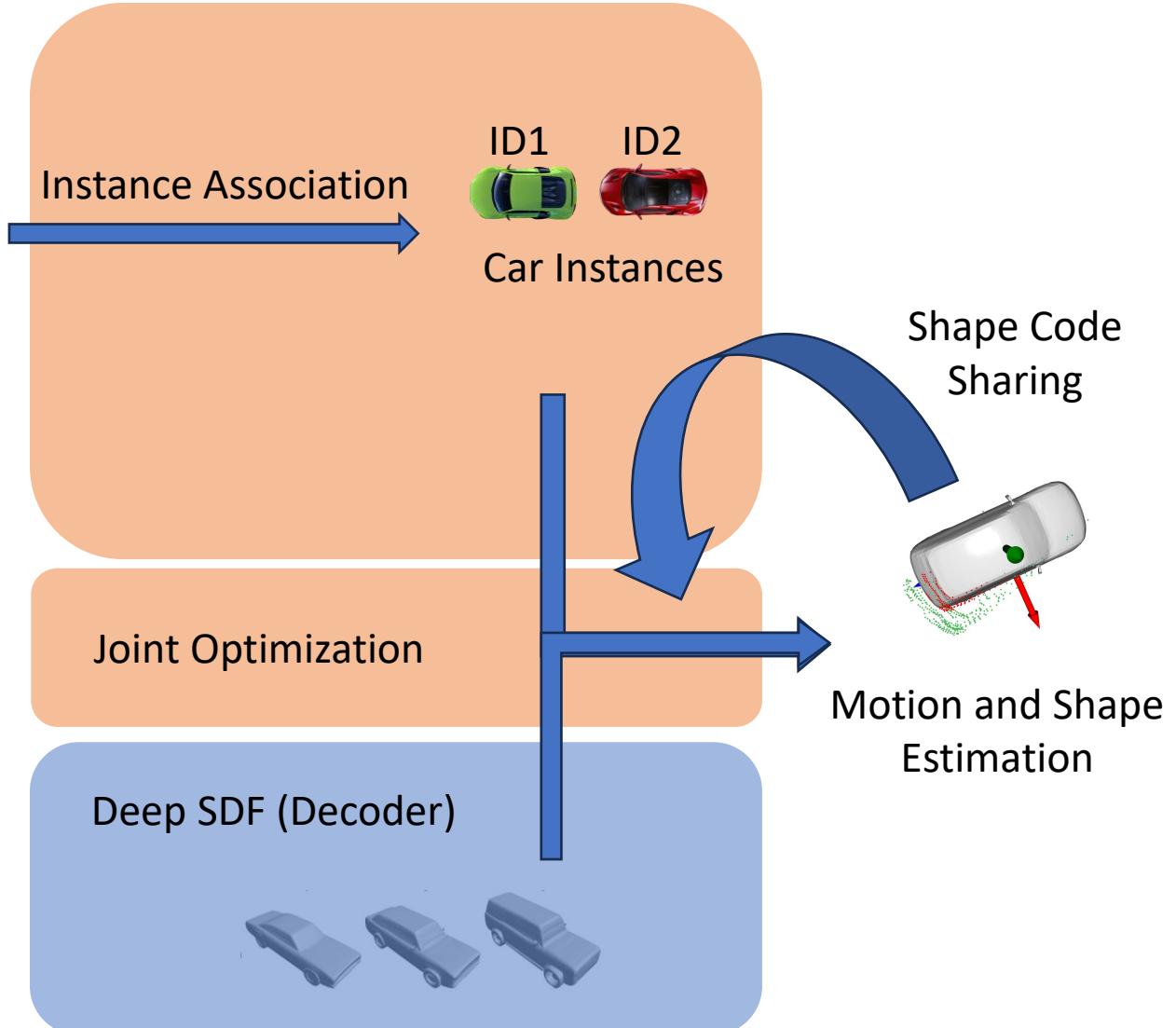
# Pipeline



Raw Lidar Sensor  
Data

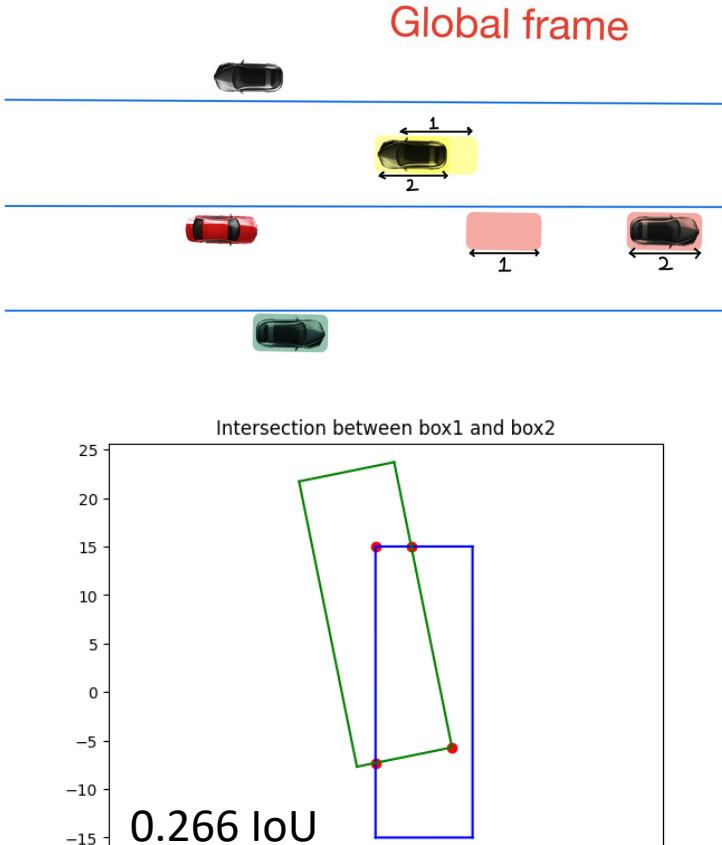


Existing Methods

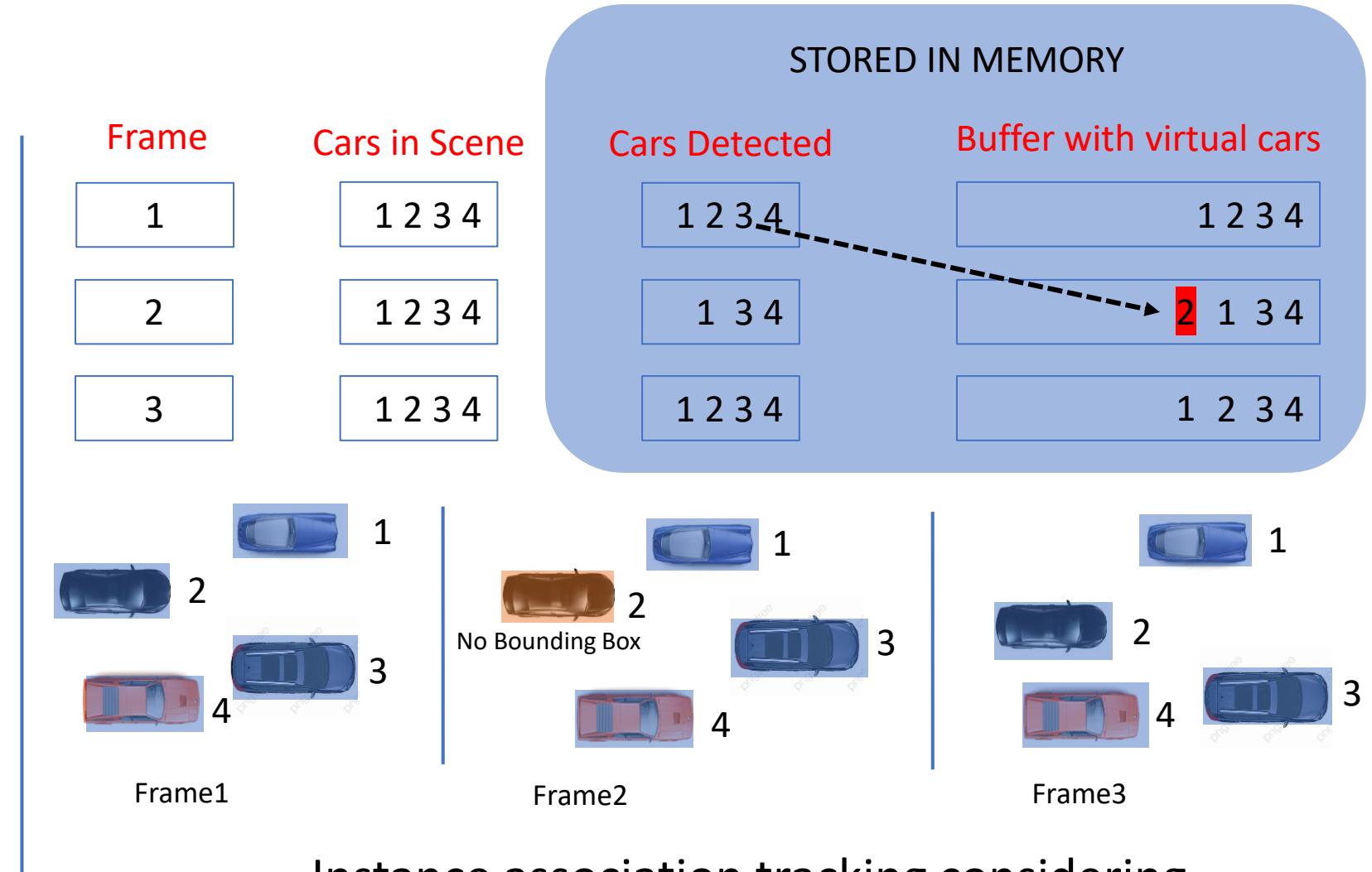


Our Contribution

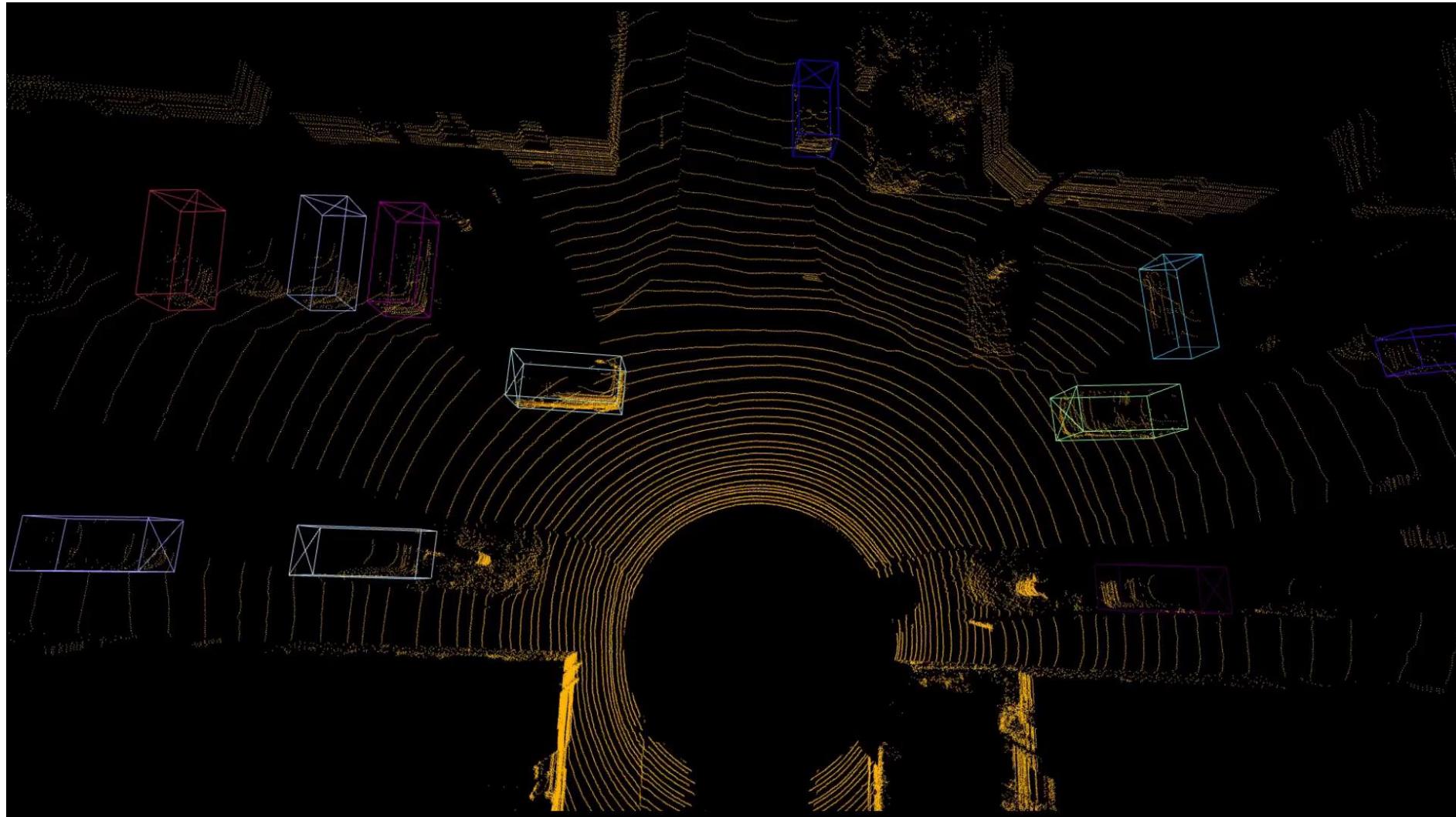
# Methodology (Instance Association)



IoU calculated in Global frame

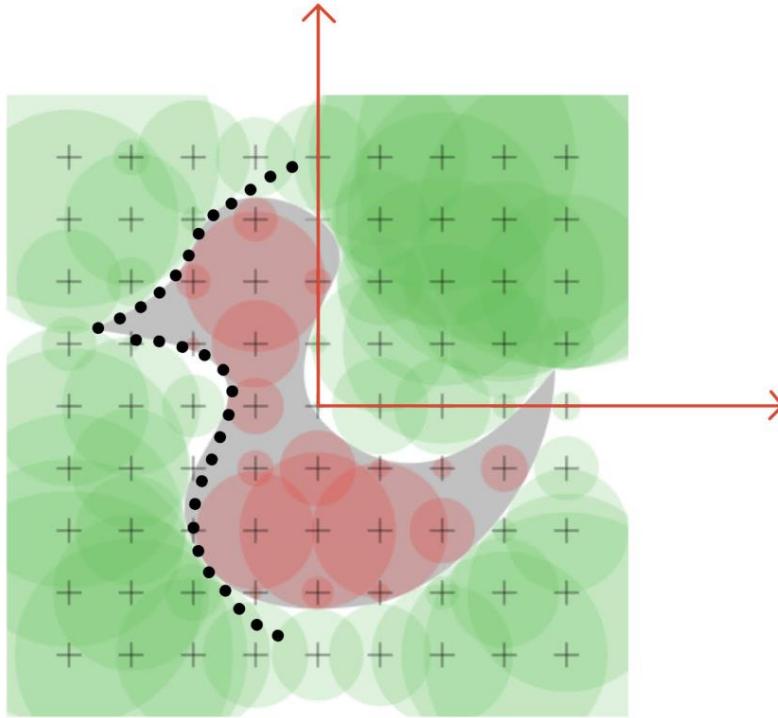


# Methodology (Instance Association)

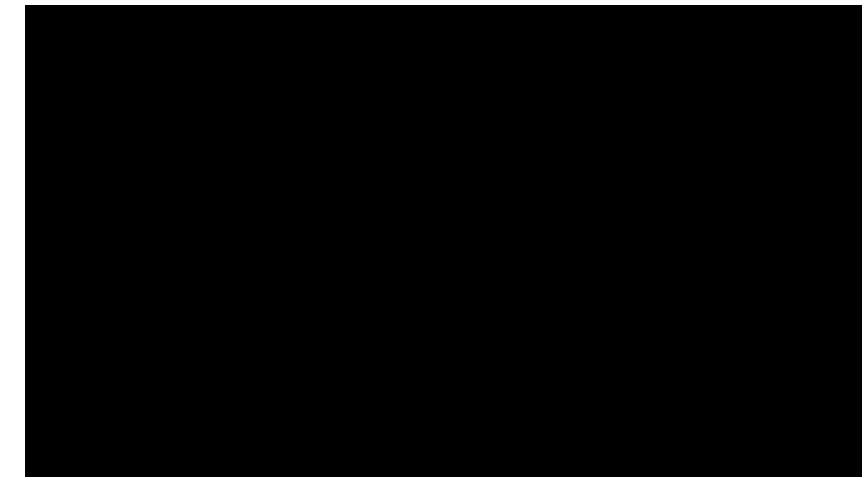
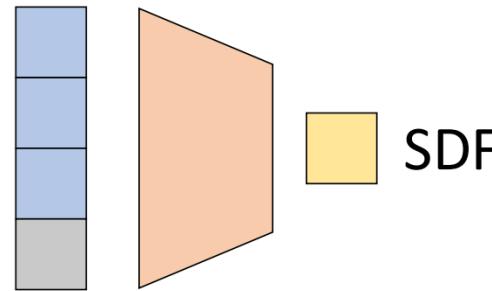


Kitti Sequence 07

# Methodology: 0 SDF



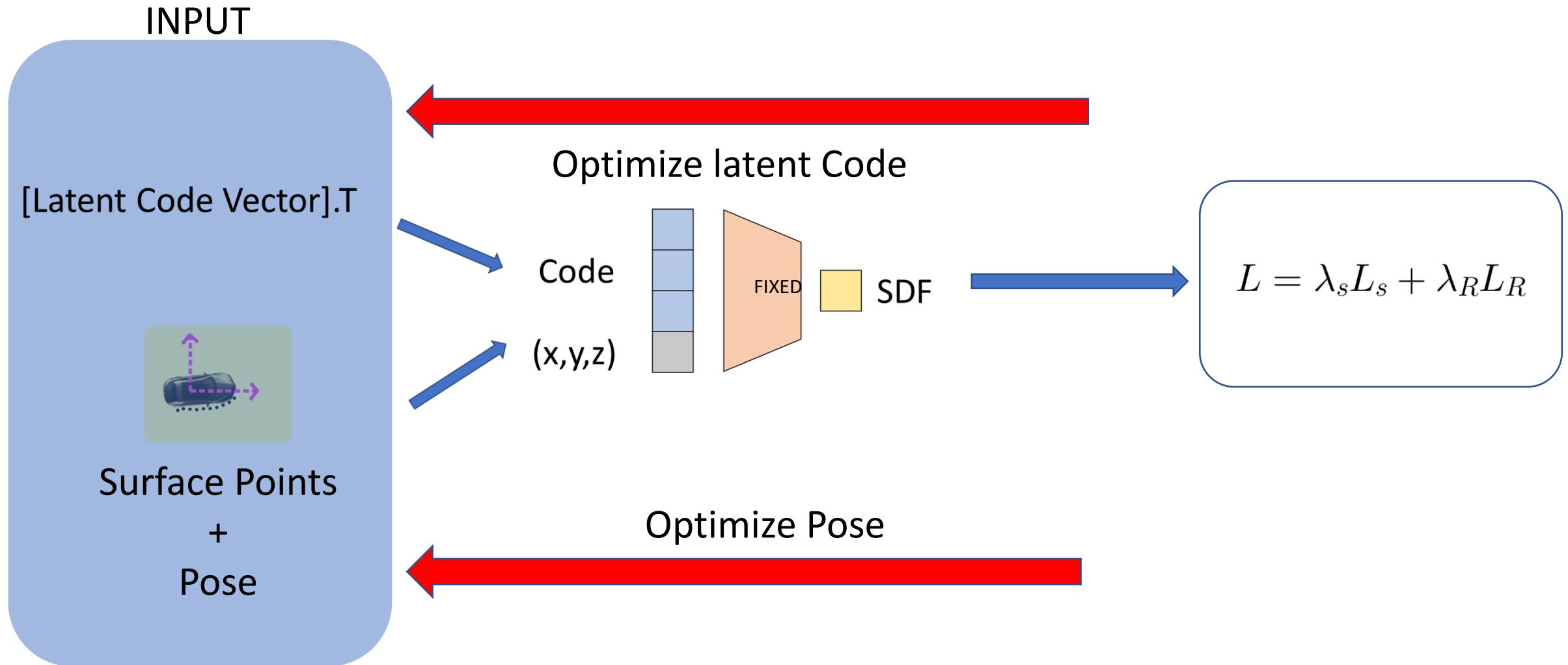
Code  
 $(x, y, z)$



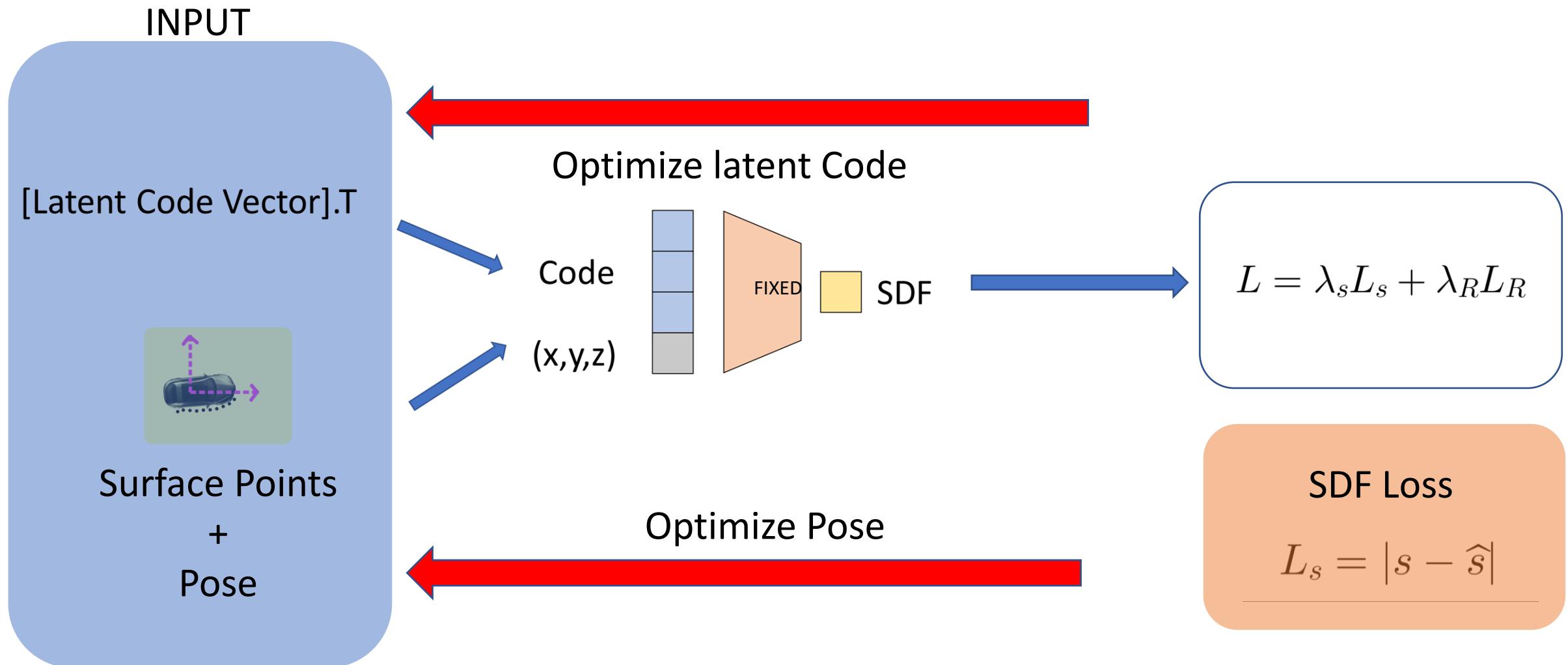
Deep SDF trained on ShapeNet

SDF value for a point on surface modelled by SDF is ZERO

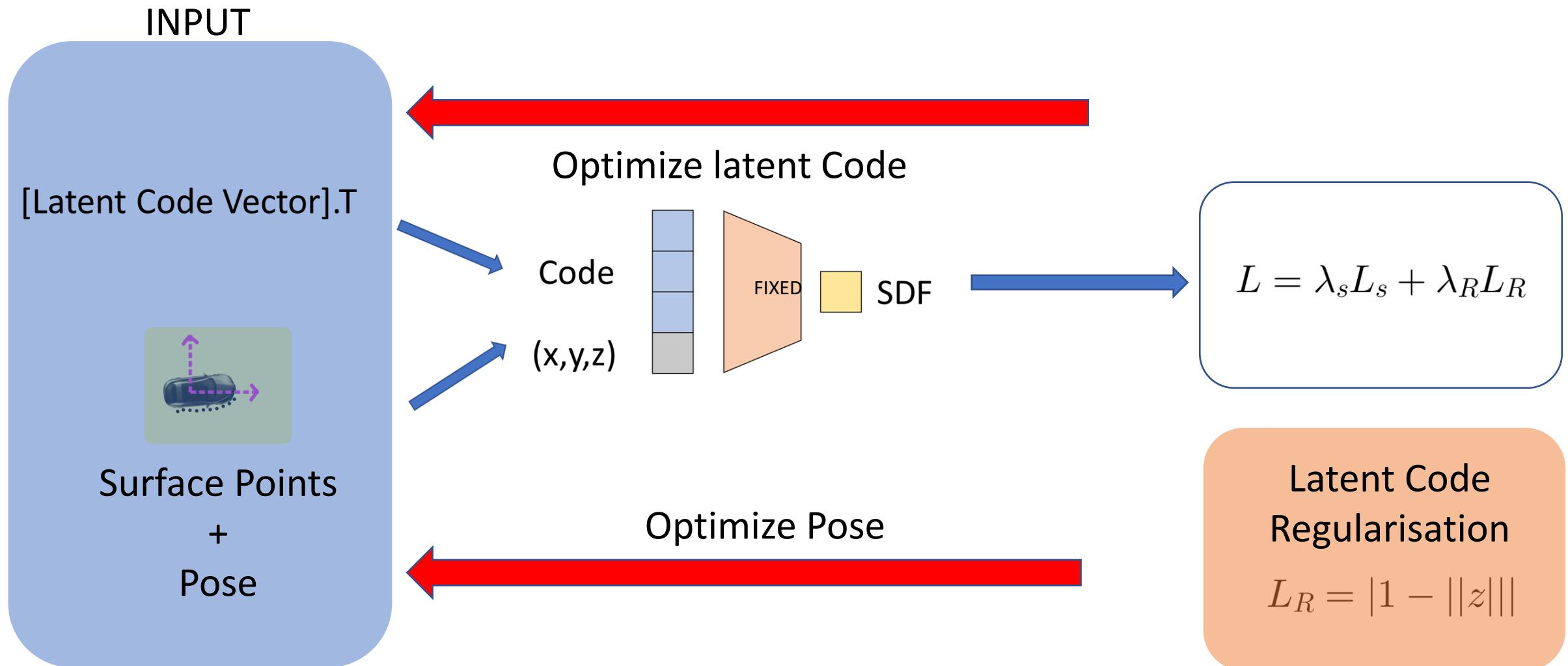
# Optimization Process



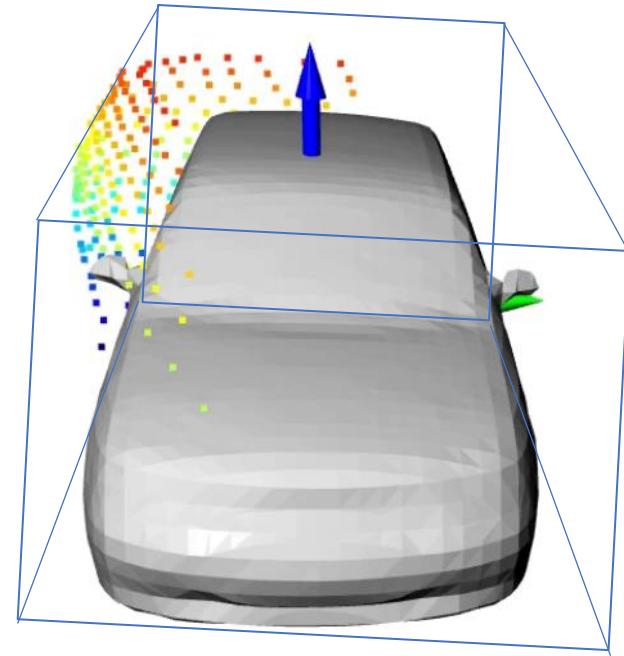
# Optimization Process



# Optimization Process



# Methodology (Joint Optimization)



**Input:** 1 frame of 1 car

**Output:** Optimized bounding box pose + Mesh

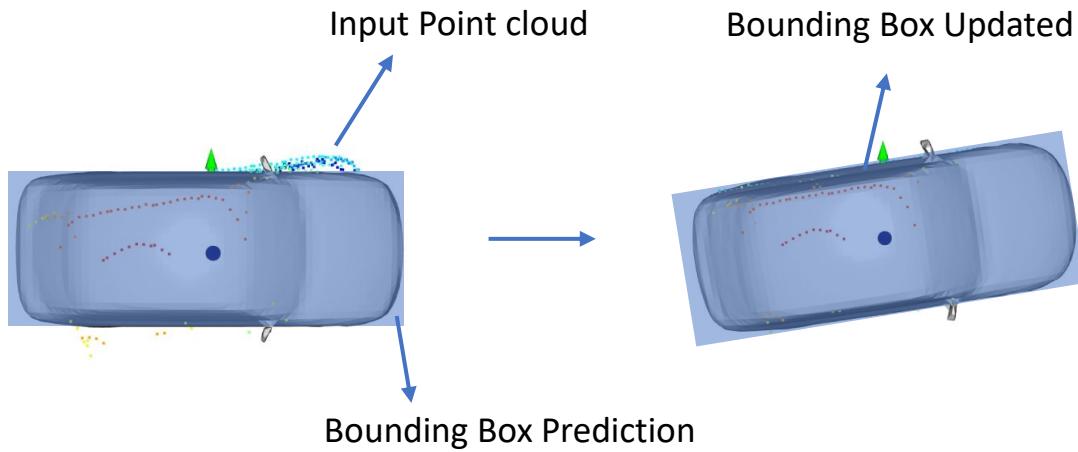
# Result

(Yaw) Pose Correction

Shape Code Estimation

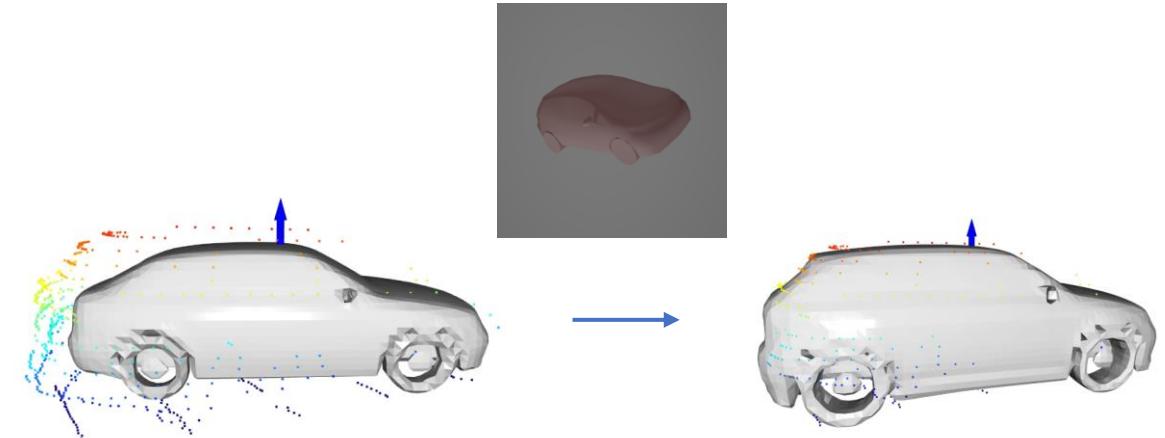
# Result

## Optimized Pose Estimation



## (Yaw) Pose Correction

## Optimized latent Code



## Shape Code Estimation

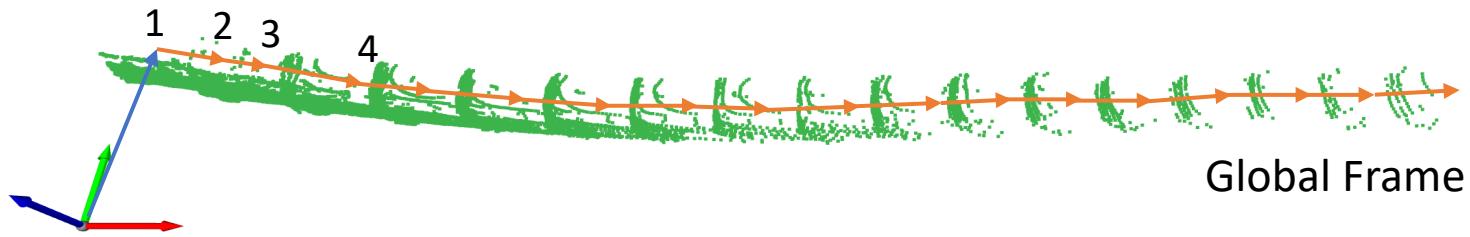
# Tracking Motion



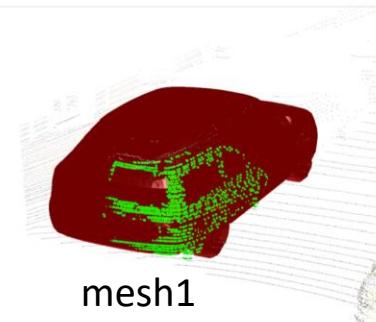
We optimize  $SE(3)$  transformation representing the movement from one timestamp to other

Optimization of 1 car instance over multiple frames

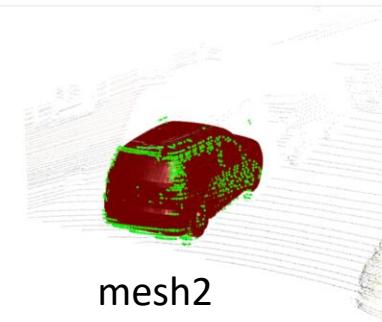
# Scale Ambiguity



Global Frame



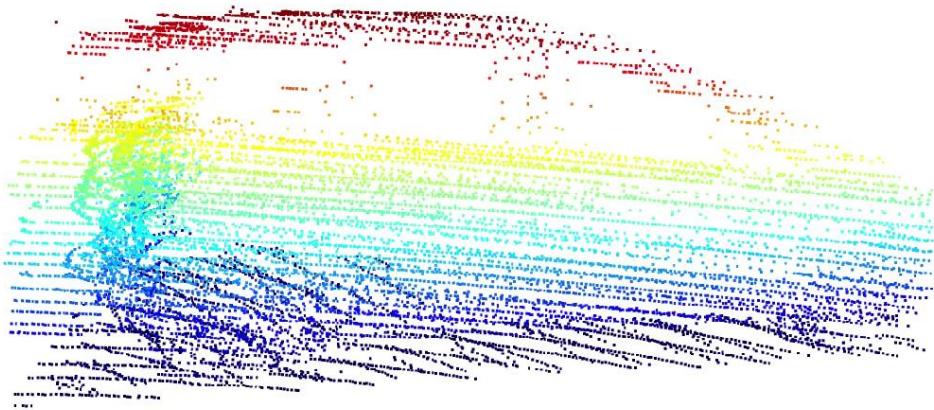
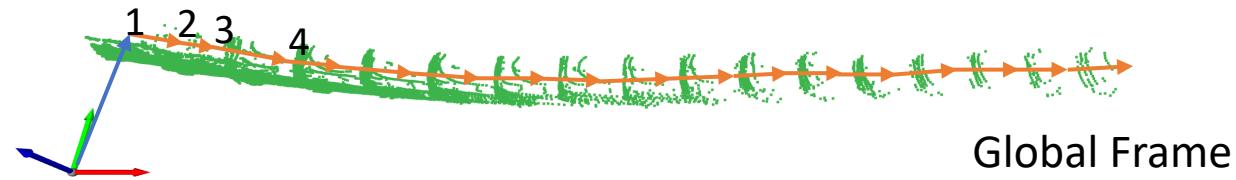
mesh1



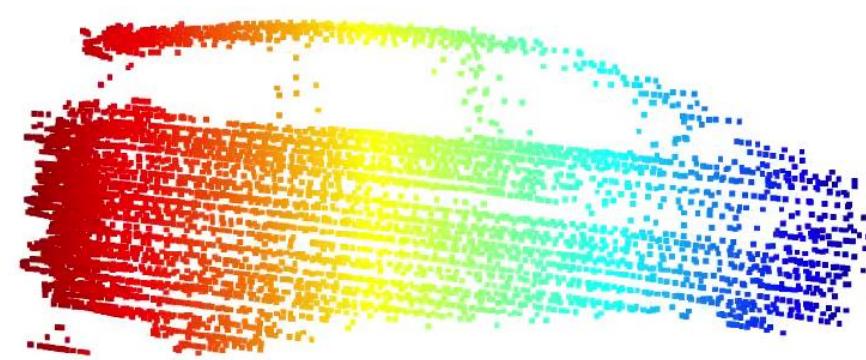
mesh2

We accumulate the points using optimized bounding box poses from one frame to the other

# Point Cloud Accumulation



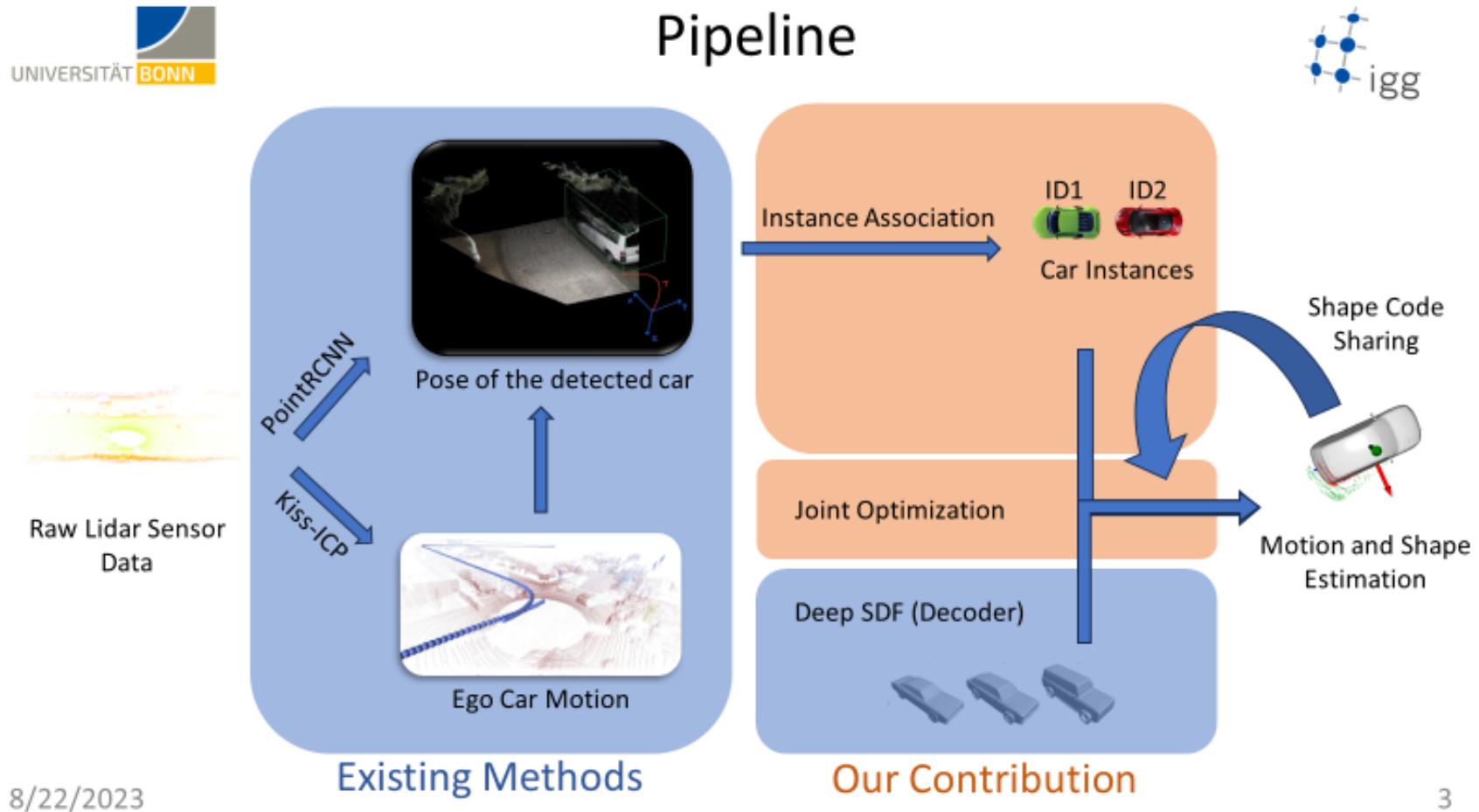
Point Cloud Accumulation based on  
bounding box predictions



Point Cloud Accumulation based on  
optimized poses for bounding boxes

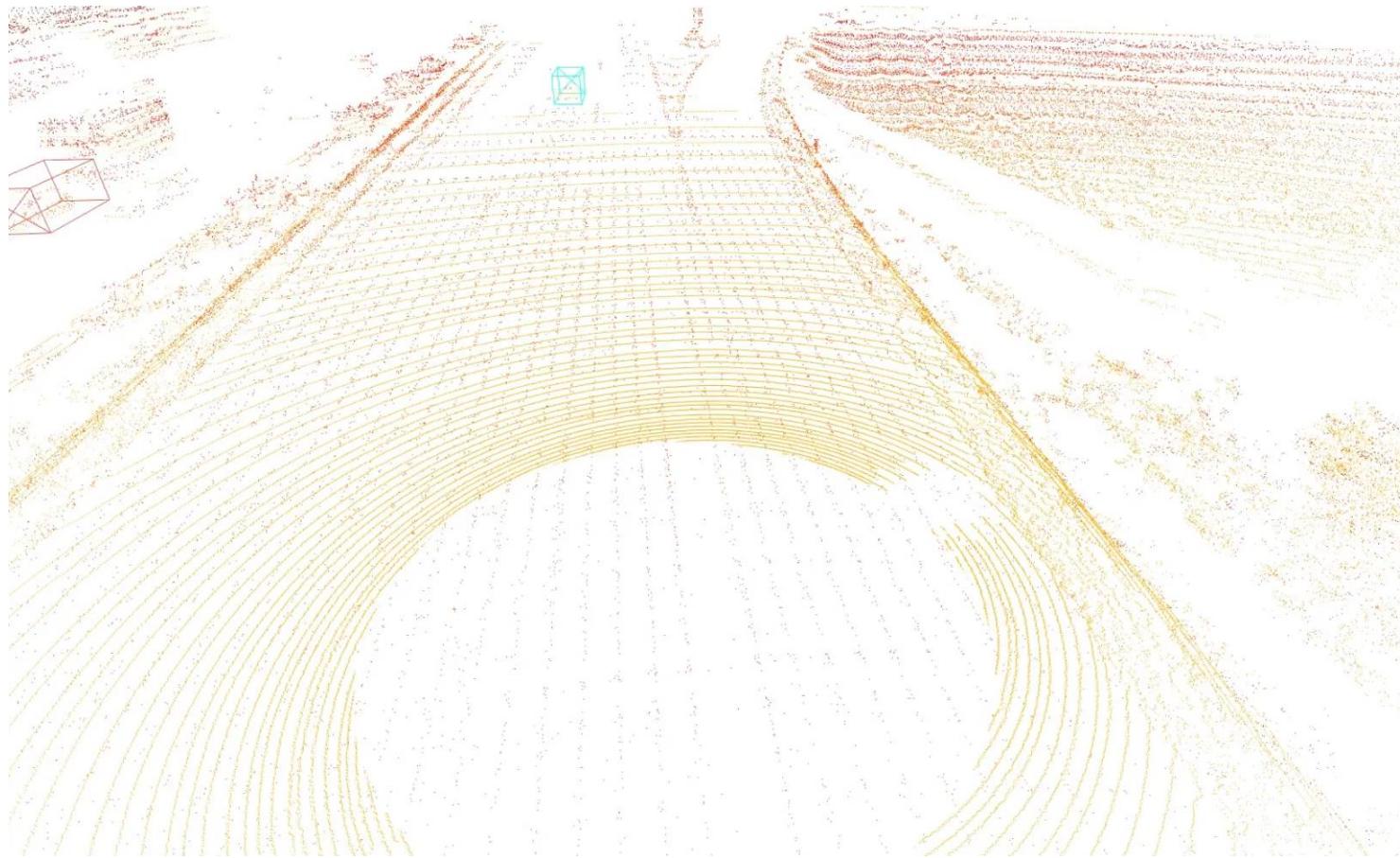
Optimization of 1 car instance over multiple frames

# Accumulation + Shape Code Sharing



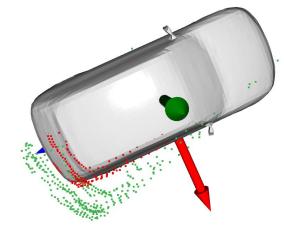
Optimization of 1 car instance over multiple frames

# Optimizing for Multiple Cars

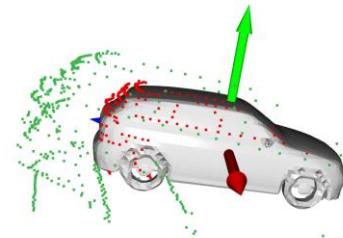


Instance Completion and Motion Estimation with Instance Association

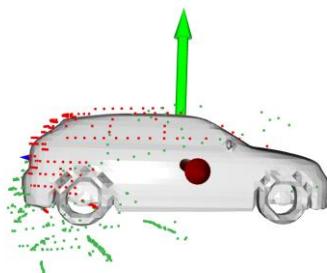
# Qualitative Study



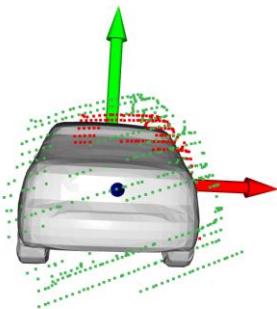
9° Yaw error correction



1.5 scale factor correction



9° Pitch error correction



9° Roll error correction

Green: Points in canonical frame before pose optimization correction

Red: Points in canonical frame after pose optimization

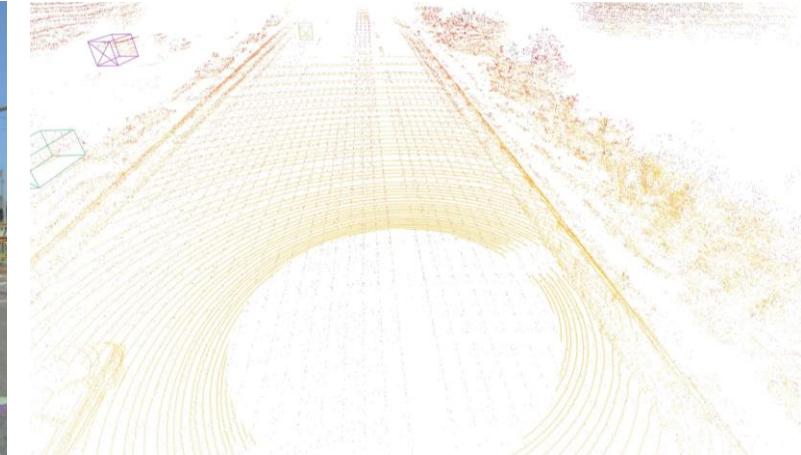
Green (Points + Pose)



Optimization

Red (Points + Pose)

# Future Work



Improve Instance Association

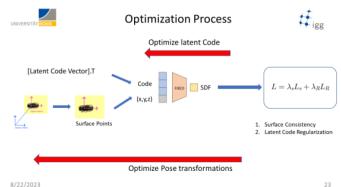
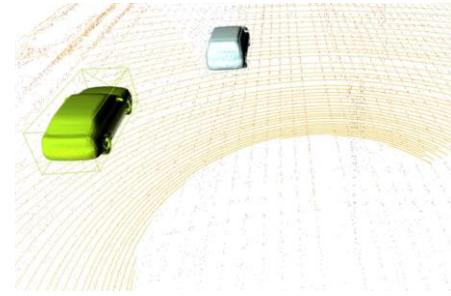
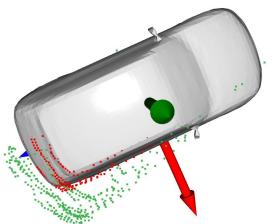
1. Use pose posterior
2. Motion Model

Evaluation using ground truth

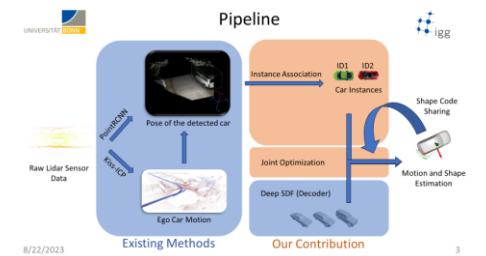
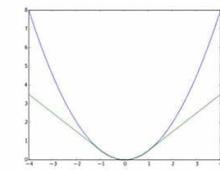
1. Argoverse
2. Simulation

Improve program  
and make it online

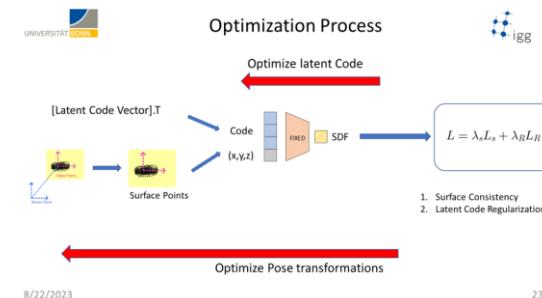
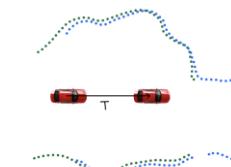
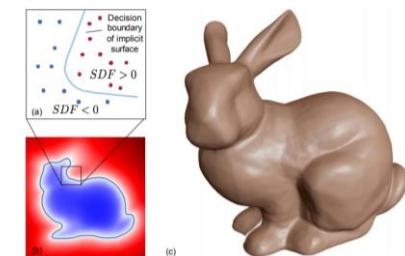
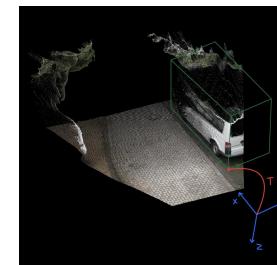
# Any Questions?



fps



# Summary of the project

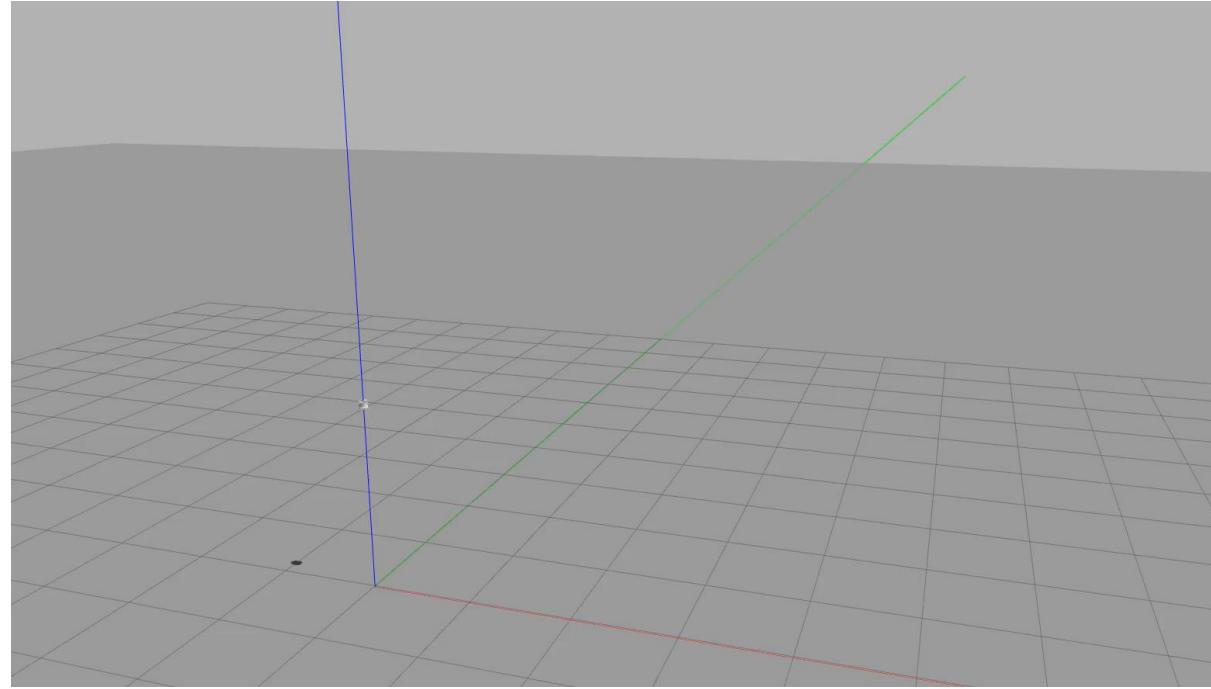


## Challenges

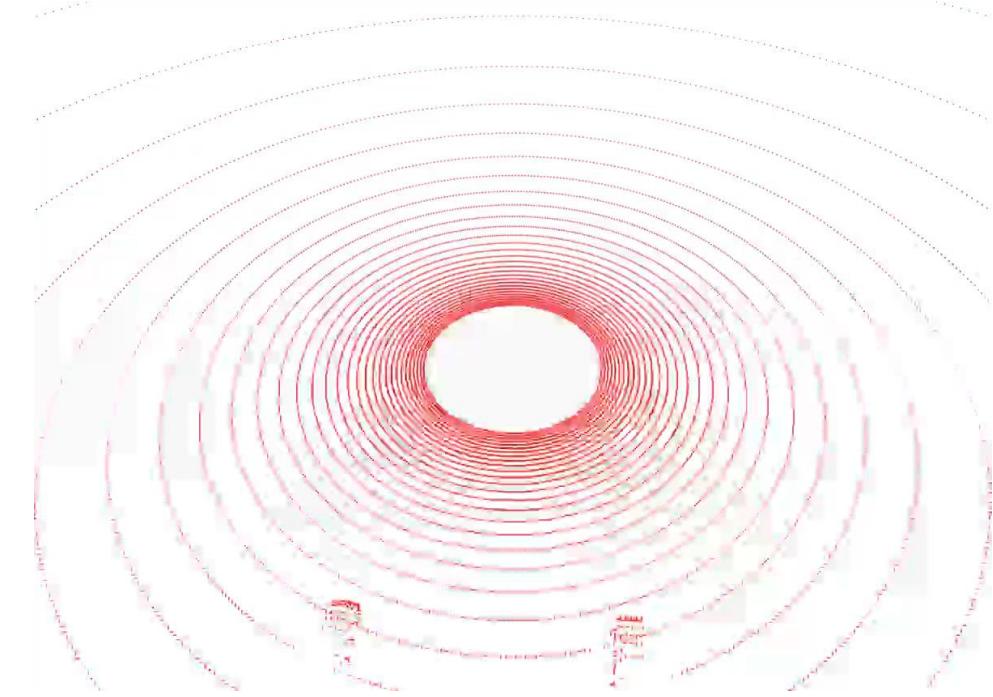
## Method

## Final Output

# Ground Truth Evaluation



Gazebo simulation provides shape and ground truth poses



Recorded 64 channel lidar data from simulation

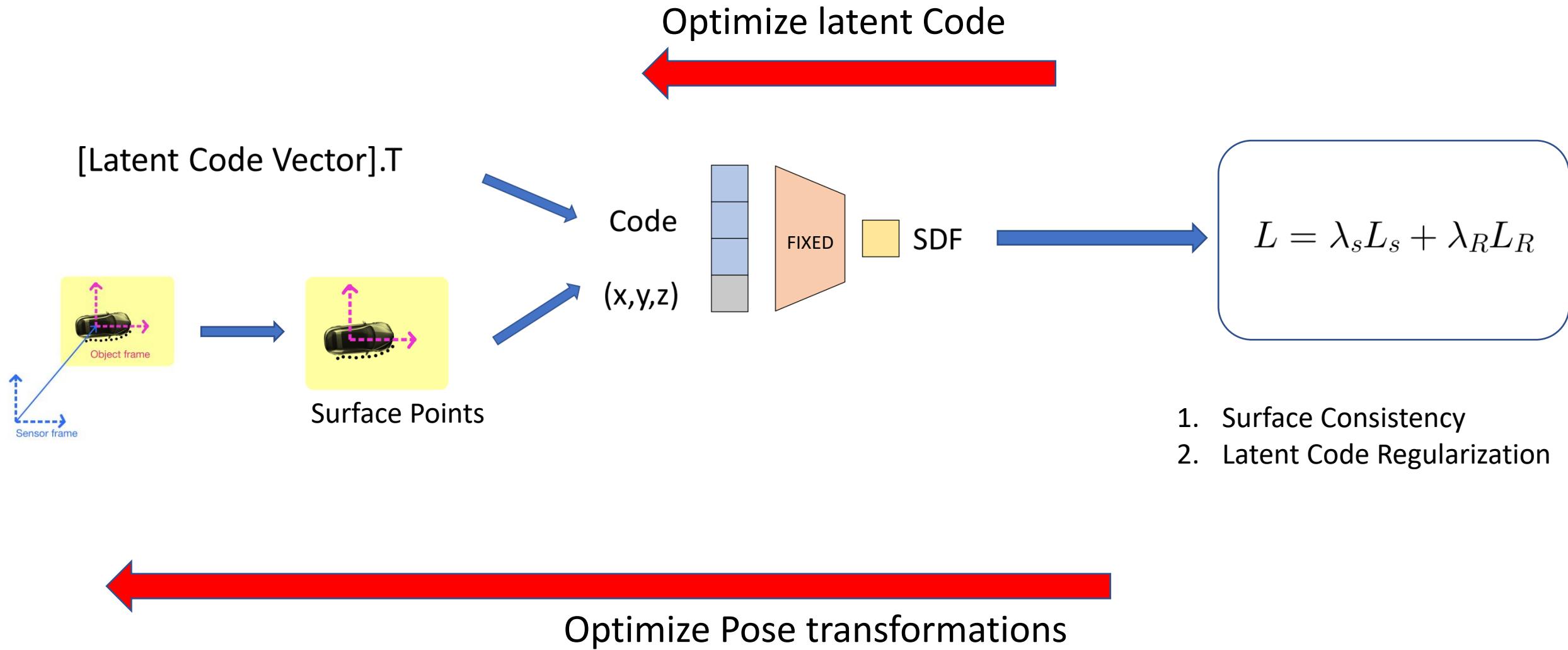
# Time Metrics

## TASK SPECIFIC FRAME RATE

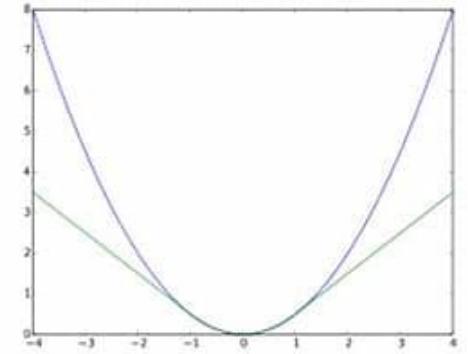
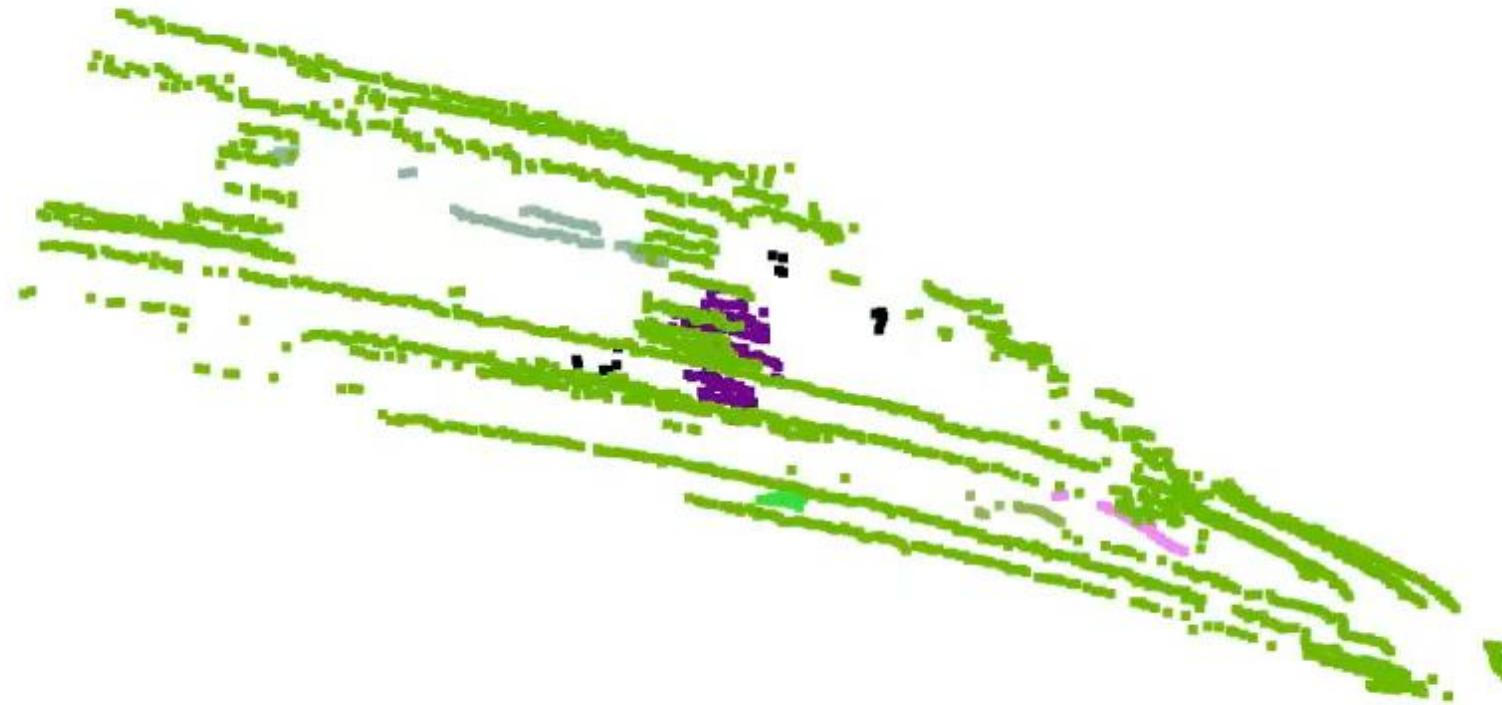
Device: Nvidia RTX 3060

NO	Task	FPS   ms
1	KissICP	60 (Tested)
2	PointRCNN	8 (Tested)
3	Code & Pose Optimization	10 (Tested)
4	Marching Cubes	125ms per mesh (Tested)
5	DSP-SLAM	10 (Paper)
6	DSP-SLAM(including mesh extraction)	5 (Tested)

# Optimization Process



# Robust Kernels



Huber loss

Transparent glass results in driver surface points

# Other Works using Shape Completion

## Panoptic Mapping with Fruit Completion and Pose Estimation for Horticultural Robots

Yue Pan  
Claus Smitt

Federico Magistri  
Chris McCool

Thomas Läbe  
Jens Behley

Elias Marks  
Cyrill Stachniss

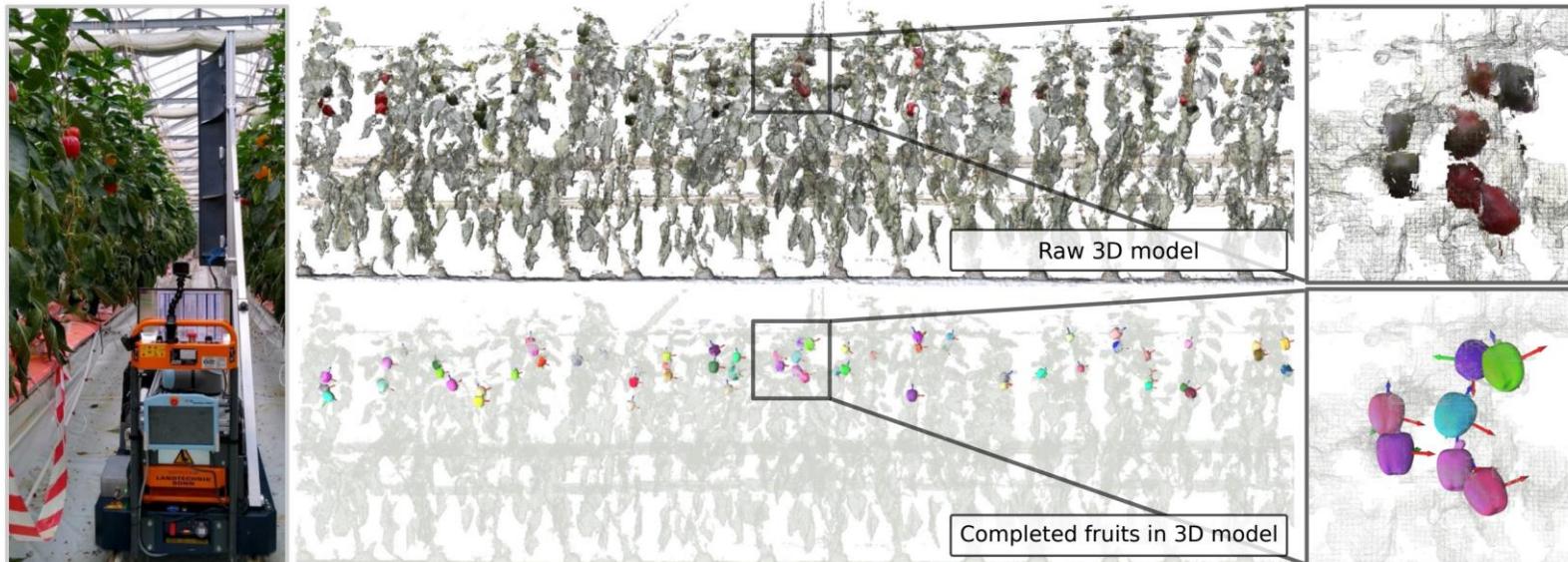
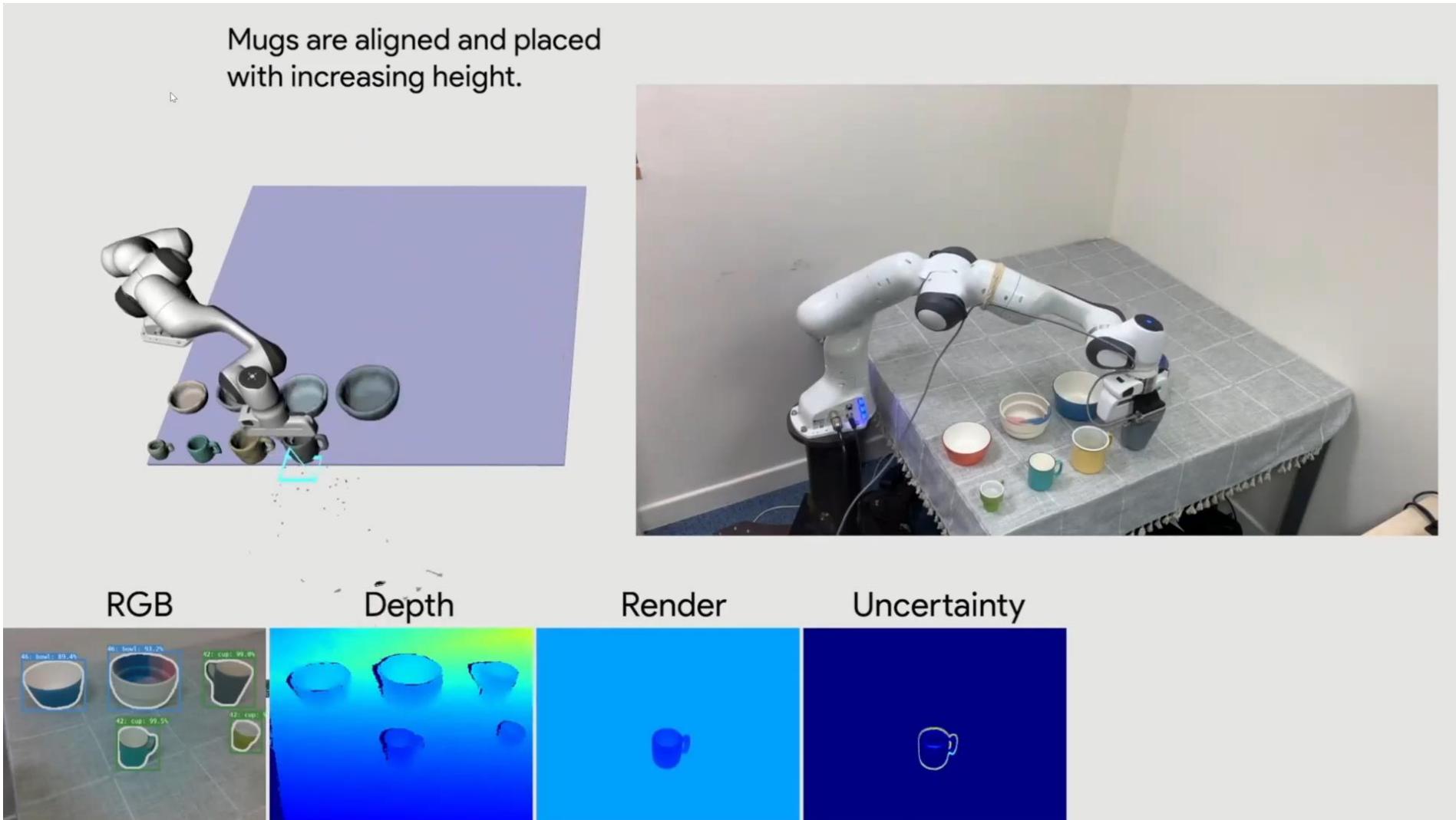


Fig. 1: Our method is able to build a multi-resolution panoptic map (top) of a challenging commercial glasshouse environment online using a mobile horticultural robot equipped with RGB-D cameras (left). Furthermore, our method manages to jointly estimate the complete shape and pose of each fruit in the map (bottom).

# Node SLAM

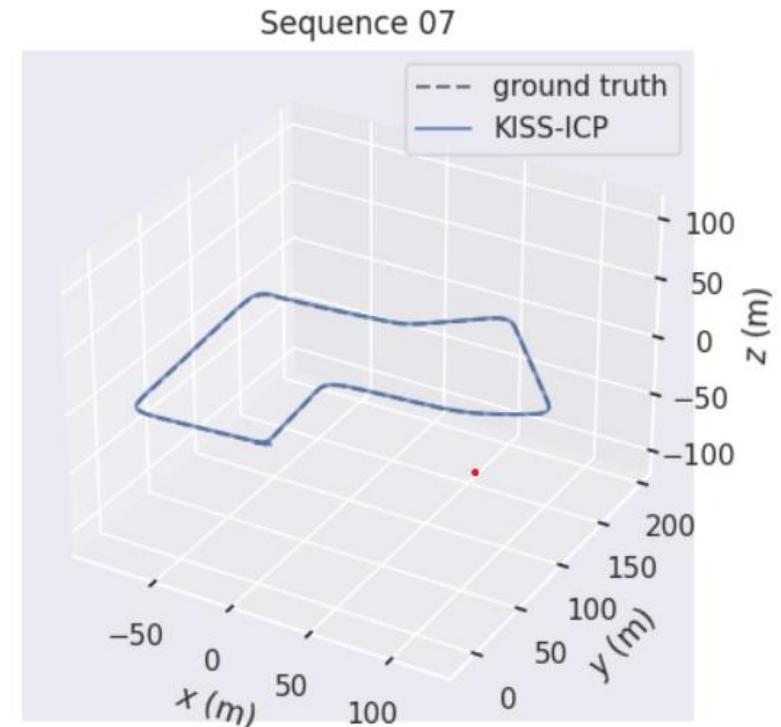


# KissICP evaluation

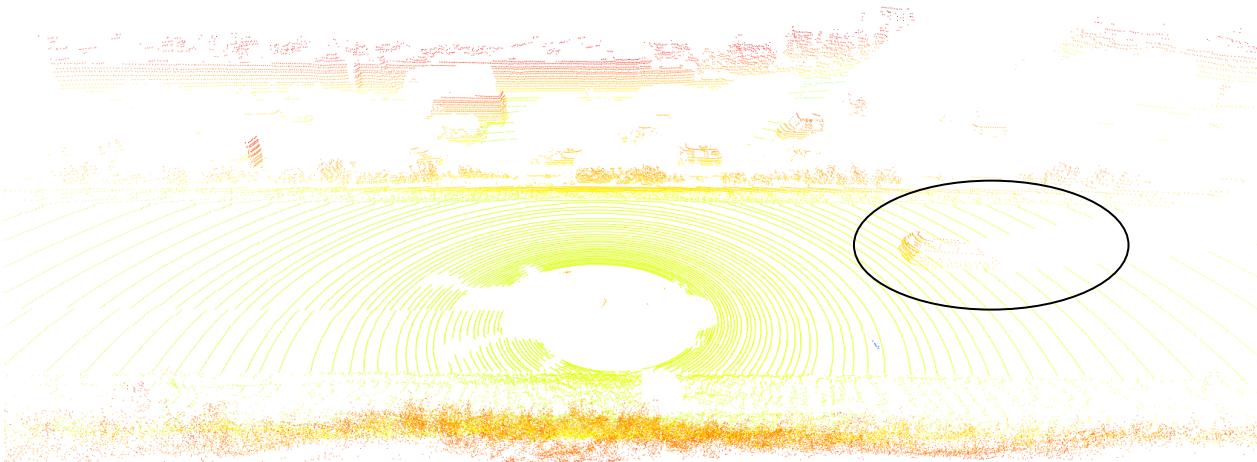
Now evaluating sequence 07

0%

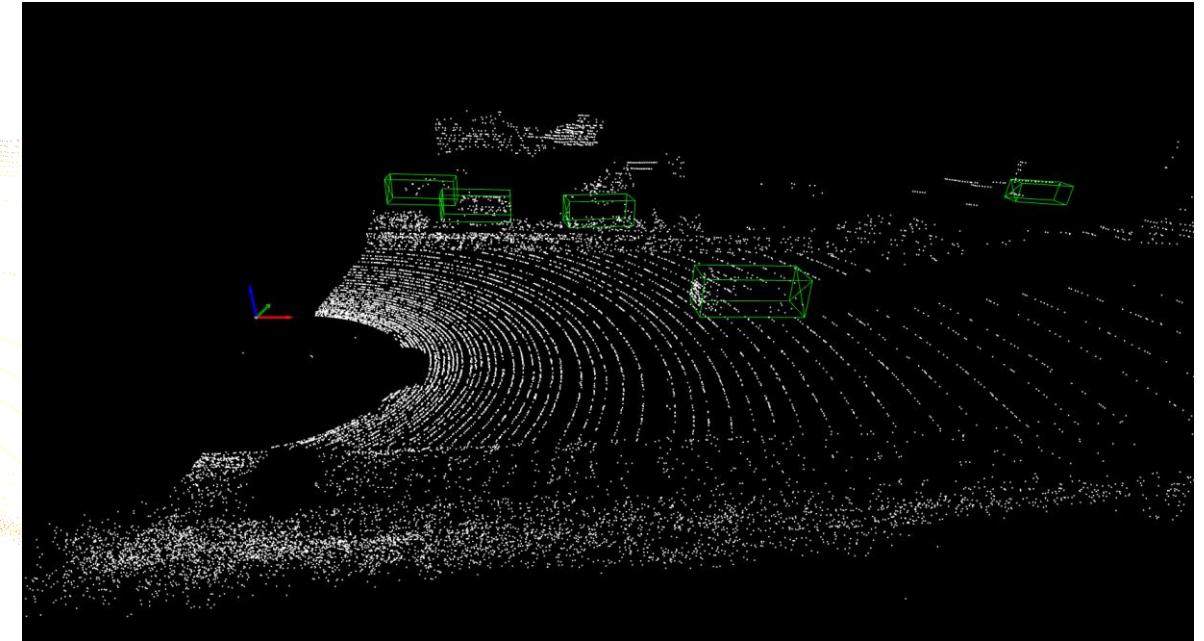
Metric	Value	Units
Average Translation Error	0.328	%
Average Rotational Error	0.164	deg/m
Absoulte Trajectory Error (ATE)	0.776	m
Absoulte Rotational Error (ARE)	0.007	rad
Average Frequency	69	Hz
Average Runtime	14	ms



# Workflow to Optimize 1 Car – 1 Frame

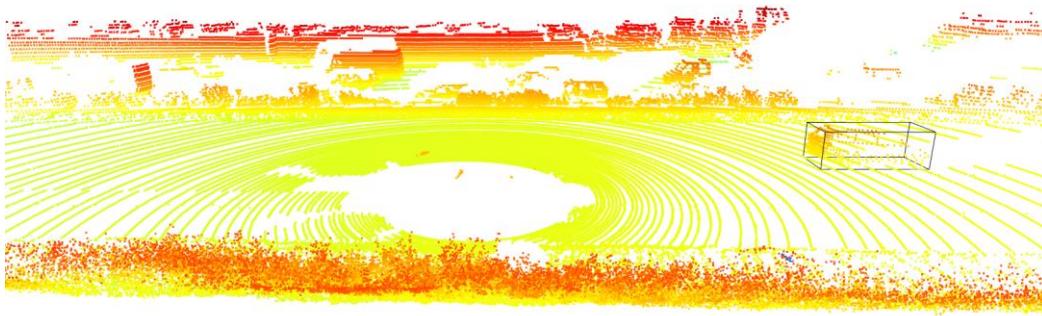


Raw Sensor Data

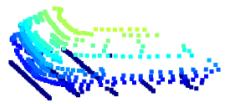
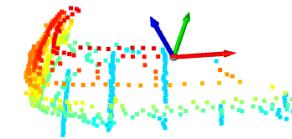


Bounding Box Predictions

# Workflow to Optimize 1 Car – 1 Frame

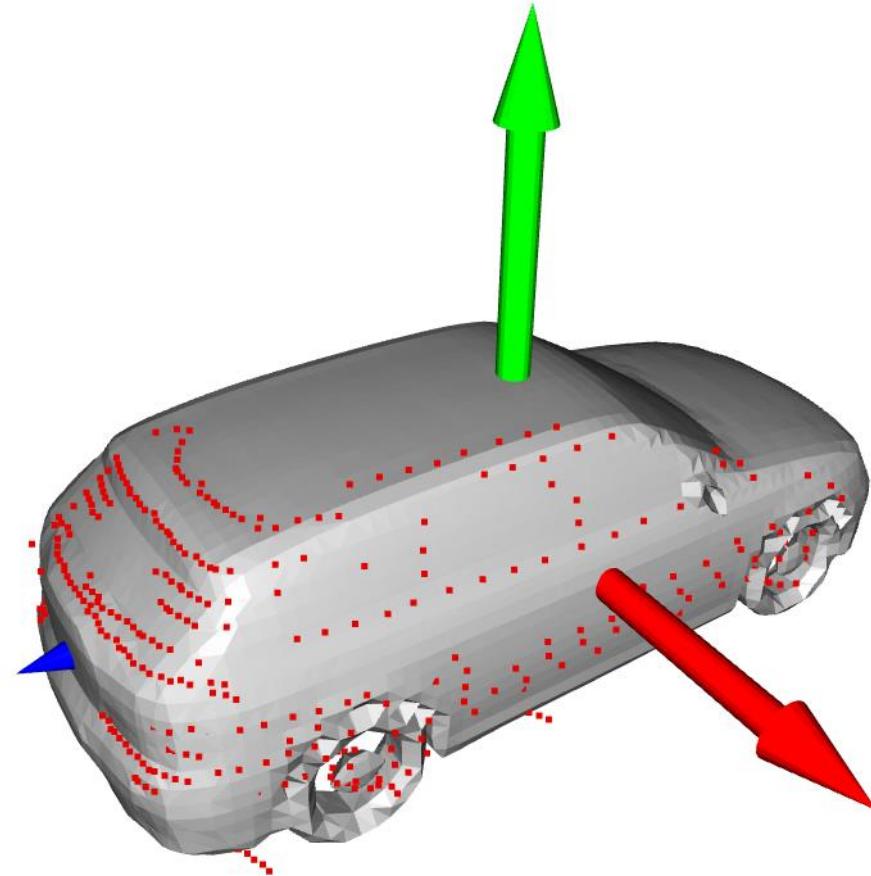


Instance Association for  
Bounding Boxes

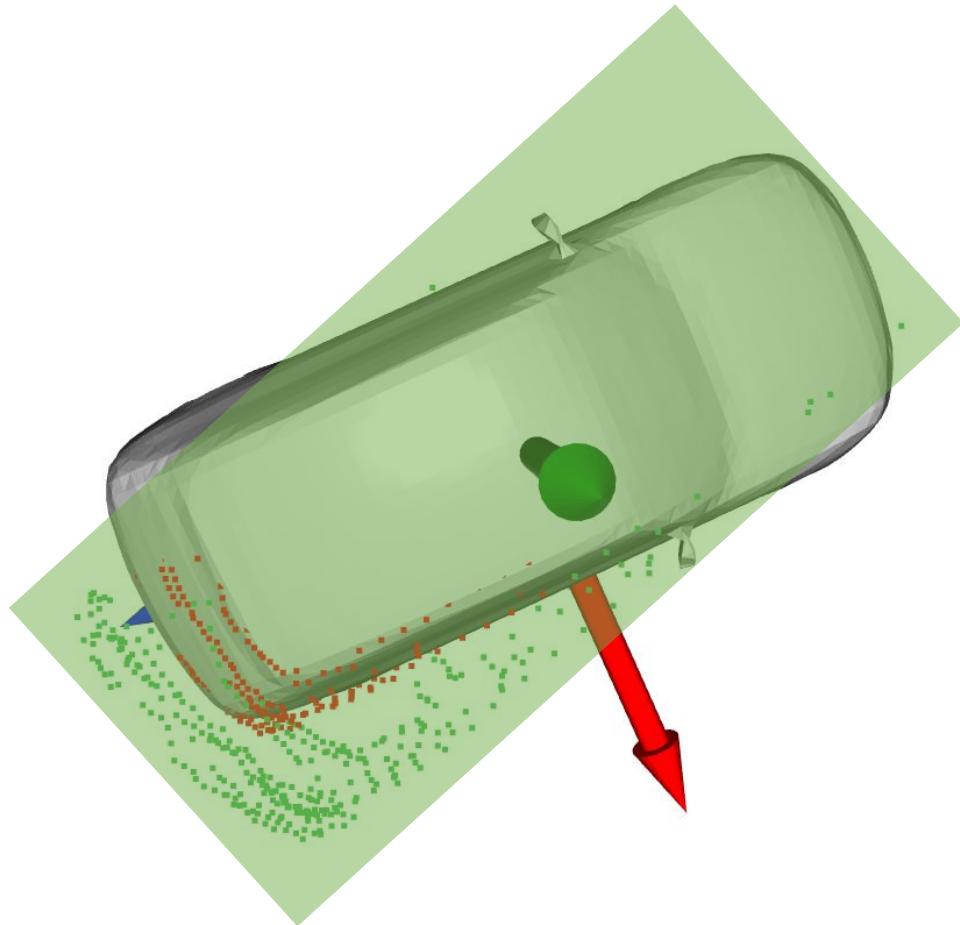


Isolating Car Surface Points

# Workflow to Optimize 1 Car – 1 Frame



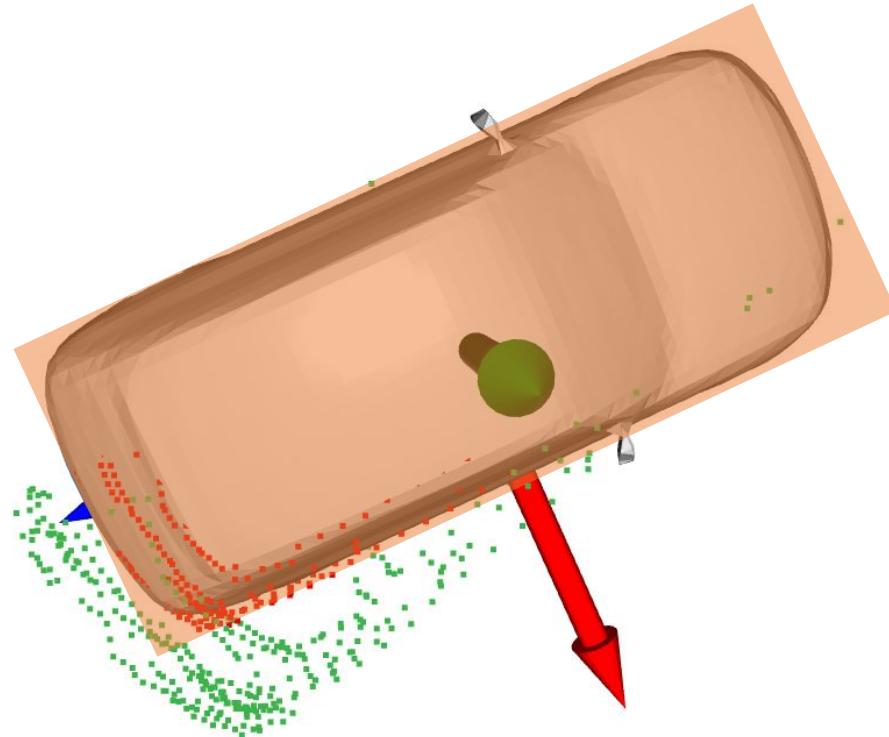
# Workflow to Optimize 1 Car – 1 Frame



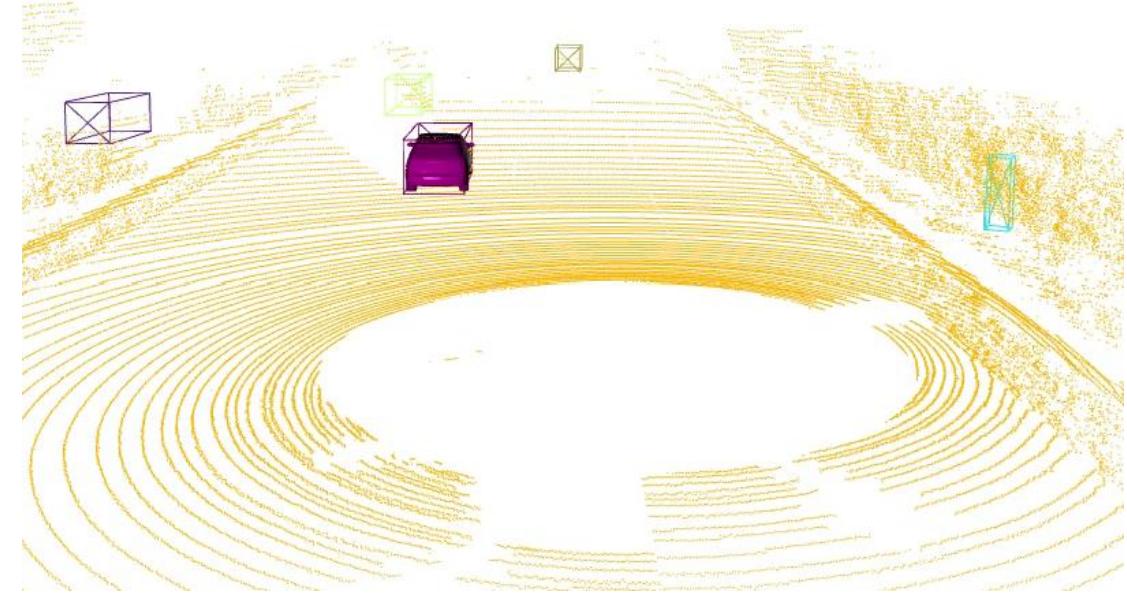
The mesh is the reconstructed object surface using marching cubes algorithm

**Sim(3) optimization**

# Workflow to Optimize 1 Car – 1 Frame

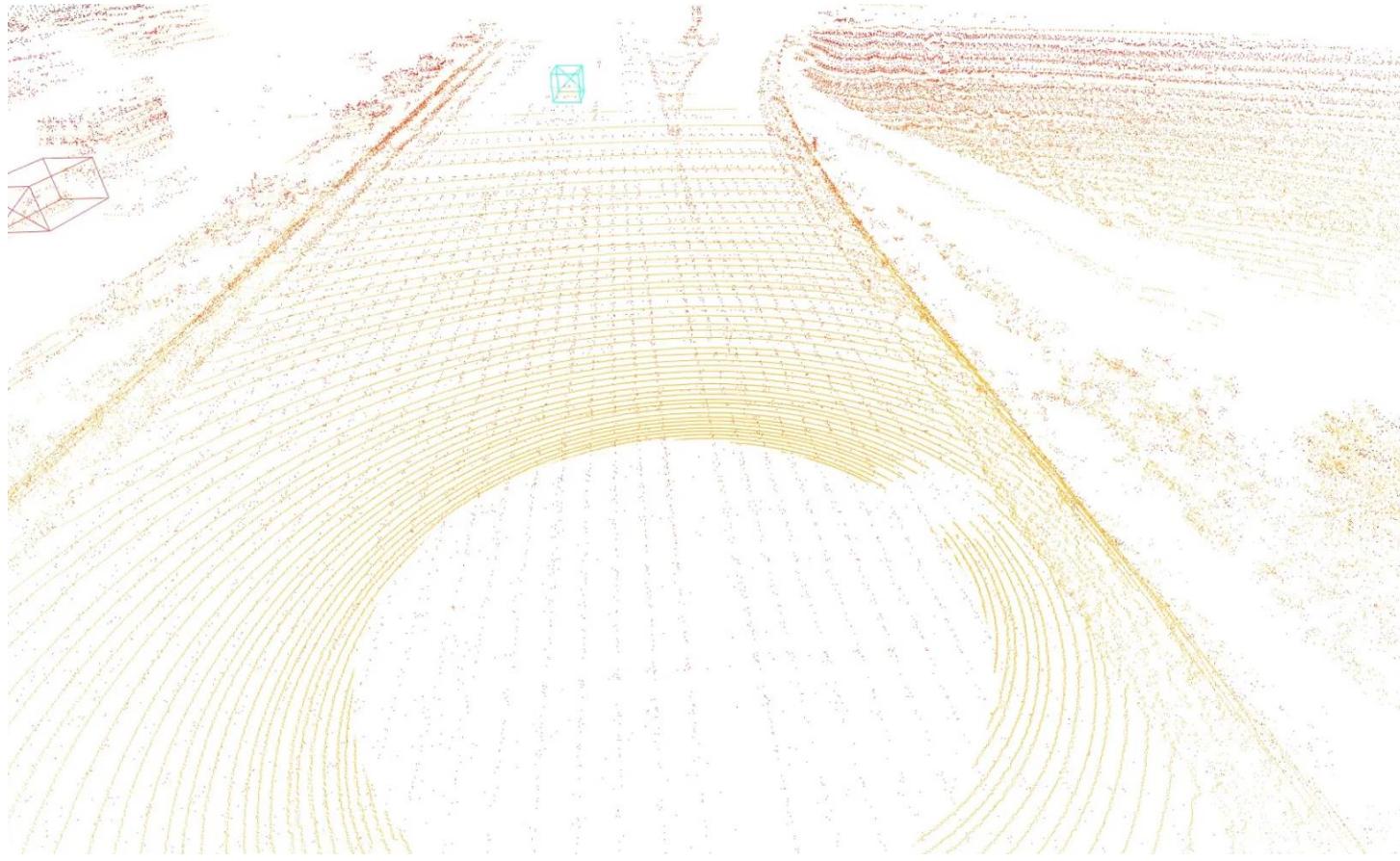


**Sim(3)** optimization



Add the mesh to the visualizer at  
the optimized pose

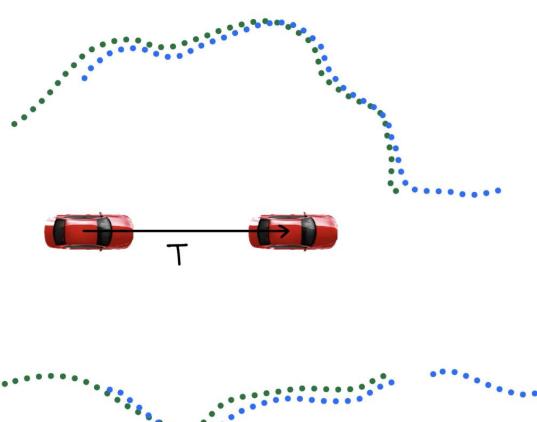
# Optimizing for Multiple Cars



Instance Completion and Motion Estimation with Instance Association

# Process of the project

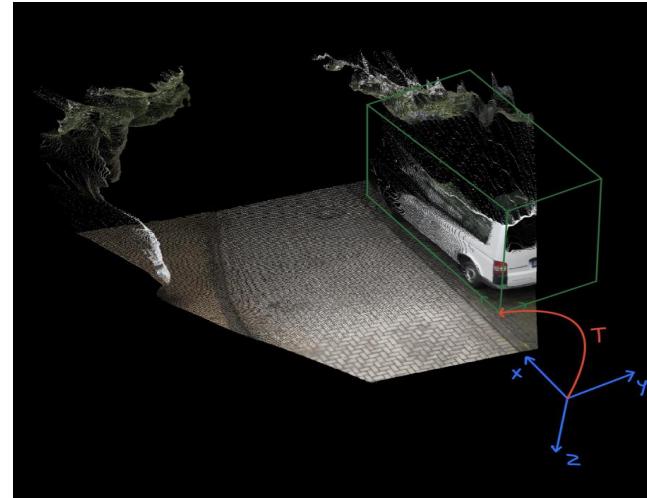
## Lidar Odometry



“Keep it small and simple”

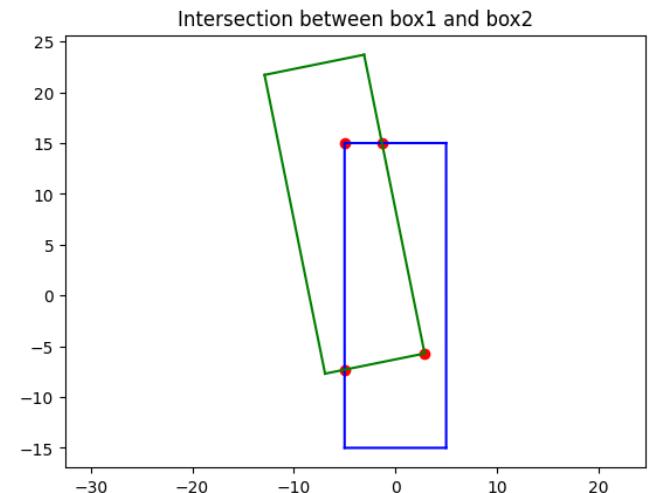
- Accurately compute a robot's pose.
- Point cloud alignment.
- A few parameters tuning.
- Assumption: The output of KISSICP is accurate.

## Bounding Box



Point RCNN gives an initial guess for bounding box in sensor frame

## Instance Association



Metrics Used:

1. Simple IOU
2. Distance in the direction of motion

# MILESTONE 1 PRESENTATION

# Instance Completion and Motion Estimation with Deep Shape Priors for Autonomous Driving

## Team members

Panyawat (Ohm) Rattana  
Shashank (Sai) Dammalapati

## Supervisors

Xingguang (Starry) Zhong  
Yue Pan

# Instance Completion and Motion Estimation with Deep Shape Priors for Autonomous Driving

## Team members

Panyawat (Ohm) Rattana  
Shashank (Sai) Dammalapati

## Supervisors

Xingguang (Starry) Zhong  
Yue Pan

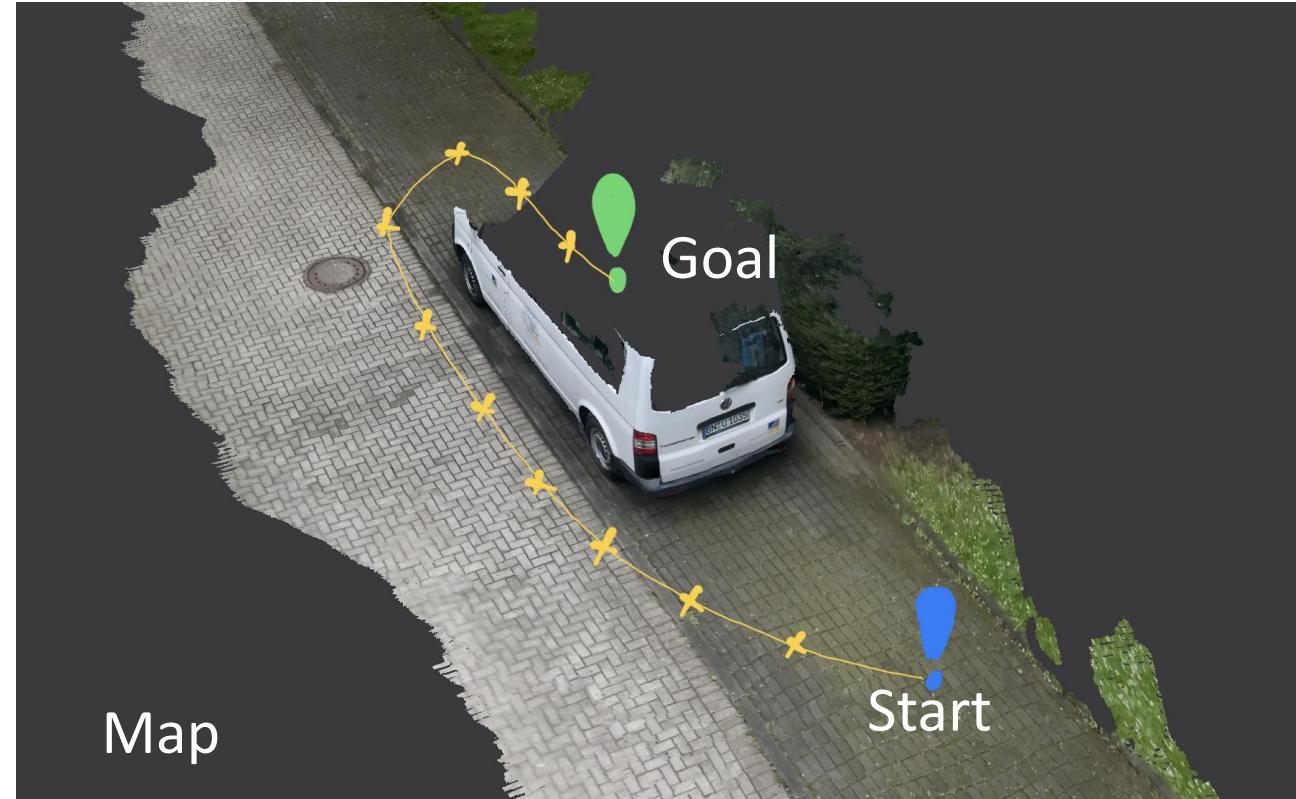
# Where is this scene from?



# Agent travelling in the safe space



# Task of a path planner

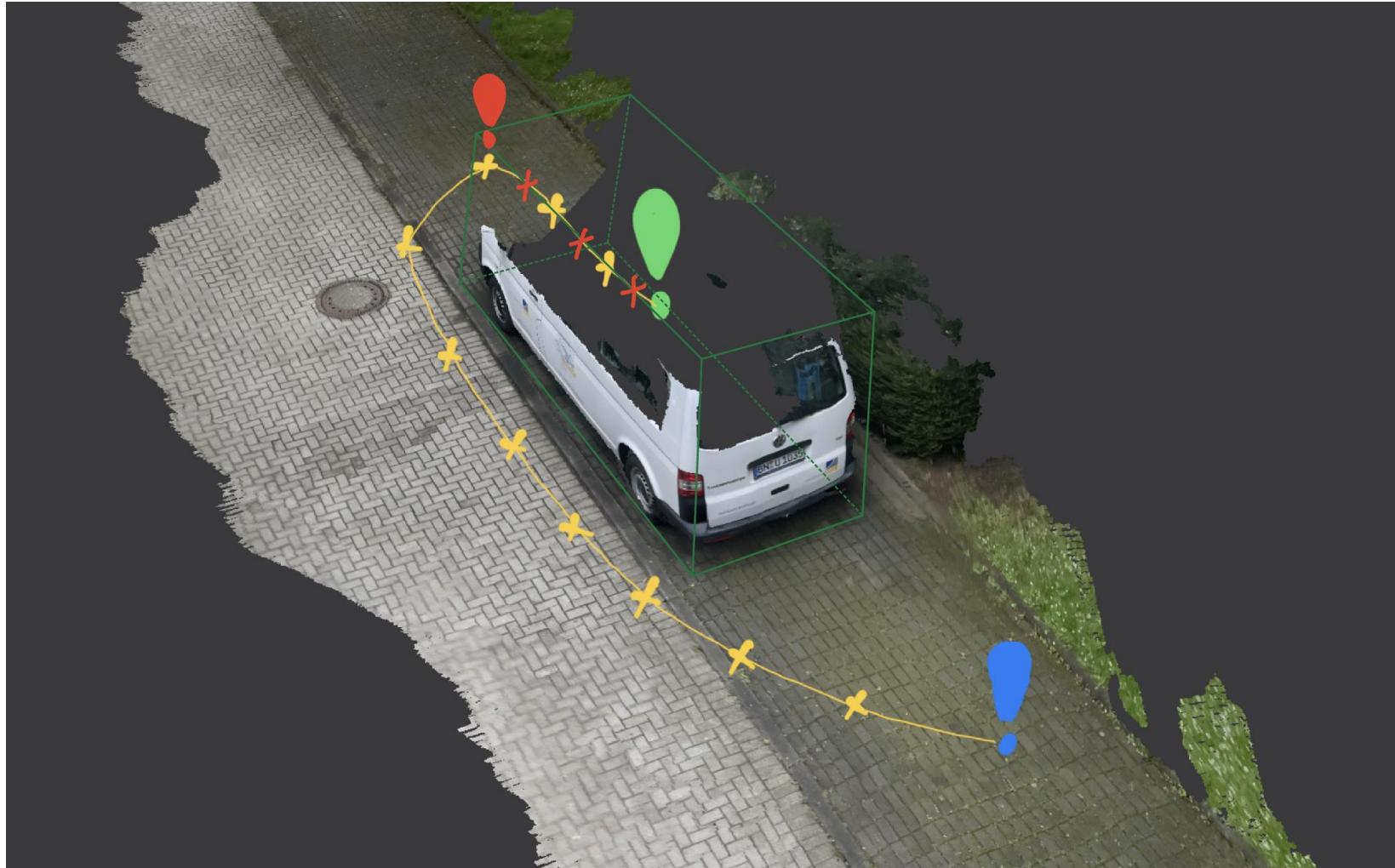


# What's wrong with this plan?



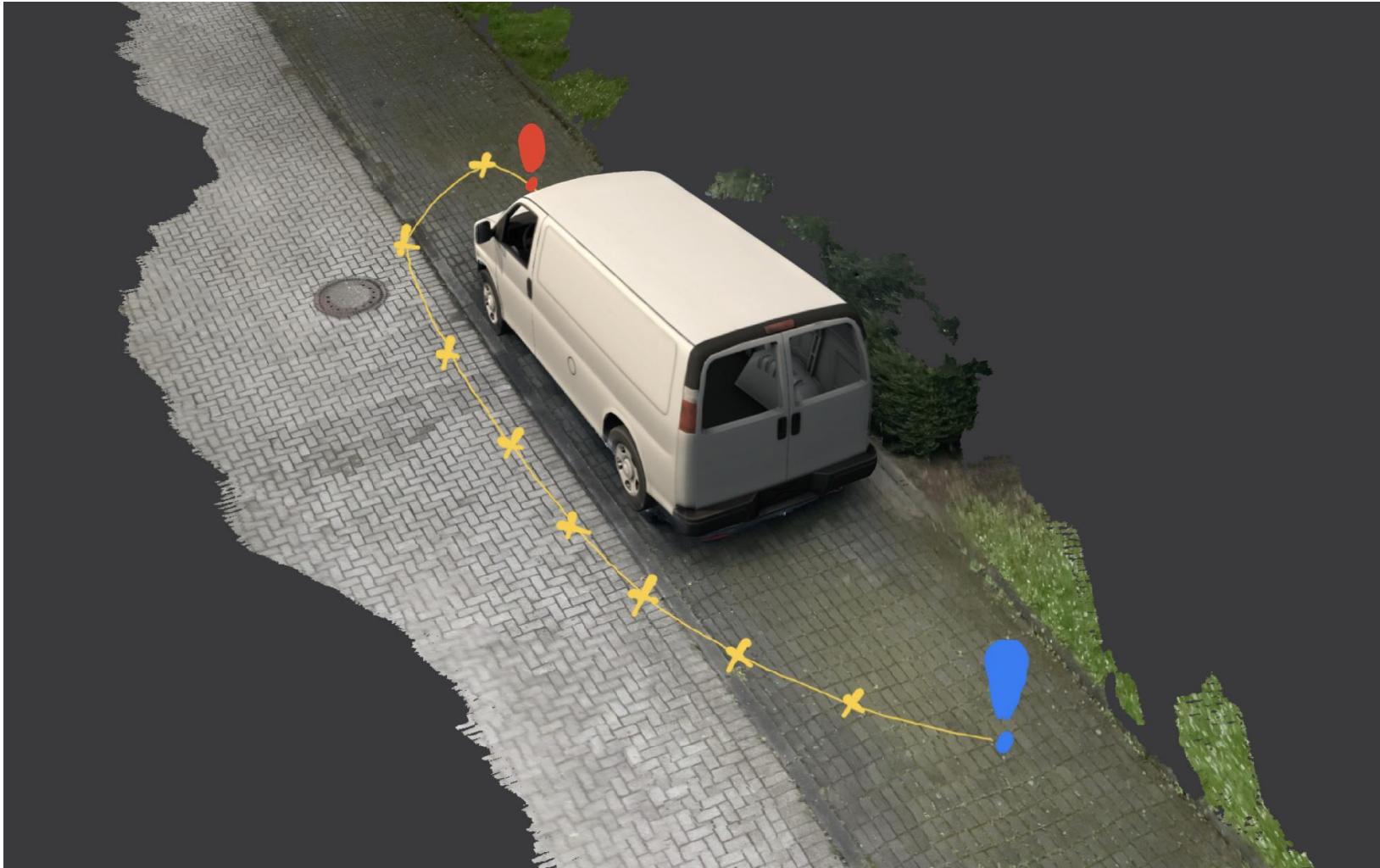
Mesh map of the scene

# This is a better representation of map



Scene with bounding boxes for objects

We argue that this representation is even better

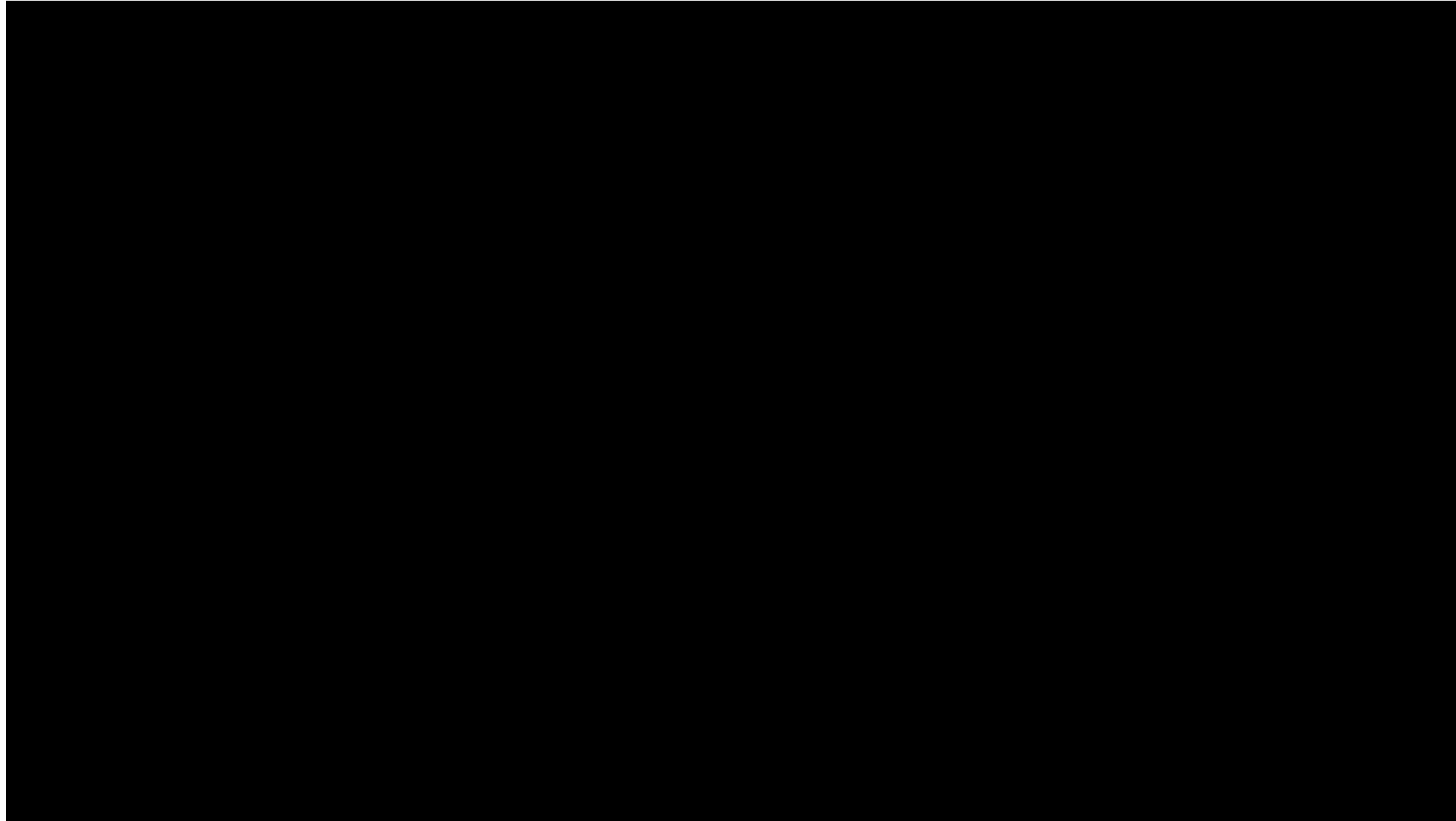


Scene with instance completion of objects



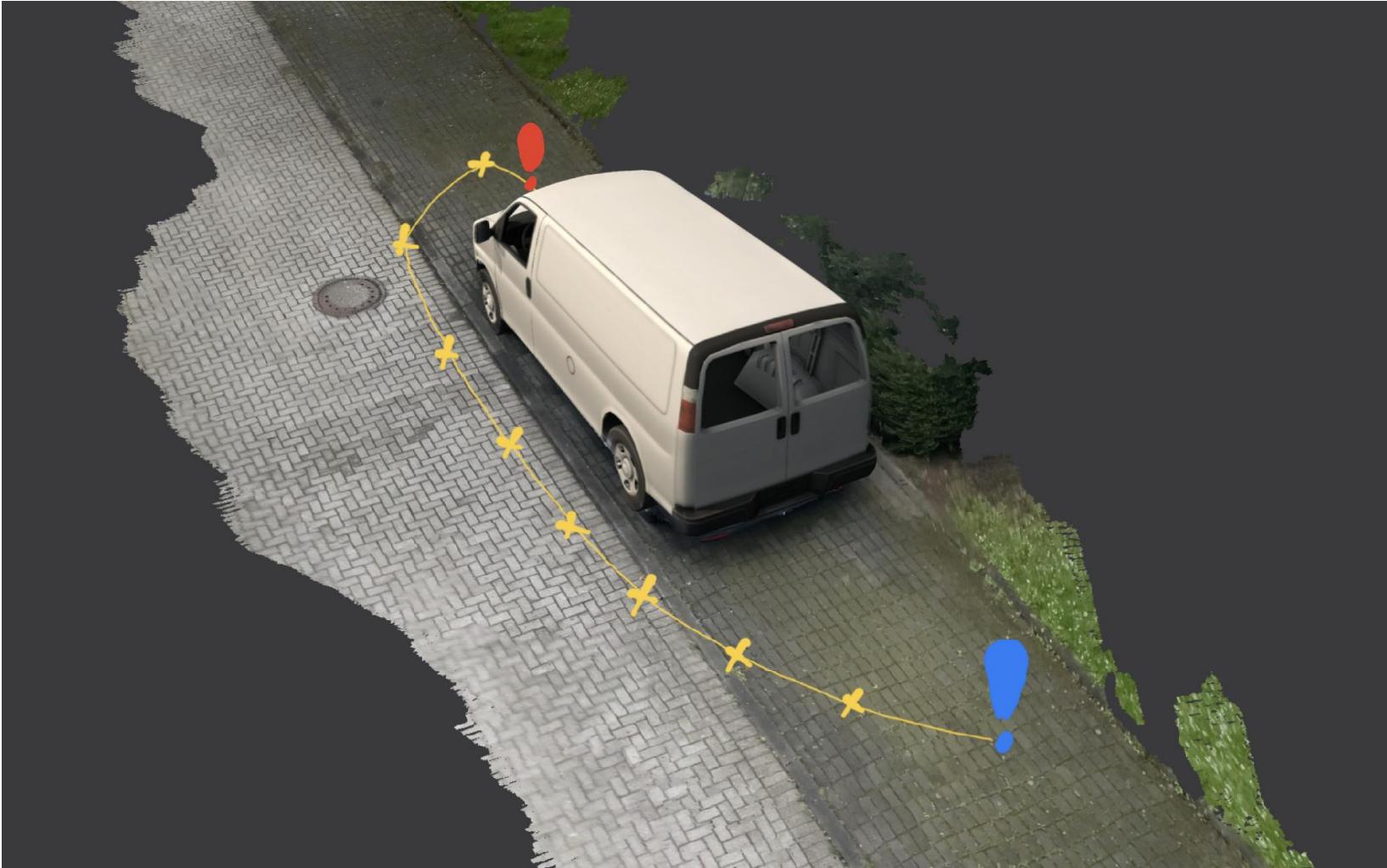
Source: Kitti Raw Dataset

# Existing method: DSP SLAM

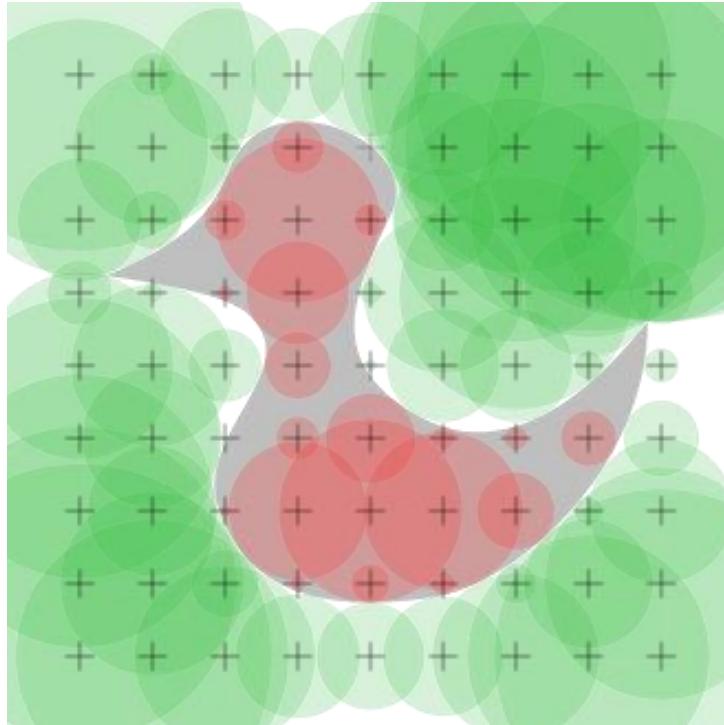


But, can we do better?

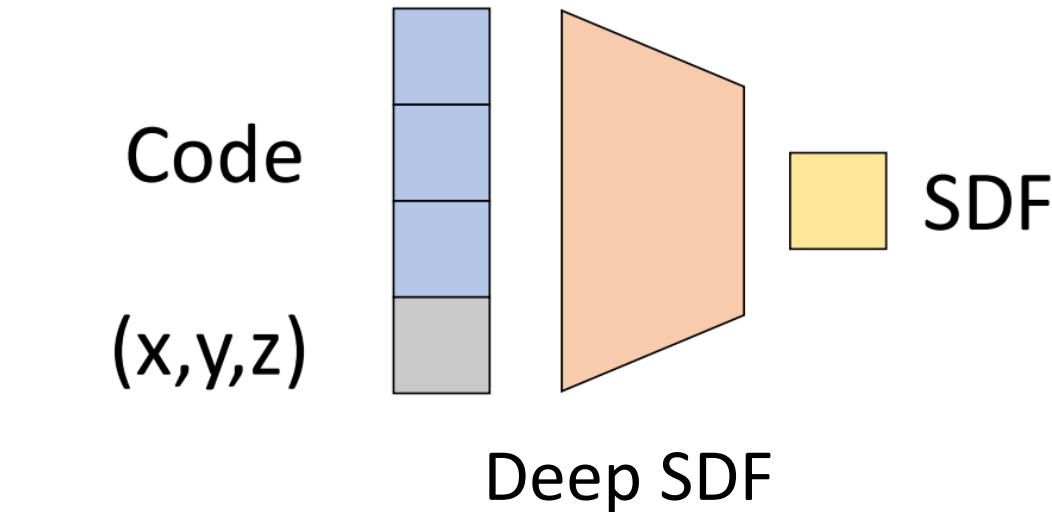
# What makes this possible?



Scene with instance completion of objects

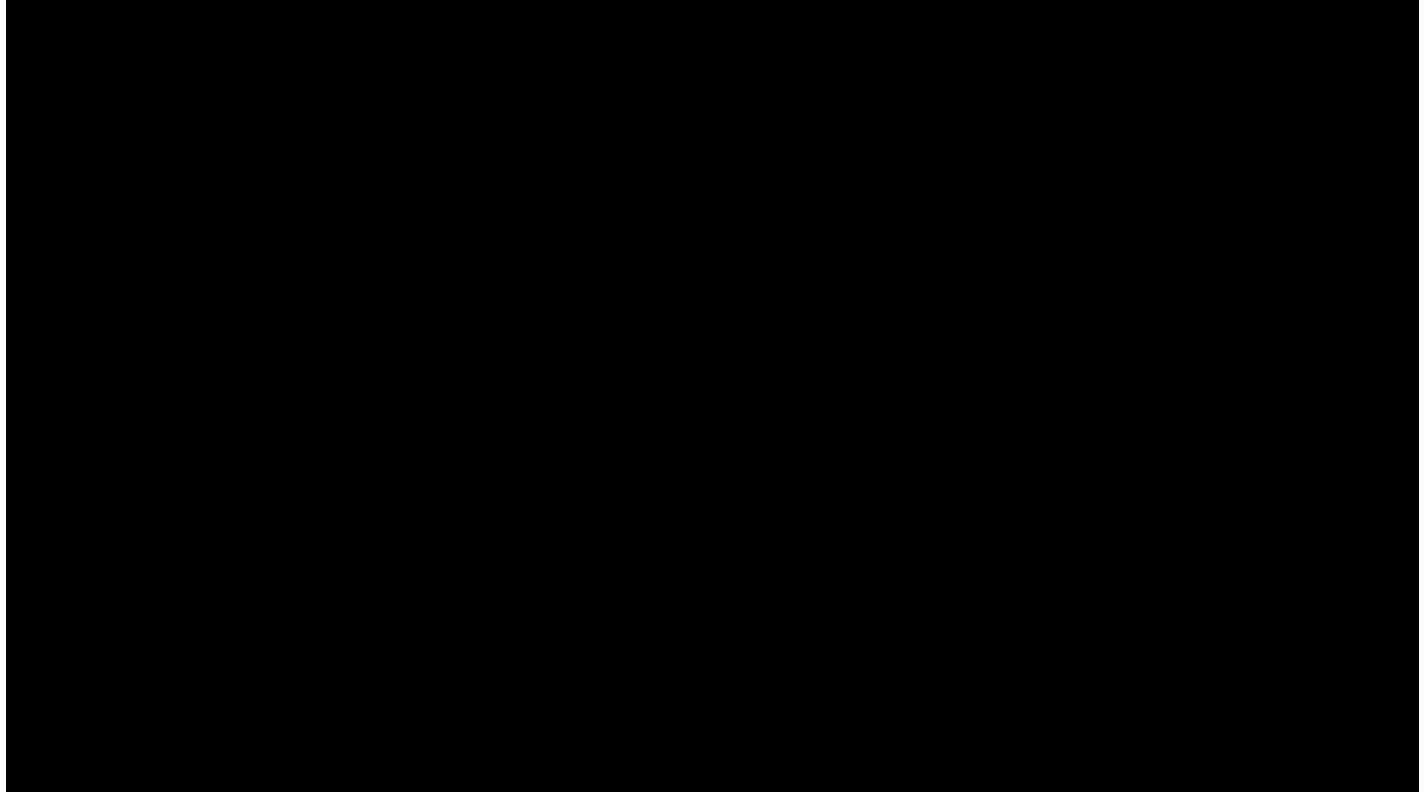
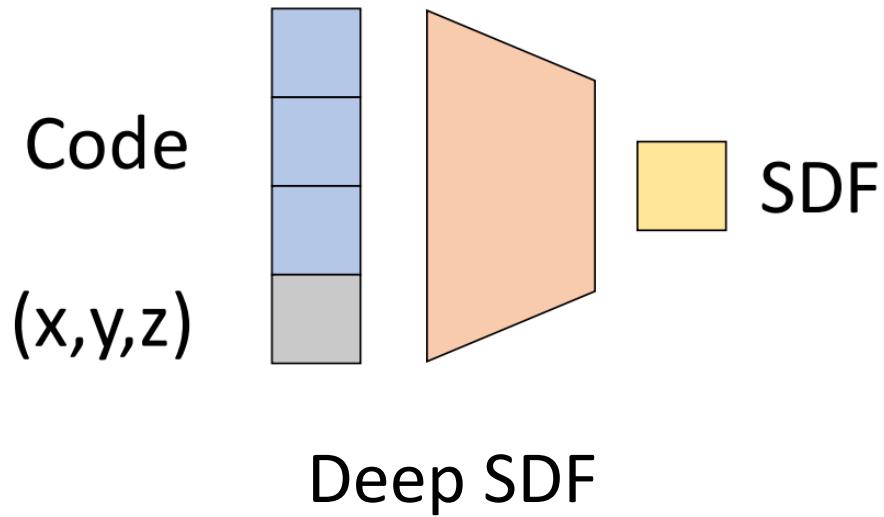


SDF  $f_{\theta}(x, y, z) \approx SDF(x, y, z)$



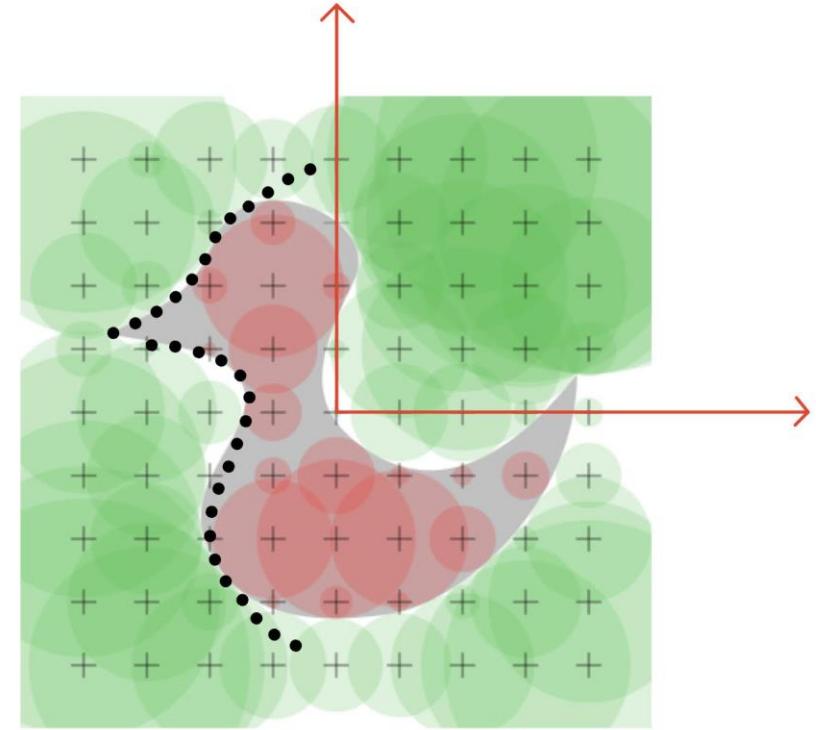
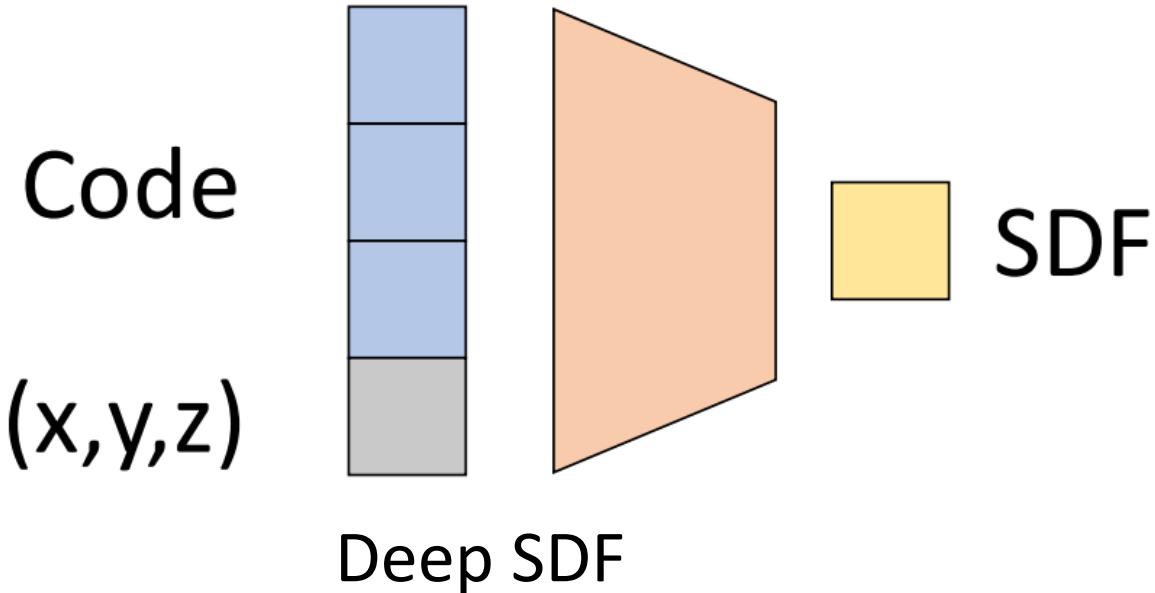
DeepSDF can implicitly model a class of objects

# Deep SDF



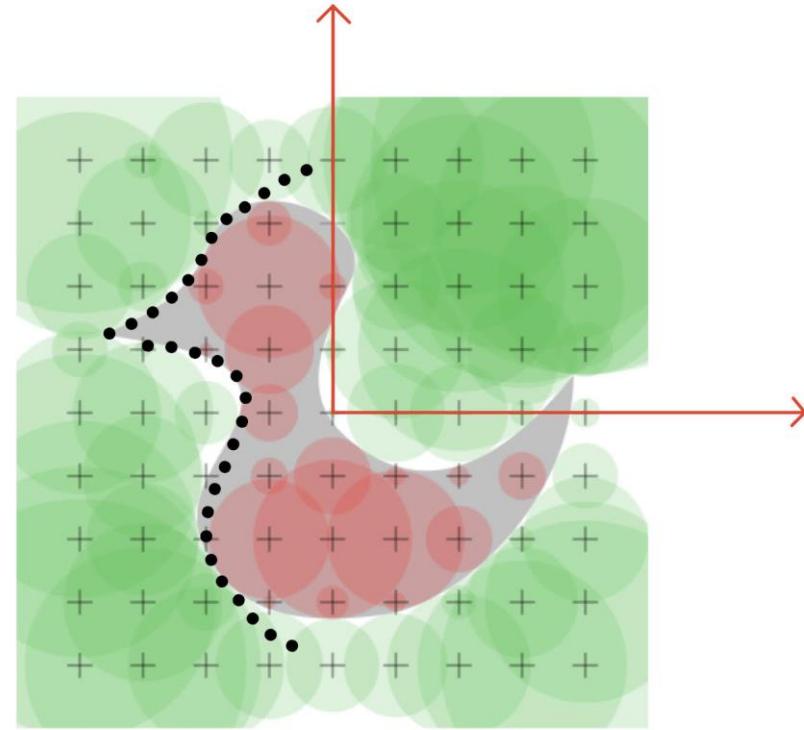
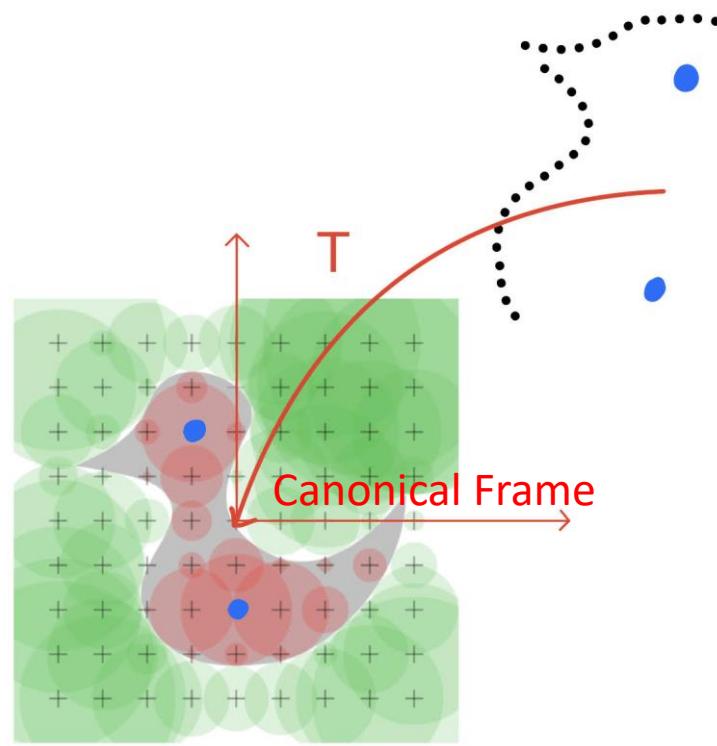
A latent code collapses a general representation into a single shape

# 0 SDF



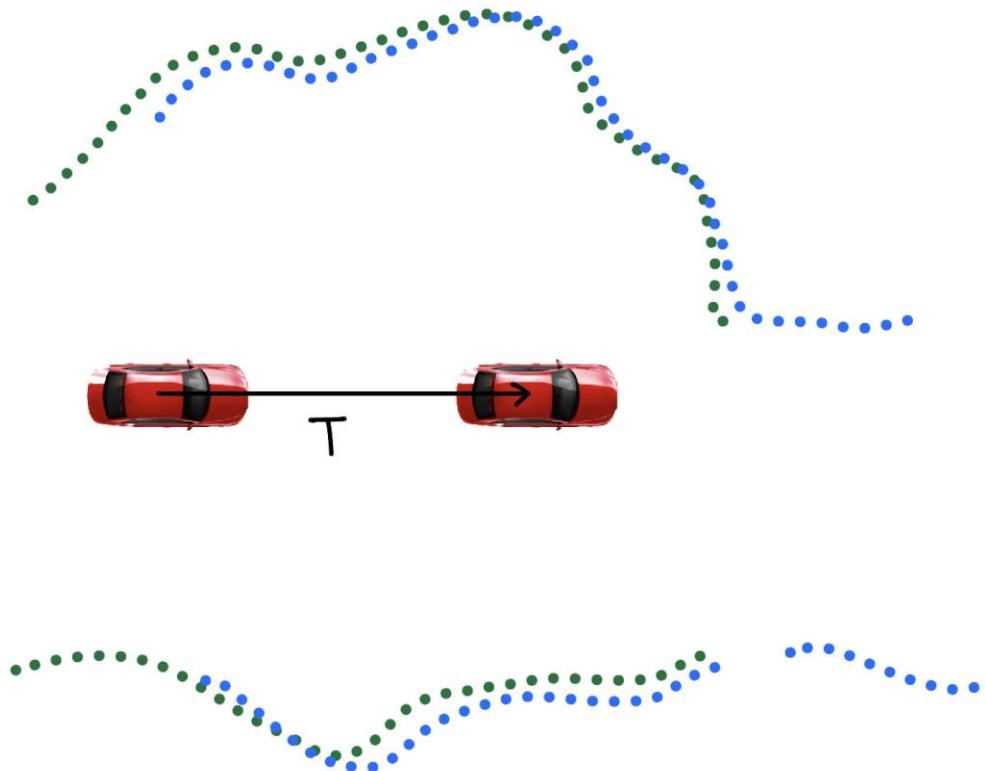
SDF value for a point on surface modelled by SDF is ZERO

# Canonical Frame



SDF presumes that the input is in Canonical Frame

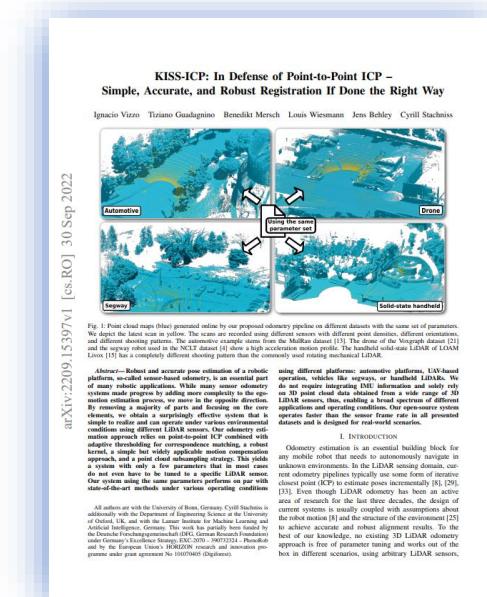
# Lidar Odometry (KISS ICP)



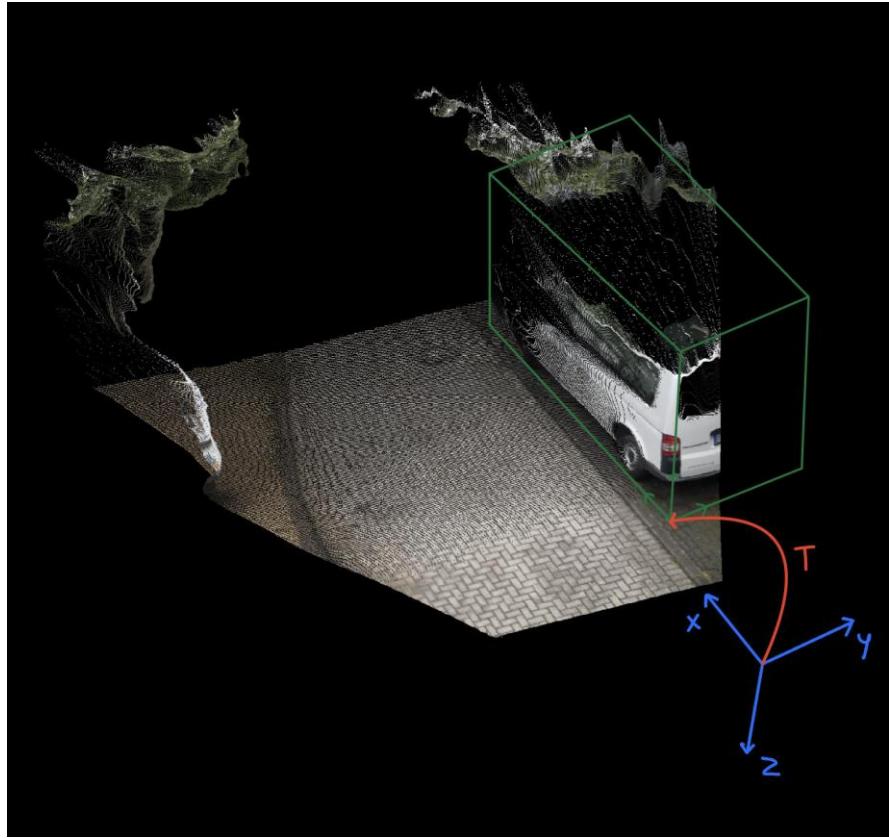
Motion Estimation of Ego Car

“Keep it small and simple”

- Accurately compute a robot’s pose.
- Point cloud alignment.
- A few parameters tuning.
- Assumption: The output of KISSICP is accurate.



# Point RCNN



Pose Estimation of surrounding cars in 1 frame

- Idea is to generate bounding boxes for cars
- Point RCNN gives us an initial guess for the relative pose of the external cars relative to ego car.

**PointRCNN: 3D Object Proposal Generation and Detection from Point Cloud**

Shaochuai Shi Xiaogang Wang Hongsheng Li  
The Chinese University of Hong Kong  
{ssshui, xgwang, hsl1}@ee.cuhk.edu.hk

**Abstract**

In this paper, we propose PointRCNN for 3D object detection from raw point cloud. The whole framework is composed of two stages: stage-1 for the bottom-up 3D proposal generation and stage-2 for refining proposals in the top-down manner to obtain the final detection results. Instead of generating proposals from RGB image or projecting point cloud to bird's view or voxels as previous methods do, our stage-1 sub-network directly generates a set of coarse and highly overlapping 3D local point cloud in a bottom-up manner via segmenting the point cloud of the whole scene into foreground points and background points. Then, the generated segments are used points of each proposal to canonical coordinates to learn better local spatial features, which is combined with global semantic features of each point learned in stage-1 for accurate 3D bounding box regression. Extensive experiments on the 3D detection benchmark of KITTI dataset show that our proposed architecture outperforms state-of-the-art methods with large margins by using only point clouds as input. The code is available at <https://github.com/shaochuai/PointRCNN>.

**1. Introduction**

Deep learning has achieved remarkable progress on 2D computer vision tasks, including object detection [8, 32, 16] and instance segmentation [6, 10, 20], etc. Beyond 2D vision, 3D perception is also increasingly indispensable for many real-world applications, such as autonomous driving and domestic robots. While recent developed 2D detection algorithms are capable of handling large variations in the scene, they still face great challenges in the detection of 3D objects with point clouds still faces great challenges from the irregular data format and large search space of 6 Degrees of Freedom (DoF) of 3D object.

In autonomous driving, the main concern is to use 3D sensors are the LiDAR sensors, which generate 3D point clouds to capture the 3D structures of the scene. The difficulty of point cloud-based 3D object detection mainly lies in irregularity of the point clouds. State-of-the-art 3D detection methods either leverage the mature 2D detection frameworks by projecting the point clouds into bird's view [14, 42, 17] (see Fig. 1(a)), to the frontal view [4, 38], or to the regular 3D voxels [34, 43], which are not optimal and suffer from heavy computation cost during the detection process.

Instead of transforming point clouds to voxels or other regular data structures for feature learning, Qi *et al.* [26, 28] proposed PointNet for learning 3D representations directly from point clouds for 3D classification and segmentation. As shown in Fig. 1(b), their follow-up work [25] applied PointNet in 3D object detection to estimate the 3D bounding boxes based on the cropped frustum point cloud from the raw point clouds. However, the performance of the method heavily relies on the 2D detection performance and cannot take the advantages of 3D information for generating robust bounding box proposals.

Unlike object detection from 2D images, 3D objects in autonomous driving series are naturally and well separated

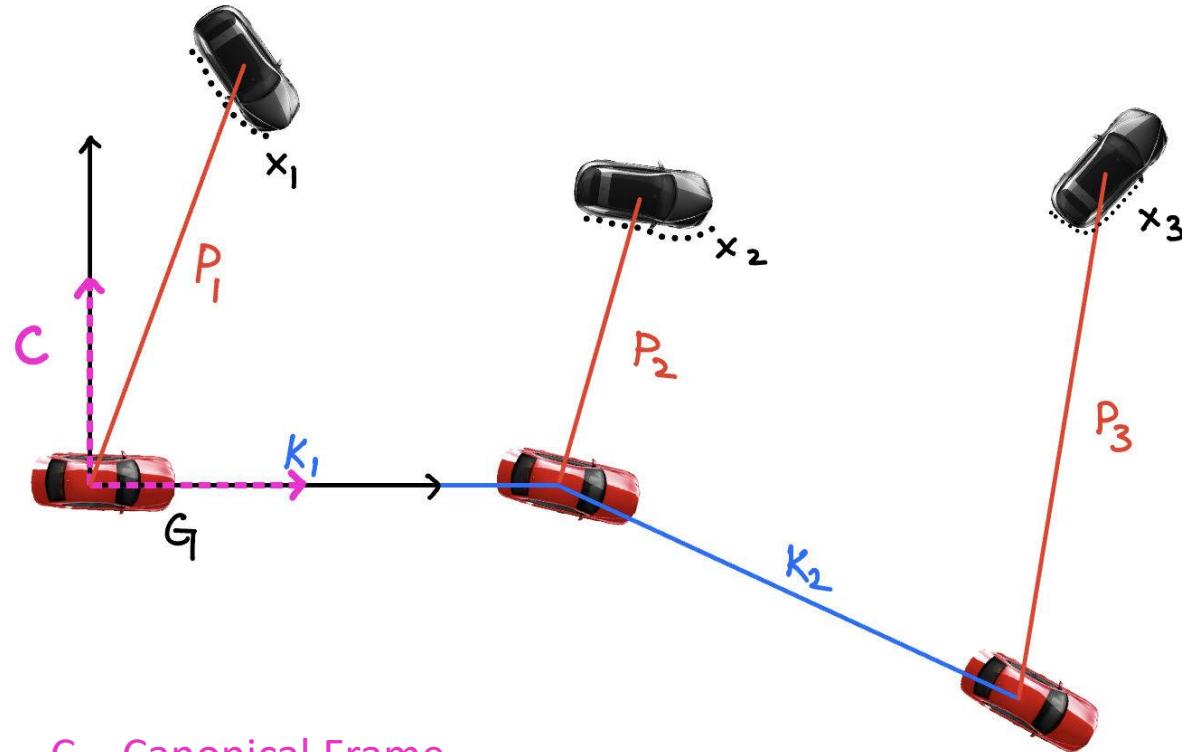
arXiv:1812.04244v2 [cs.CV] 16 May 2019

**a. Aggregate View Object Detection (AVOD)**

**b. Frustum-Pooling**

**c. Our approach (PointRCNN)**

# Workflow: Transformation



C – Canonical Frame

G – Global Frame

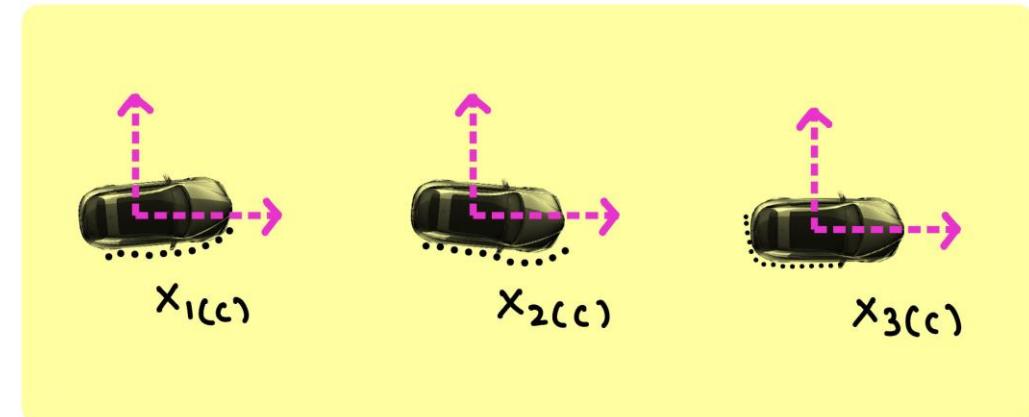
P – Transformation from Point RCNN

K – Transformation from KISS ICP

$$X_1[C] = P_1 X_1$$

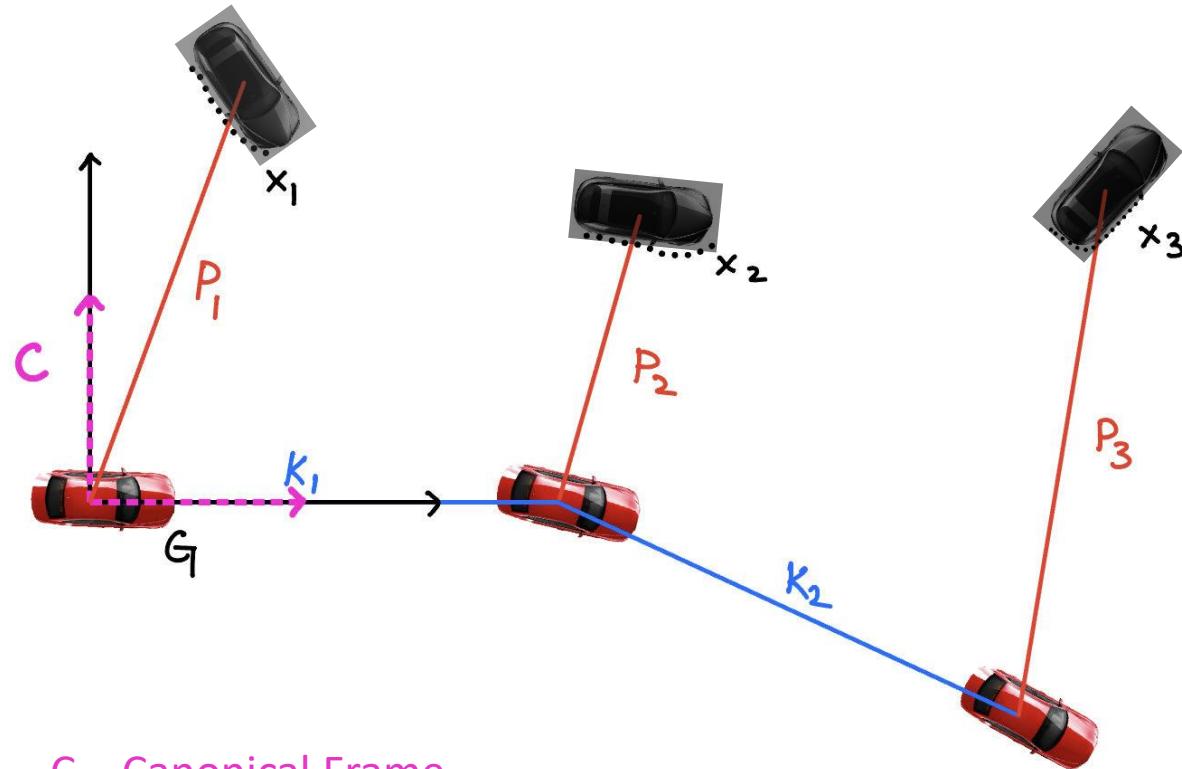
$$X_2[C] = K_1 P_2 X_2$$

$$X_3[C] = K_1 K_2 P_3 X_3$$



Canonical Space

# Workflow: Transformation



C – Canonical Frame

G – Global Frame

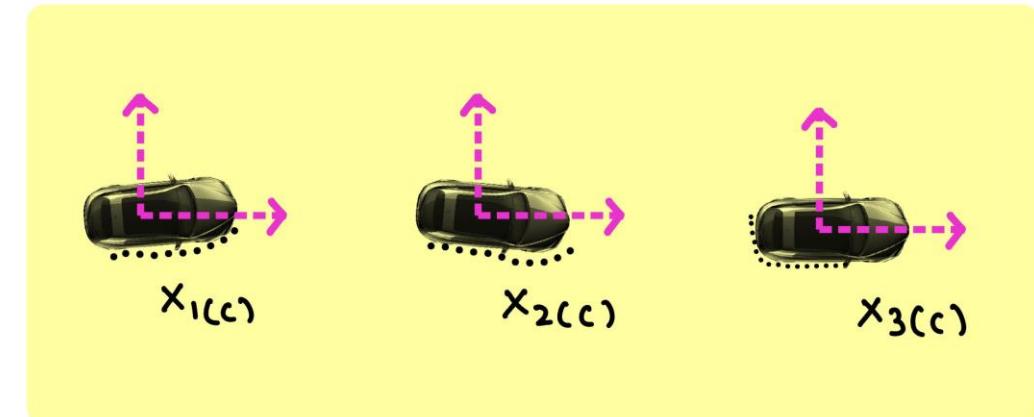
P – Transformation from Point RCNN

K – Transformation from KISS ICP

$$X_1[C] = P_1 X_1$$

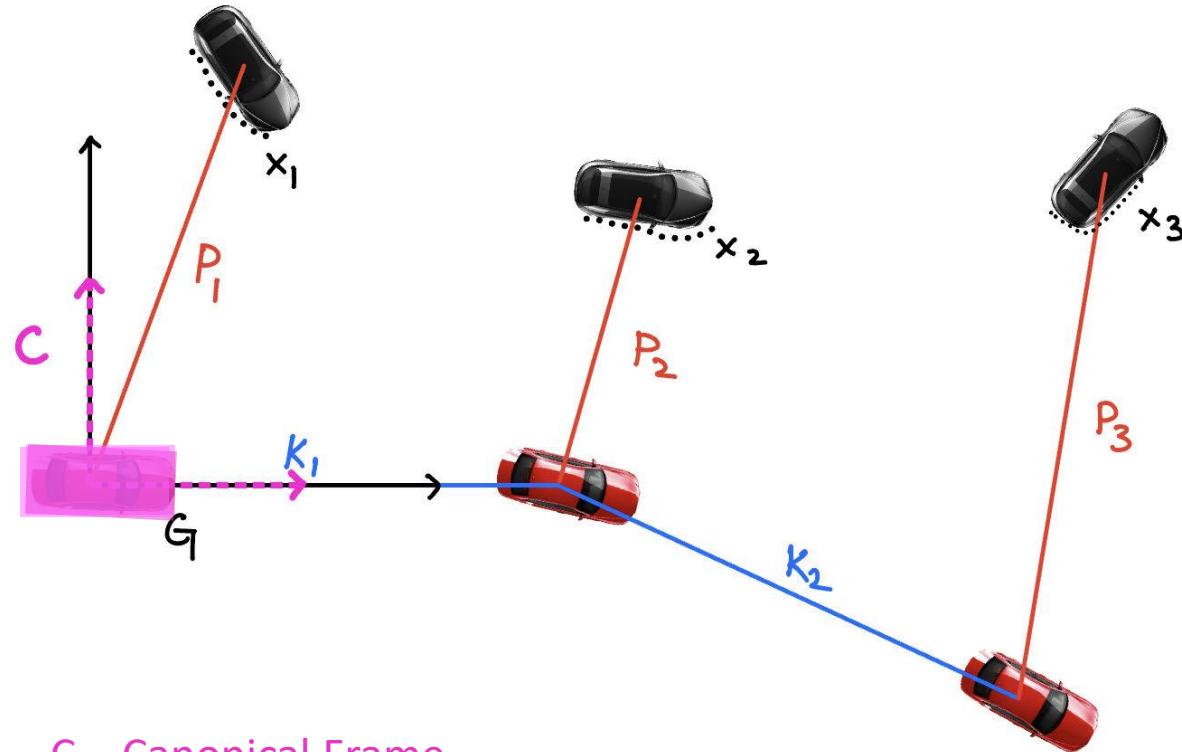
$$X_2[C] = K_1 P_2 X_2$$

$$X_3[C] = K_1 K_2 P_3 X_3$$



Canonical Space

# Workflow: Transformation



C – Canonical Frame

G – Global Frame

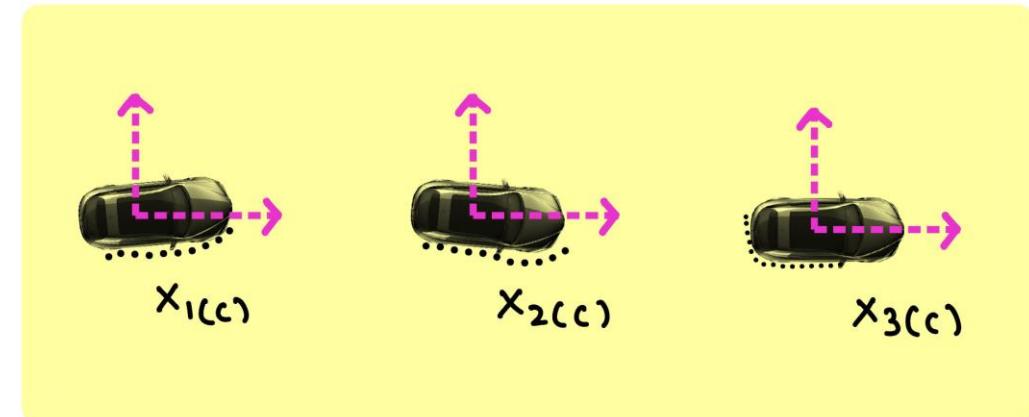
P – Transformation from Point RCNN

K – Transformation from KISS ICP

$$X_1[C] = P_1 X_1$$

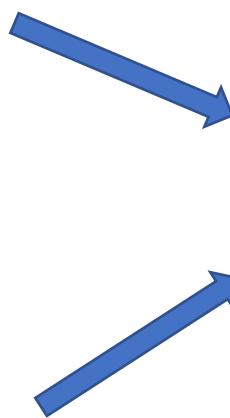
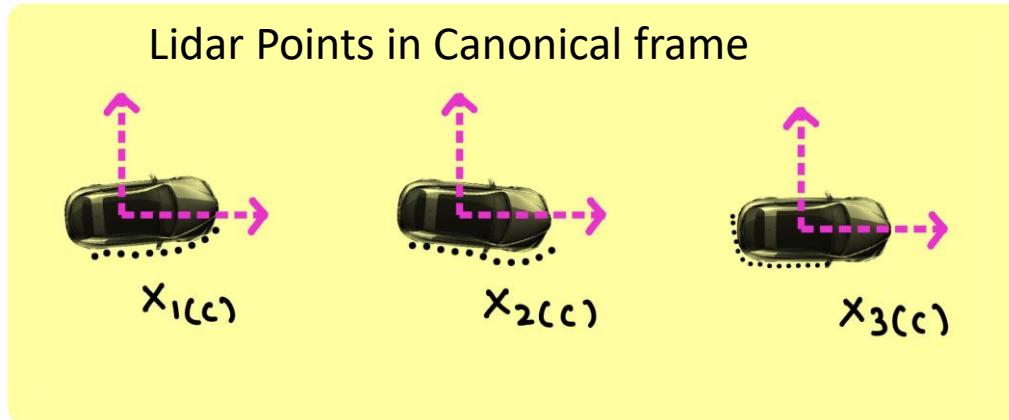
$$X_2[C] = K_1 P_2 X_2$$

$$X_3[C] = K_1 K_2 P_3 X_3$$



Point Clouds in Canonical Space

[Latent Code Vector].T

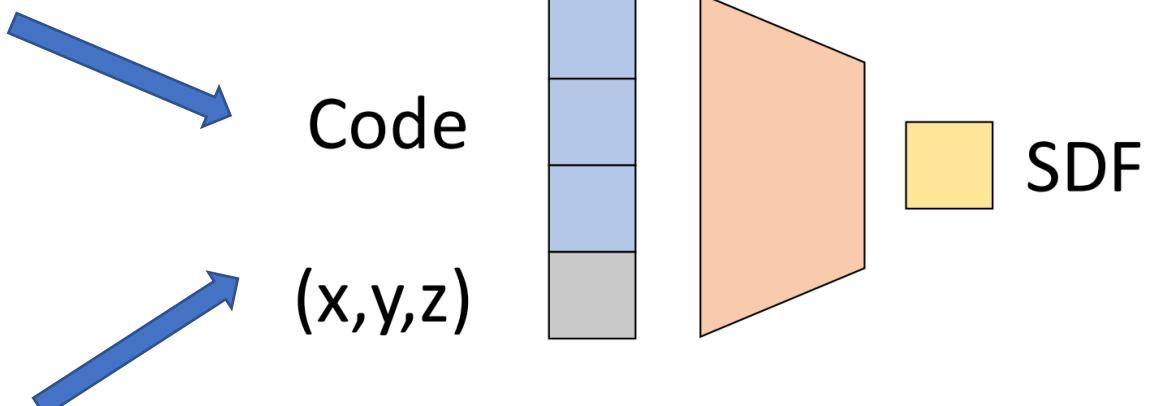
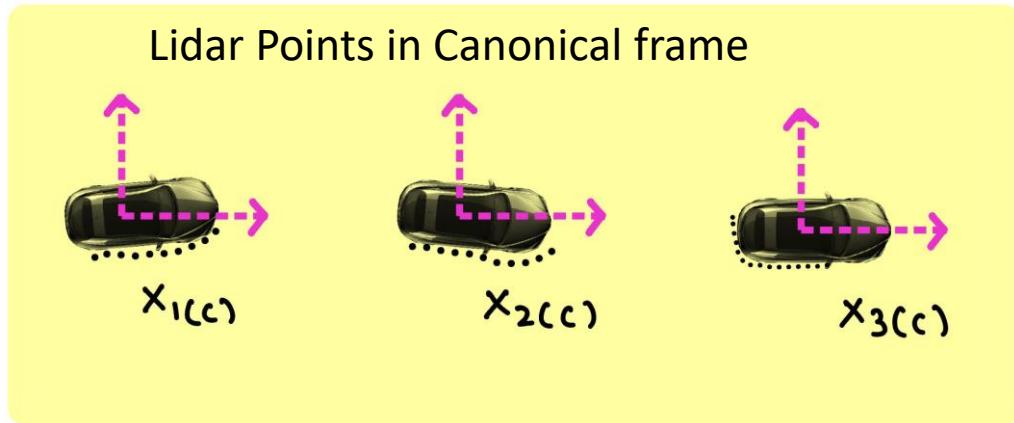


Code  
 $(x,y,z)$



SDF

[Latent Code Vector].T



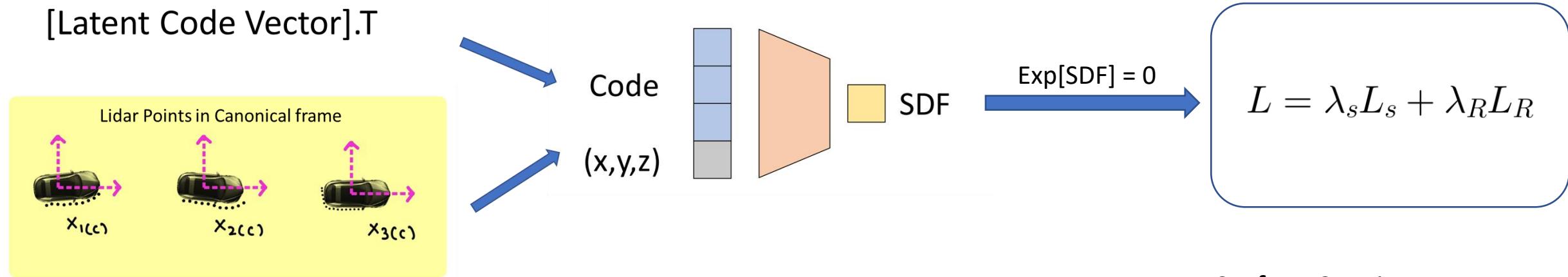
$$X_1[C] = P_1 X_1$$

$$X_2[C] = K_1 P_2 X_2$$

$$X_3[C] = K_1 K_2 P_3 X_3$$



Remember, points in canonical frame are a function of transformation predicted by PointRCNN

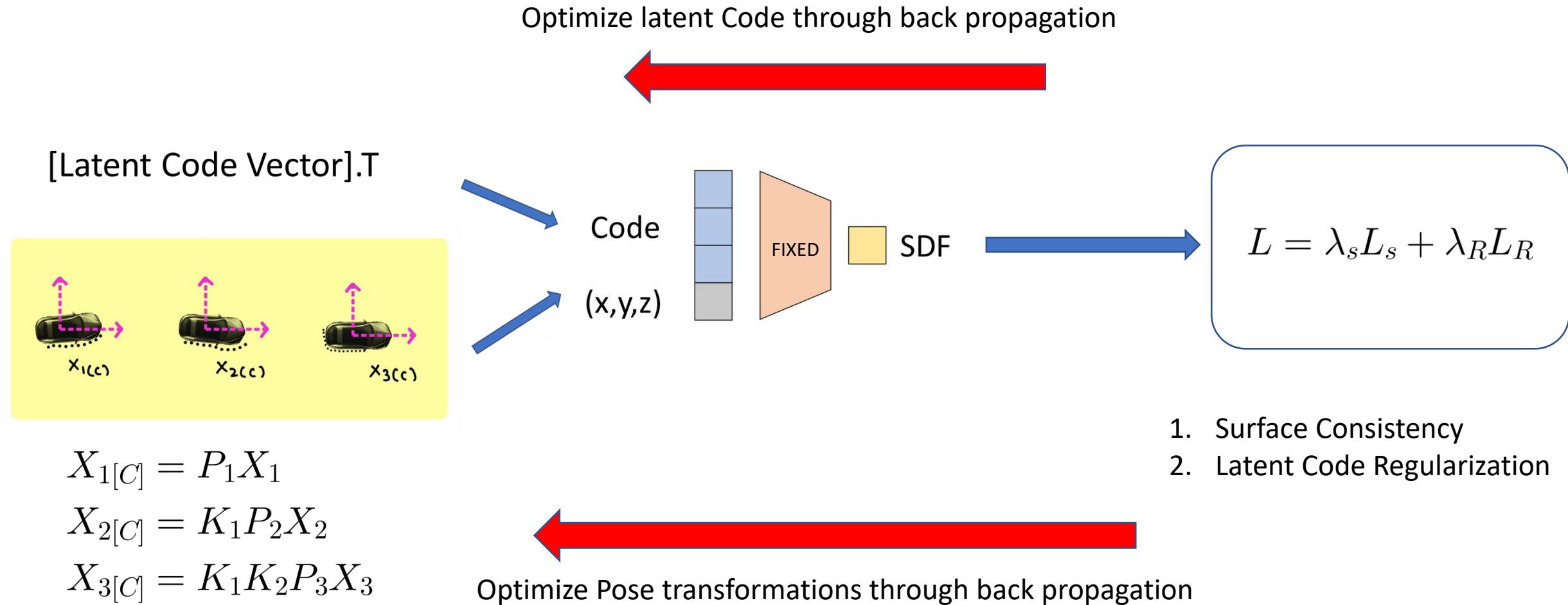


$$X_{1[C]} = P_1 X_1$$

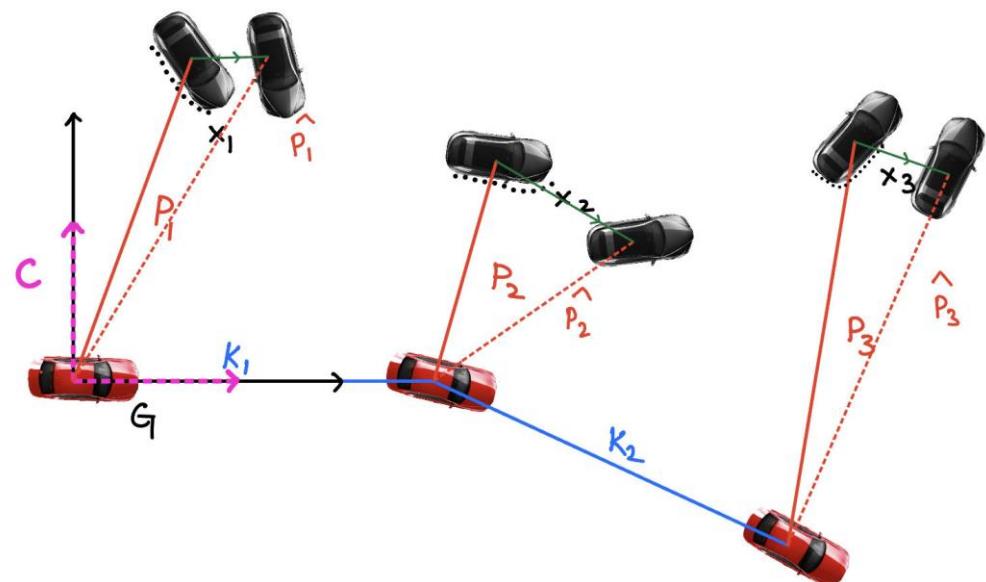
$$X_{2[C]} = K_1 P_2 X_2$$

$$X_{3[C]} = K_1 K_2 P_3 X_3$$

1. Surface Consistency
2. Latent Code Regularization

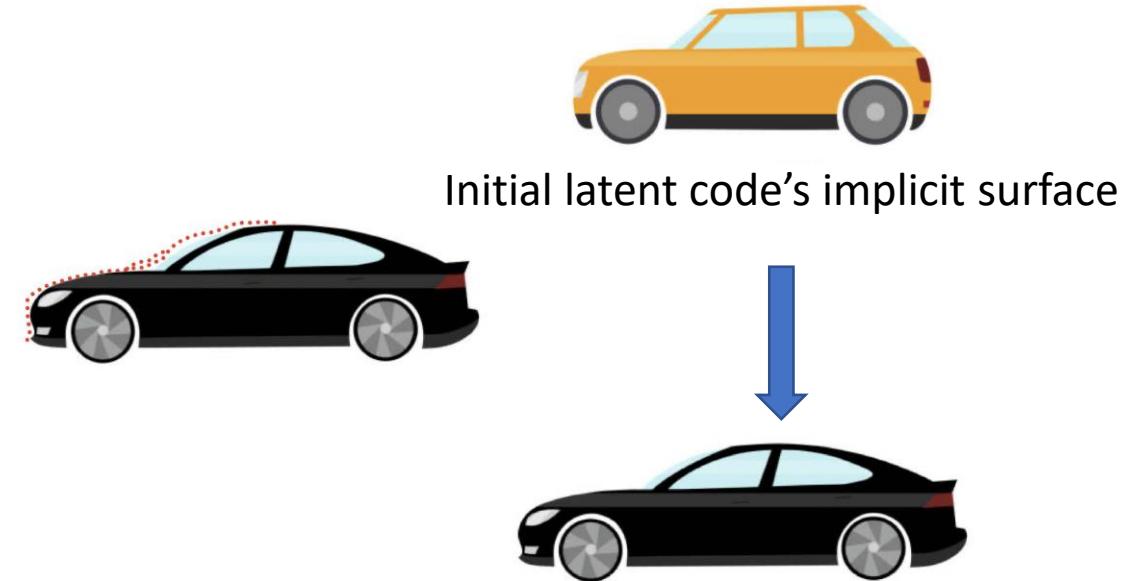


## Optimized Pose Estimation



Optimized poses of dynamic obstacles

## Optimized latent Code



Optimized latent code's implicit surface

# DataSets and Evaluation

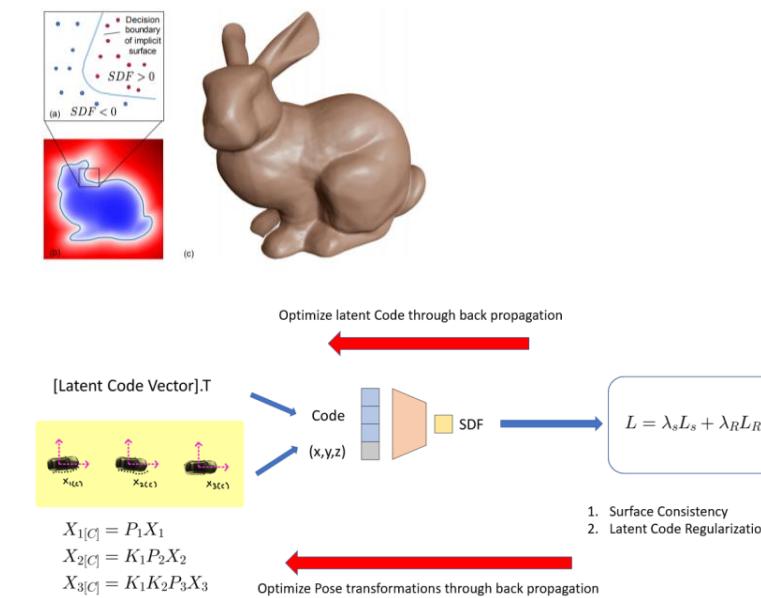
Dynamic Object Pose Estimation



# Conclusion



Challenges



Method



Final Output

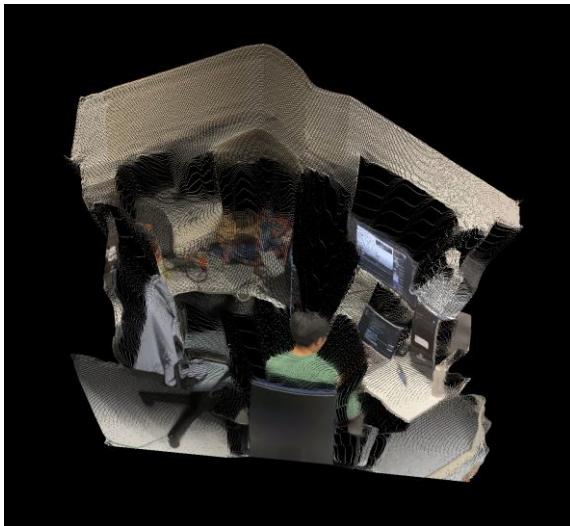
# Thank You



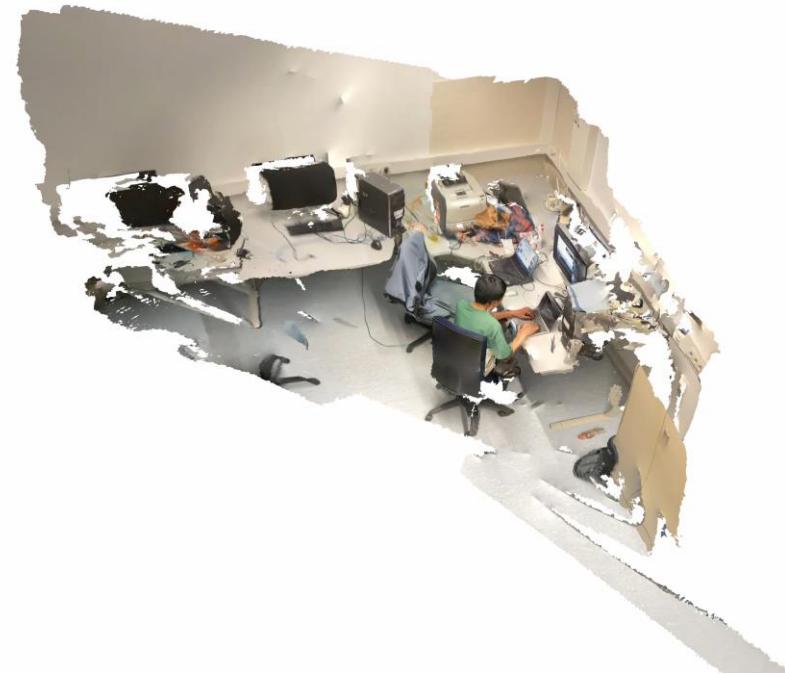
# 3D Reconstruction



Multiple Images

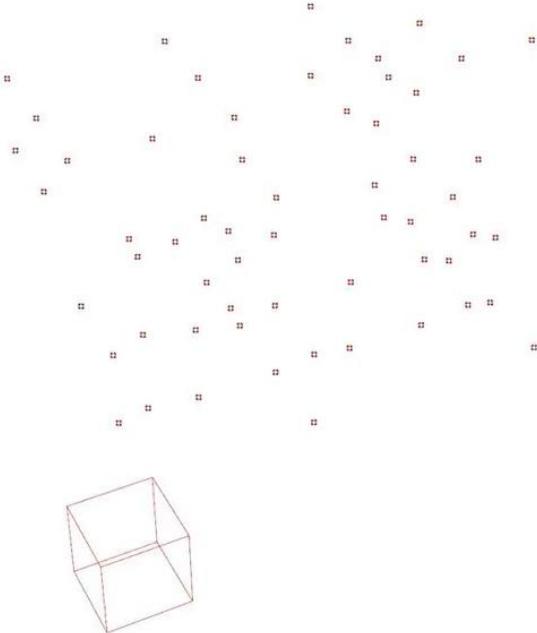
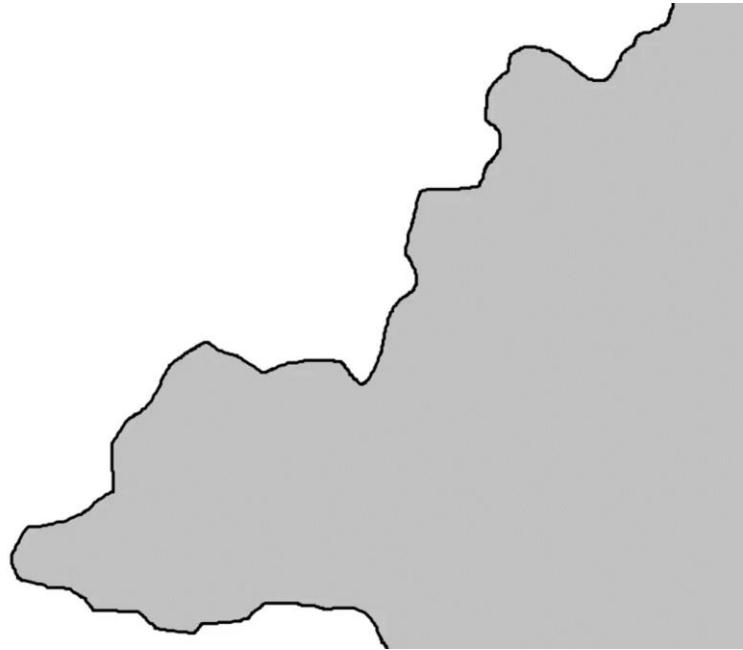


Partial Point Clouds



3D Scene

Read reasons for why shape completion is important from  
papers



+ :: M1: Basic knowledge learning and paper reading (~4 weeks)

- Pytorch learning: finish [\[A1\]](#) [\[A2\]](#) [\[A3\]](#) Assignments of eecs498. (Only if you know nothing about Pytorch before)
- Read and understand the following papers:
  - Cluster-VO [\[code\]](#)[\[paper\]](#)[\[video\]](#)
  - + :: ◦ DeepSDF [\[code\]](#)[\[paper\]](#)
  - Nerf [\[paper\]](#)
  - DSP-SLAM [\[code\]](#)[\[paper\]](#)[\[video\]](#)
  - Weakly Supervised Learning of Rigid 3D Scene Flow [\[code\]](#)[\[paper\]](#)
- Get familiar with the following datasets we will use:
  1. KITTI
  2. Nuscenes
  3. Waymo

EECS 498.008 / 598.008  
Deep Learning for Computer Vision  
Winter 2022

## Course Description

Computer Vision has become ubiquitous in our society, with applications in search, image understanding, apps, mapping, medicine, drones, and self-driving cars. Core to many of these applications are visual recognition tasks such as image classification and object detection. Recent developments in neural network approaches have greatly advanced the performance of these state-of-the-art visual recognition systems. This course is a deep dive into details of neural-network based deep learning methods for computer vision. During this course, students will learn to implement, train and debug their own neural networks and gain a detailed understanding of cutting-edge research in computer vision. We will cover learning algorithms, neural network architectures, and practical engineering tricks for training and fine-tuning networks for visual recognition tasks.

## + :: M2: Code Reading and demo test (~2 weeks)

- Read the code of DSP-SLAM and test it in different datasets, especially in dynamic scenes, and observe failure situations.
- Read the instance association part of Cluster-VO.

## M3: Build the main program framework (~4 weeks)

- Finish the data preprocessing part, including detection network and lidar odometry.
- Finish 2 in Methodology based on DSP-SLAM.
- Finish 3 in Methodology with the simplest method (bounding box overlap).

## M4: Make some contribution (~6 weeks)

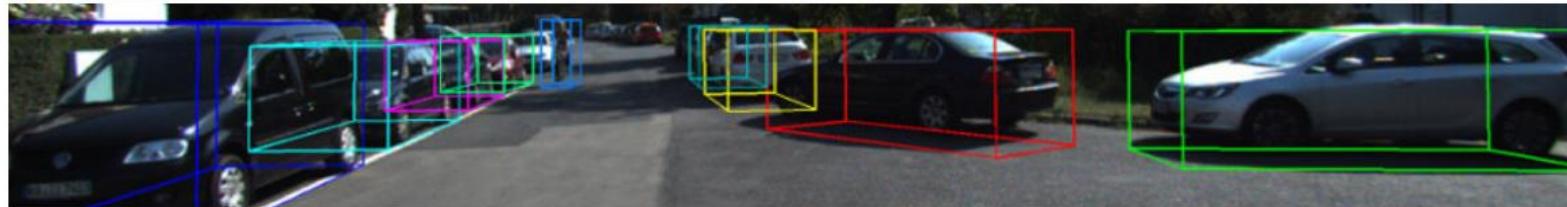
- Finish 4 in Methodology by designing a proper loss function.
- Explore more robust methods to achieve instance association.
- Explore the complete differentiable process from instance association to motion estimation.

## M5: Final Result: Evaluation (~4 weeks)

- Use scene flow dataset and metric to evaluate motion estimation.
- Design another metric to evaluate instance association with panoptic segmentation dataset.

# The KITTI Vision Benchmark Suite

## 3D Object Detection Evaluation 2017



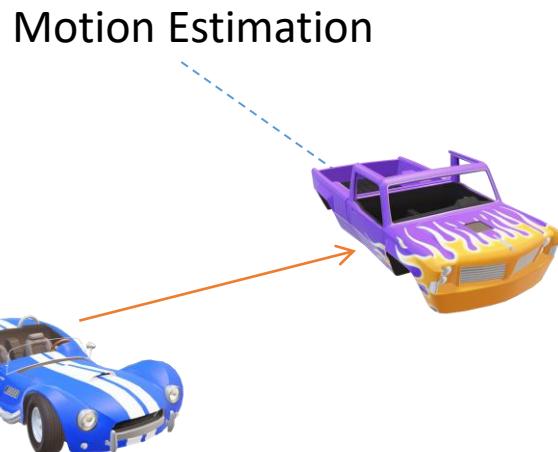
The 3D object detection benchmark consists of 7481 training images and 7518 test images as well as the corresponding point clouds, comprising a total of 80.256 labeled objects. For evaluation, we compute precision-recall curves. To rank the methods we compute average precision. We require that all methods use the same parameter set for all test pairs. Our development kit provides details about the data format as well as MATLAB / C++ utility functions for reading and writing the label files.

- [Download left color images of object data set \(12 GB\)](#)
- [Download right color images, if you want to use stereo information \(12 GB\)](#)
- [Download the 3 temporally preceding frames \(left color\) \(36 GB\)](#)
- [Download the 3 temporally preceding frames \(right color\) \(36 GB\)](#)
- [Download Velodyne point clouds, if you want to use laser information \(29 GB\)](#)
- [Download camera calibration matrices of object data set \(16 MB\)](#)
- [Download training labels of object data set \(5 MB\)](#)
- [Download object development kit \(1 MB\) \(including 3D object detection and bird's eye view evaluation code\)](#)
- [Download pre-trained LSVM baseline models \(5 MB\) used in Joint 3D Estimation of Objects and Scene Layout \(NIPS 2011\). These models are referred to as LSVM-MDPM-sv \(supervised version\) and LSVM-MDPM-us \(unsupervised version\) in the tables below.](#)
- [Download reference detections \(L-SVM\) for training and test set \(800 MB\)](#)
- Qianli Liao (NYU) has put together [code to convert from KITTI to PASCAL VOC file format](#) (documentation included, requires Emacs).
- Karl Rosaen (U.Mich) has released [code to convert between KITTI, KITTI tracking, Pascal VOC, Udacity, CrowdAI and AUTTI formats](#).
- Jonas Heylen (TRACE vzw) has released [pixel accurate instance segmentations](#) for all 7481 training images.
- We thank [David Stutz](#) and [Bo Li](#) for developing the 3D object detection benchmark.

# MILESTONE 2 PRESENTATION



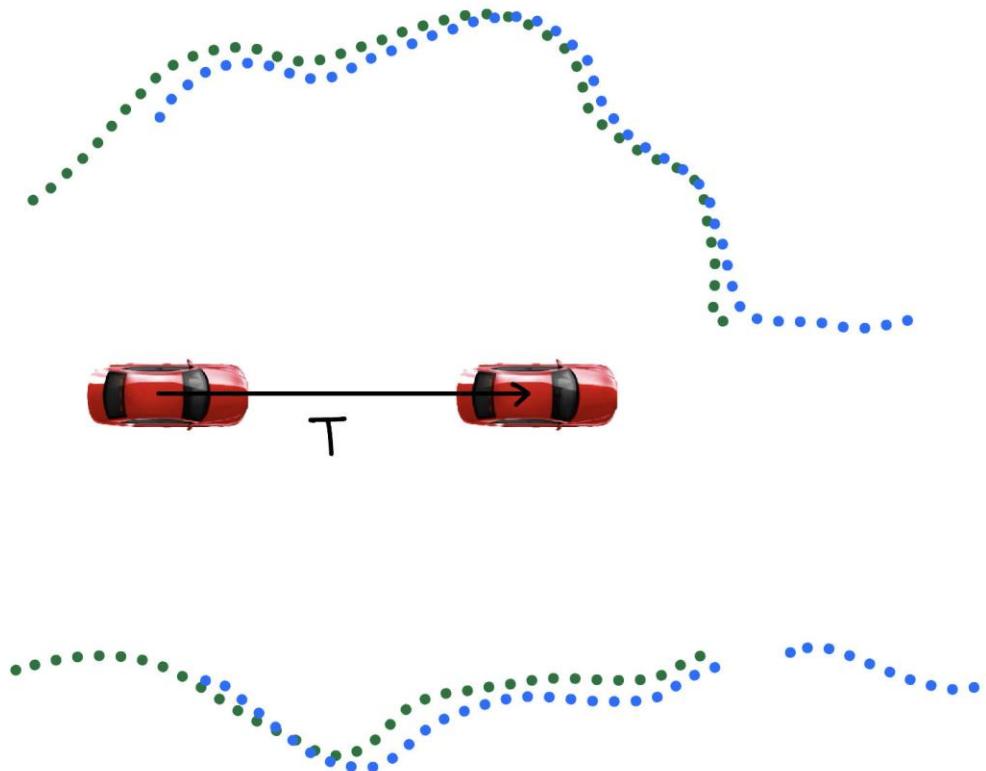
# Instance Completion and Motion Estimation with Deep Shape Priors for Autonomous Driving



**Team members**  
Panyawat (Ohm) Rattana  
Shashank (Sai) Dammalapati

**Supervisors**  
Xingguang (Starry) Zhong  
Yue Pan

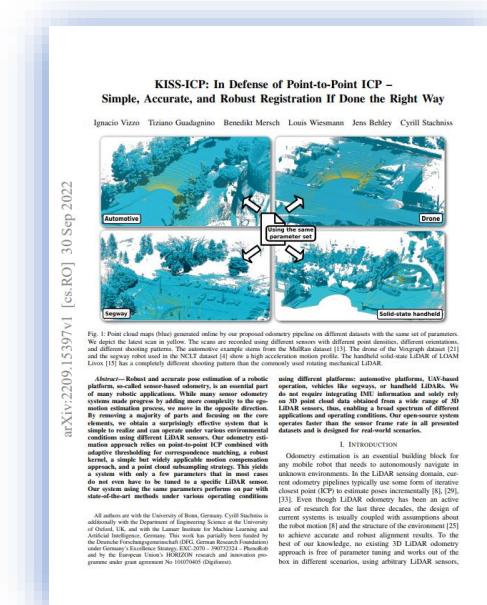
# Lidar Odometry (KISS ICP)



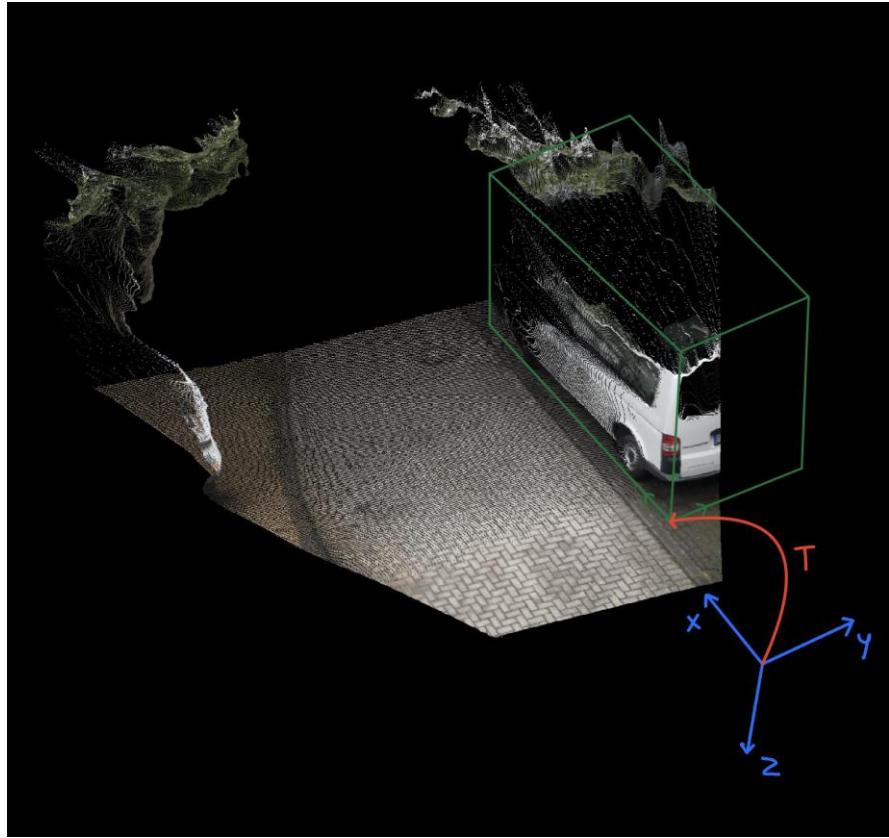
Motion Estimation of Ego Car

“Keep it small and simple”

- Accurately compute a robot’s pose.
- Point cloud alignment.
- A few parameters tuning.
- Assumption: The output of KISSICP is accurate.



# Point RCNN



Pose Estimation of surrounding cars in 1 frame

- Idea is to generate bounding boxes for cars
- Point RCNN gives us an initial guess for the relative pose of the external cars relative to ego car.

**PointRCNN: 3D Object Proposal Generation and Detection from Point Cloud**

Shaochai Shi Xiaogang Wang Hongsheng Li  
The Chinese University of Hong Kong  
{ssshi, xgwang, hsl1}@ee.cuhk.edu.hk

**Abstract**

In this paper, we propose PointRCNN for 3D object detection from raw point cloud. The whole framework is composed of two stages: stage-1 for the bottom-up 3D proposal generation and stage-2 for refining proposals in the top-down 3D object detection and classification results. Instead of generating proposals from RGB image or projecting point cloud to bird's view or voxels as previous methods do, our stage-1 sub-network directly generates a set of coarse and highly overlapping 3D local point cloud in a bottom-up manner via segmenting the point cloud of the whole scene into foreground points. The local point cloud is then processed to extract the canonical points of each proposal to canonical coordinates to learn better local spatial features, which is combined with global semantic features of each point learned in stage-1 for accurate classification. In stage-2, the refined 3D bounding boxes are generated by a 3D multi-scale feature fusion network. Extensive experiments on the 3D detection benchmark of KITTI dataset show that our proposed architecture outperforms state-of-the-art methods with large margins by using only point clouds as input. The code is available at <https://github.com/shaochai/PointRCNN>.

**1. Introduction**

Deep learning has achieved remarkable progress on 2D computer vision tasks, including object detection [8, 32, 16] and instance segmentation [6, 10, 20], etc. Beyond 2D vision, 3D perception is also increasingly indispensable for many real-world applications, while recent developed 2D detection algorithms are capable of handling large variations in the scene. However, the 3D detection of 3D objects with point clouds still faces great challenges from the irregular data format and large search space of 6 Degrees of Freedom (DoF) of 3D objects.

In autonomous driving, the sensor commonly used 3D sensors are the LiDAR sensors, which generate 3D point clouds to capture 3D structures of the scene. The difficulty of point cloud-based 3D object detection mainly lies in irregularity of the point clouds. State-of-the-art 3D detection methods either leverage the mature 2D detection frameworks by projecting the point clouds into bird's view [14, 42, 17] (see Fig. 1(a)), to the frontal view [4, 38], or to the regular 3D voxels [34, 43], which are not optimal and suffer from low resolution during the projection process.

Instead of transforming point clouds to voxels or other regular data structures for feature learning, Qi *et al.* [26, 28] proposed PointNet for learning 3D representations directly from point clouds for 3D classification and segmentation. As shown in Fig. 1(b), their follow-up work [25] applied PointNet in 3D object detection to estimate the 3D bounding boxes based on the cropped frustum point cloud from the LiDAR sensor. However, the performance of the method heavily relies on the 2D detection performance and cannot take the advantages of 3D information for generating robust bounding box proposals.

Unlike object detection from 2D images, 3D objects in autonomous driving series are naturally and well separated

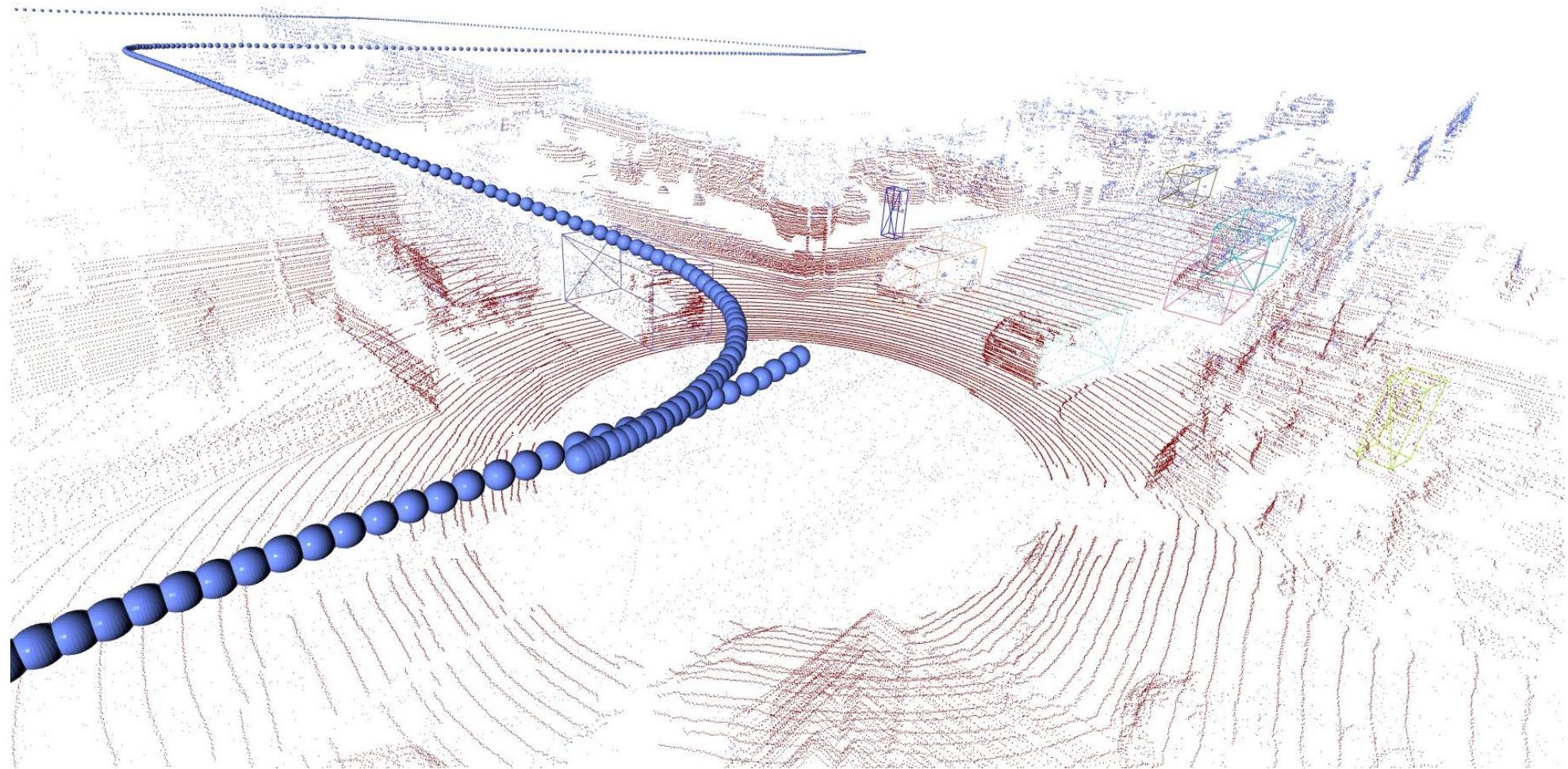
arXiv:1812.04244v2 [cs.CV] 16 May 2019

**a. Aggregate View Object Detection (AVOD)**

**b. Frustum-Pooling**

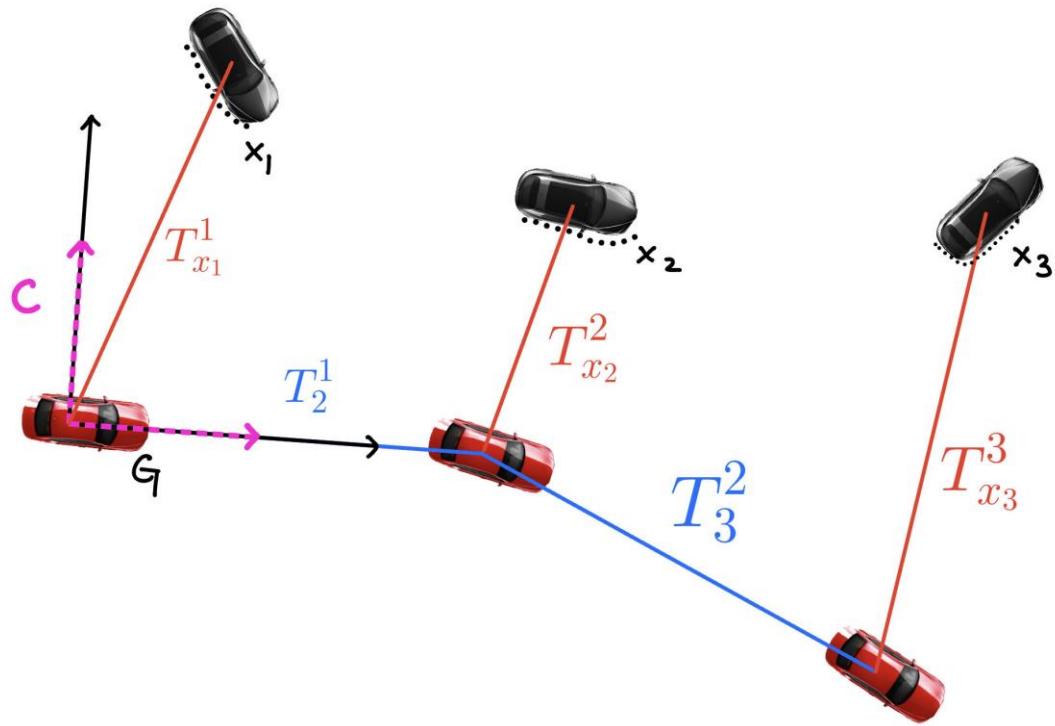
**c. Our approach (PointRCNN)**

# Lidar Odometry and 3D Bounding Box Detection



KISS ICP + Point RCNN

# Workflow: Transformation



C – Canonical Frame

G – Global Frame

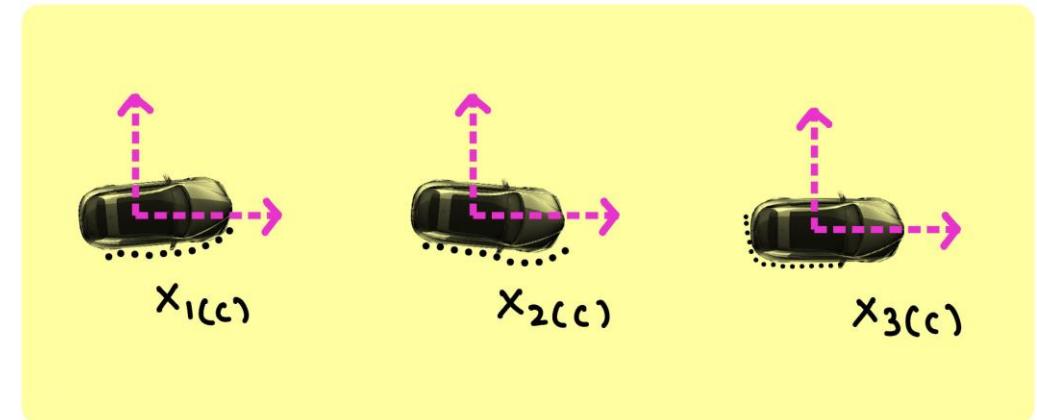
$T_{x_N}^N$  Transformation from Point RCNN

$T_N^{N-1}$  Transformation from KISS ICP

$$X_1[C] = T_{x_1}^1 X_1$$

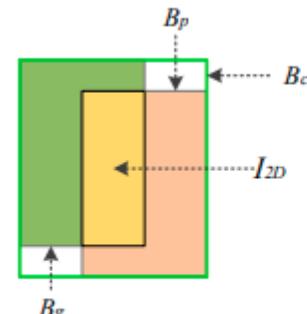
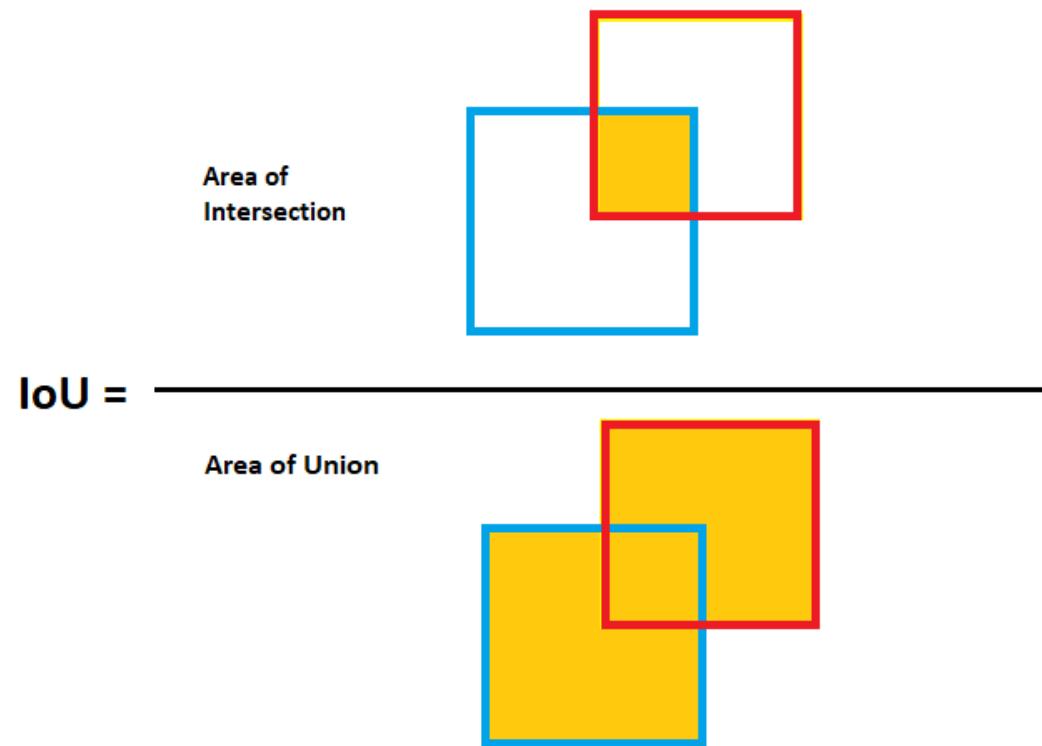
$$X_2[C] = T_2^1 T_{x_2}^2 X_2$$

$$X_3[C] = T_2^1 T_3^2 T_{x_3}^3 X_3$$

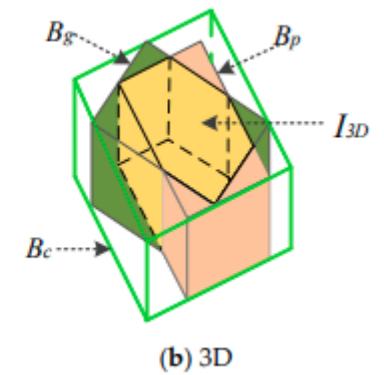


Point Clouds in Canonical Space

# IOU: Definition



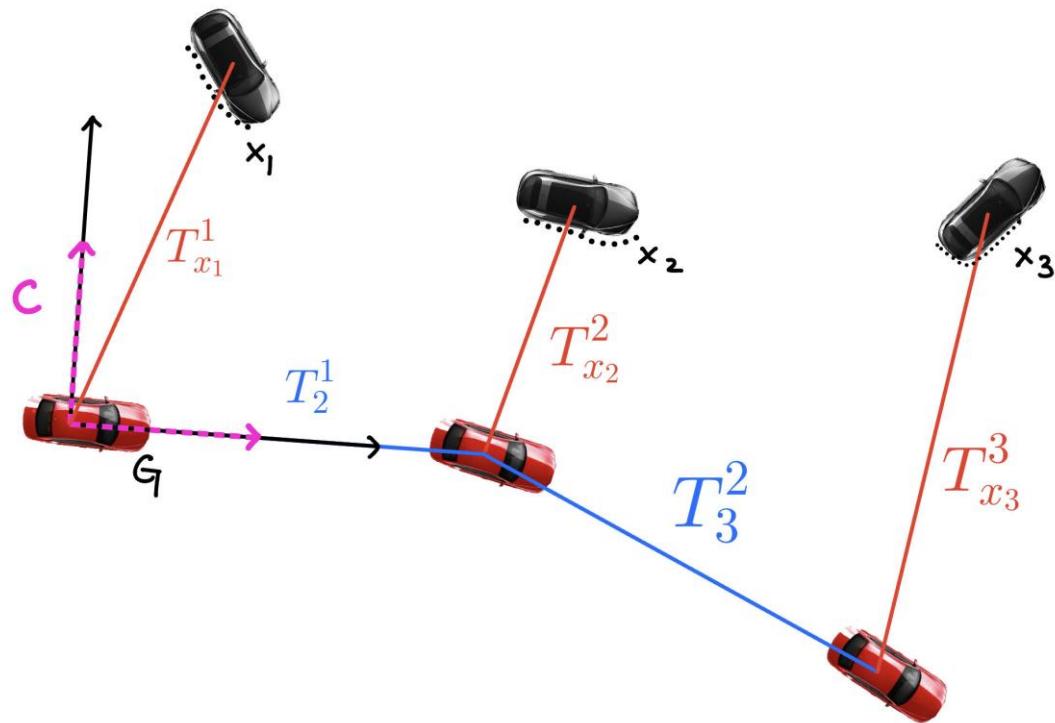
(a) 2D



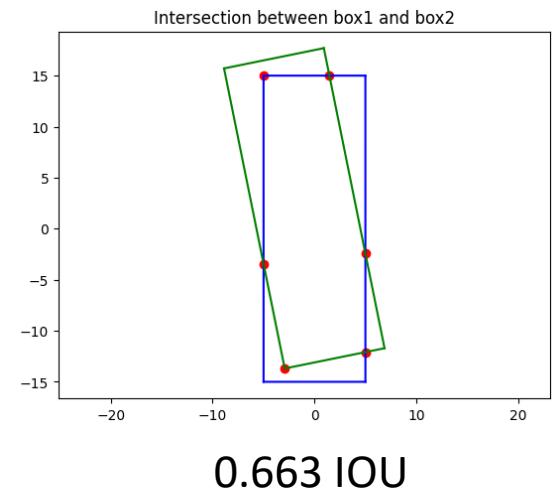
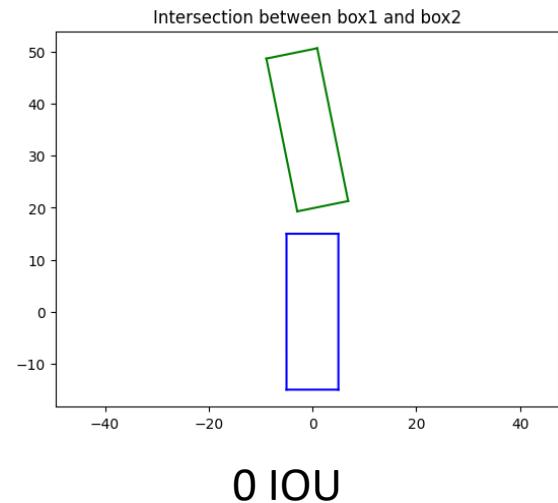
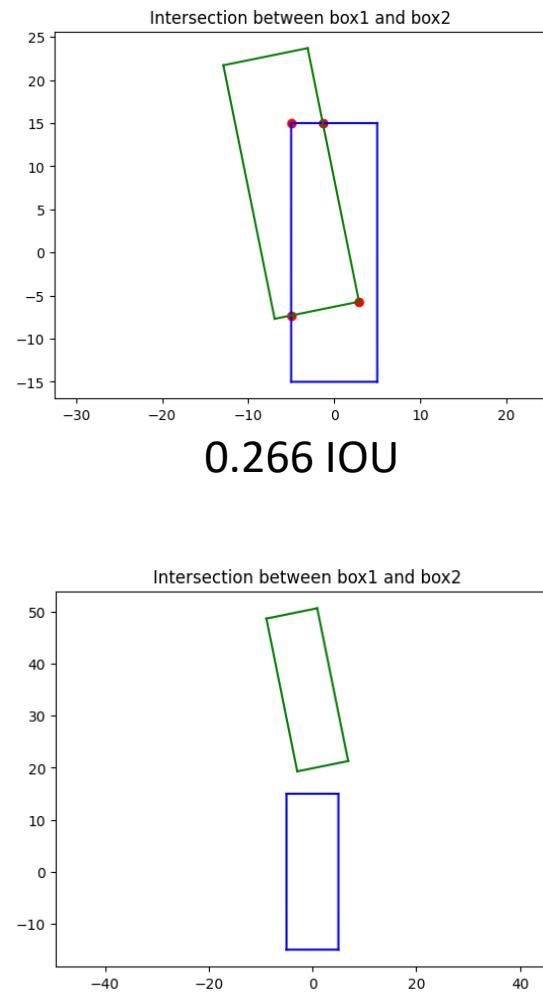
# Video



# Tracking Multiple Objects

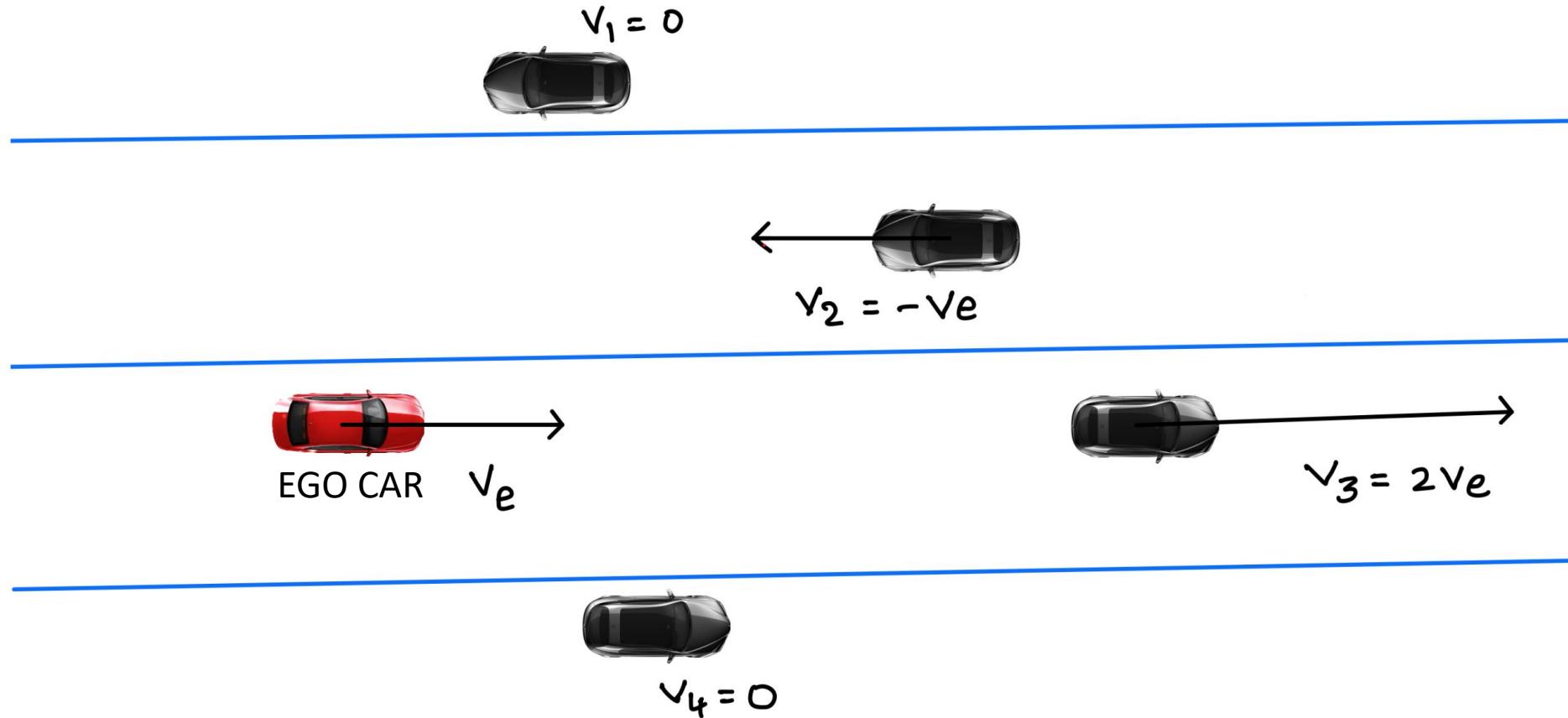


Example of tracking one car using lidar  
odometry and bounding box measurements



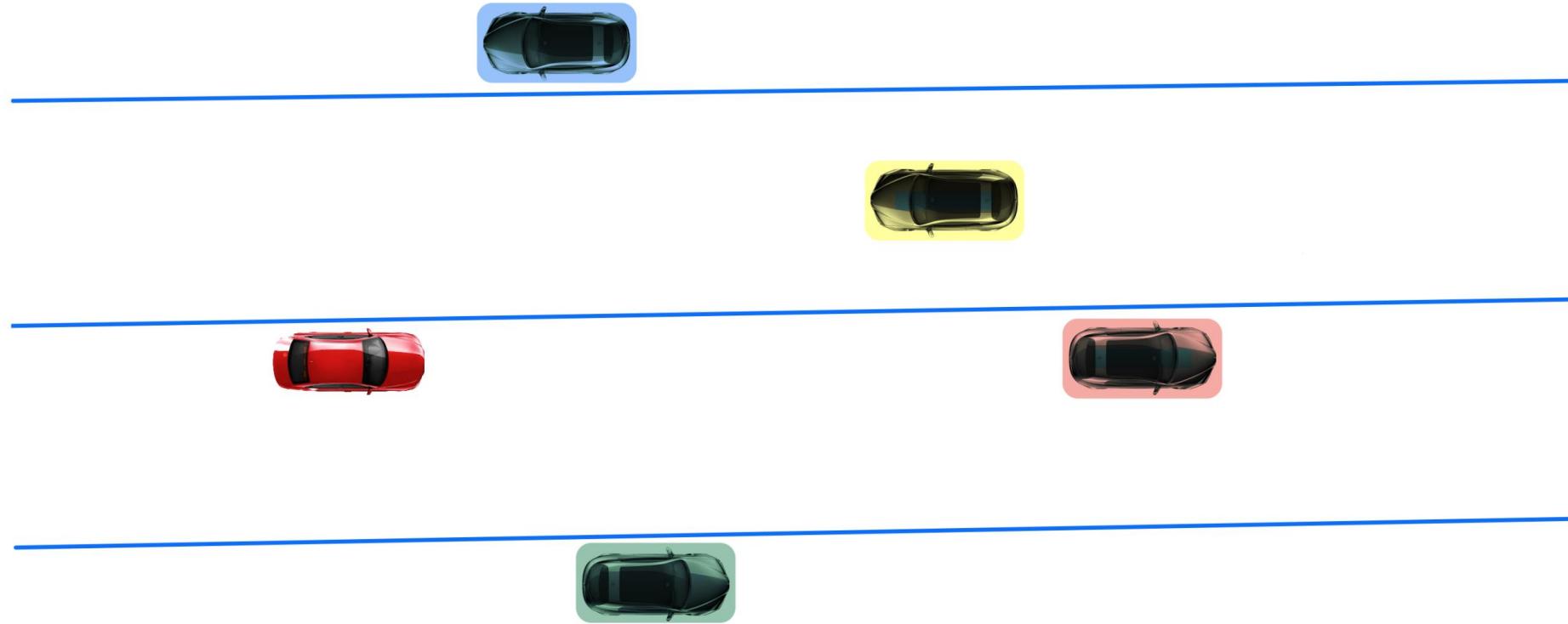
# Tracking Multiple Objects

Global frame



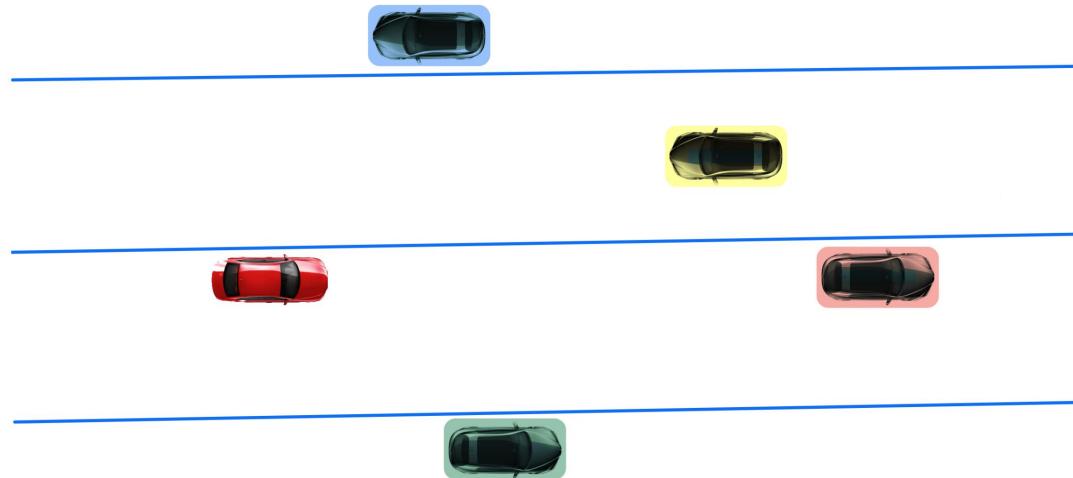
# Prediction of Bounding Boxes

Global frame

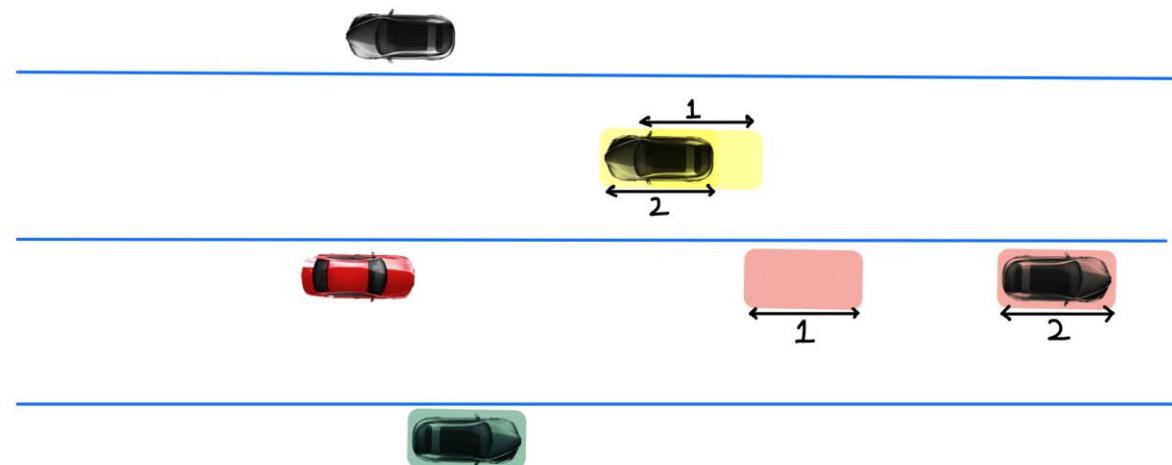


# Basic IOU problems

Global frame



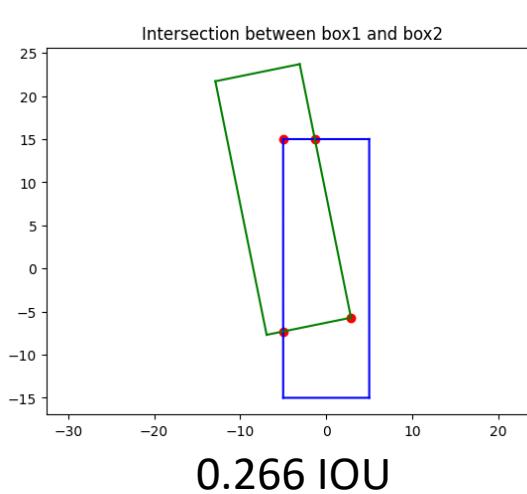
Global frame



1. Sometimes PointRCNN detection fails to detect cars
2. Cars moving with high speeds may have 0 IOU, failing instance tracking

How can we improve instance tracking?

# How can we improve instance tracking?



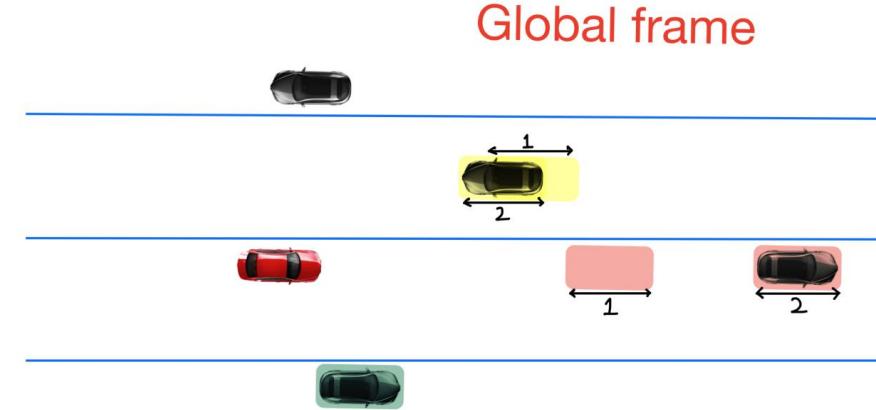
Frame	Cars in Scene	Cars Detected	Buffer
1	1 2 3	1 2 3	1 2 3
2	1 2 3 4	1 3 4	1 3 4
3	2 3 4	2 3 4	2 3 4

If (Buffer = Detections)



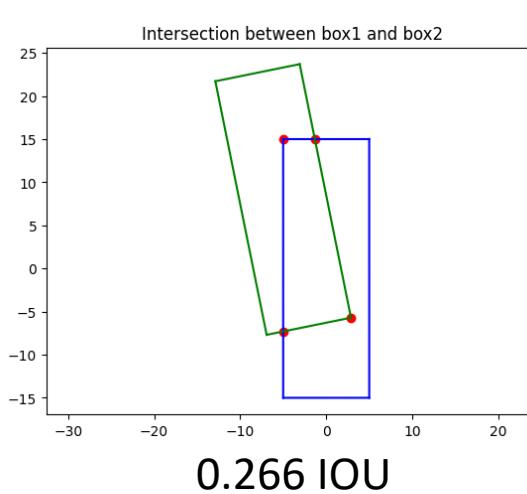
PyTorch3D IoU

Improving IoU calculation



More robust instance association using motion models

# How can we improve instance tracking?

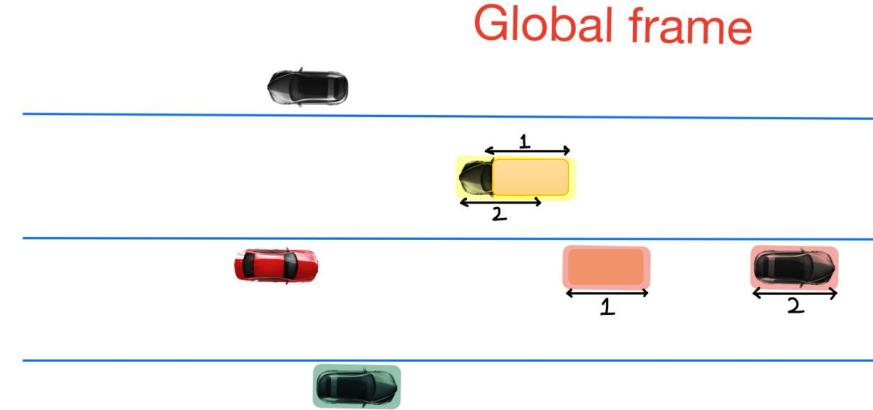


Frame	Cars in Scene	Cars Detected	Buffer
1	1 2 3	1 2 3	1 2 3
2	1 2 3 4	1 3 4	1 3 4
3	2 3 4	2 3 4	2 3 4

If (Buffer = Detections)

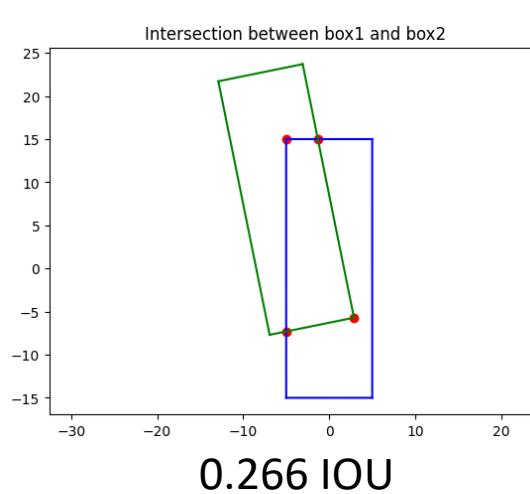


Improving IoU calculation



More robust instance association track cars' velocities

# How can we improve instance tracking?

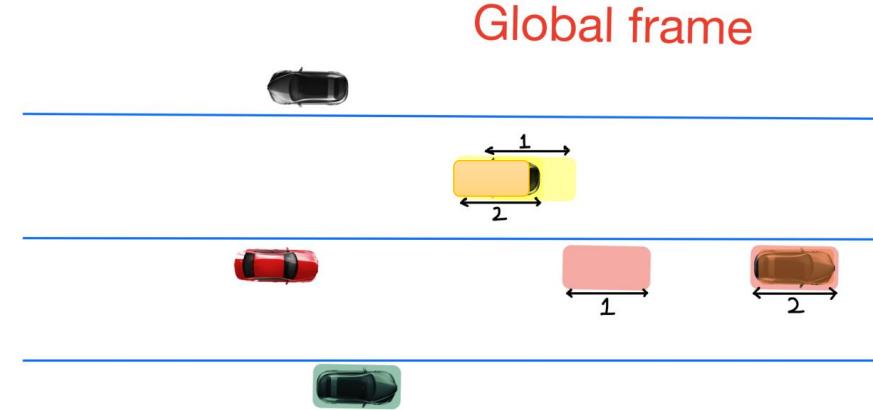


Frame	Cars in Scene	Cars Detected	Buffer
1	1 2 3	1 2 3	1 2 3
2	1 2 3 4	1 3 4	1 3 4
3	2 3 4	2 3 4	2 3 4

If (Buffer = Detections)

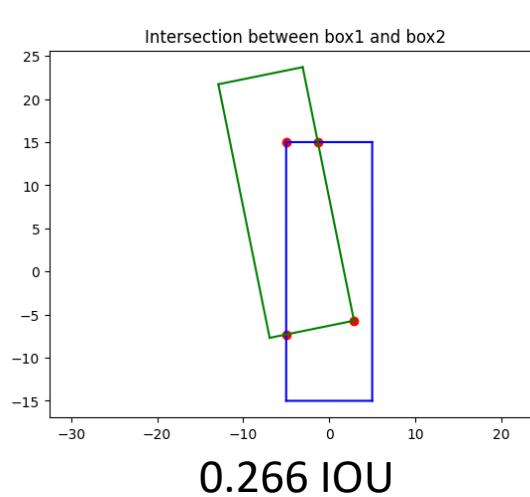


Improving IoU calculation



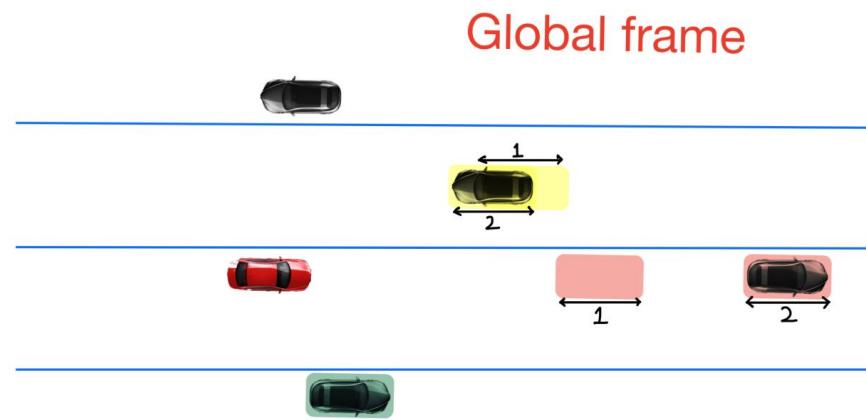
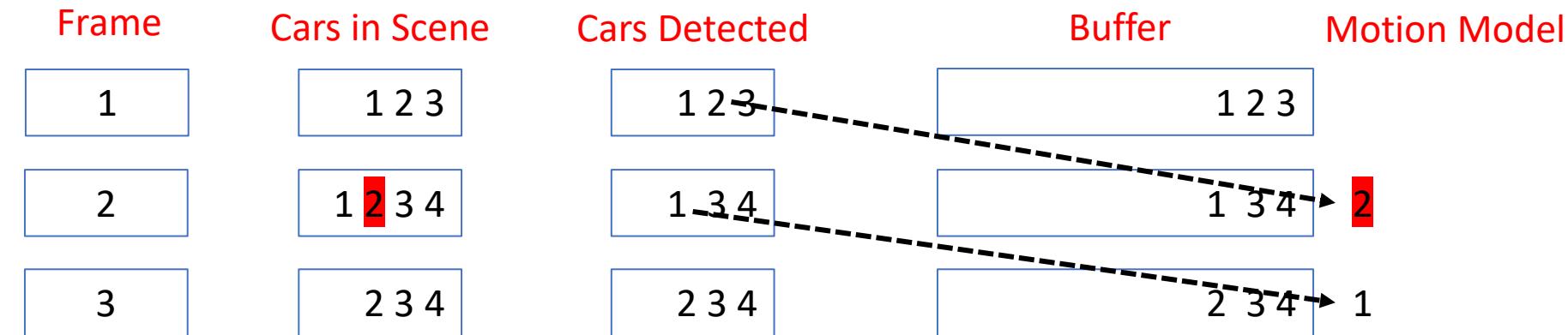
More robust instance association track cars' velocities

# How can we improve instance tracking?



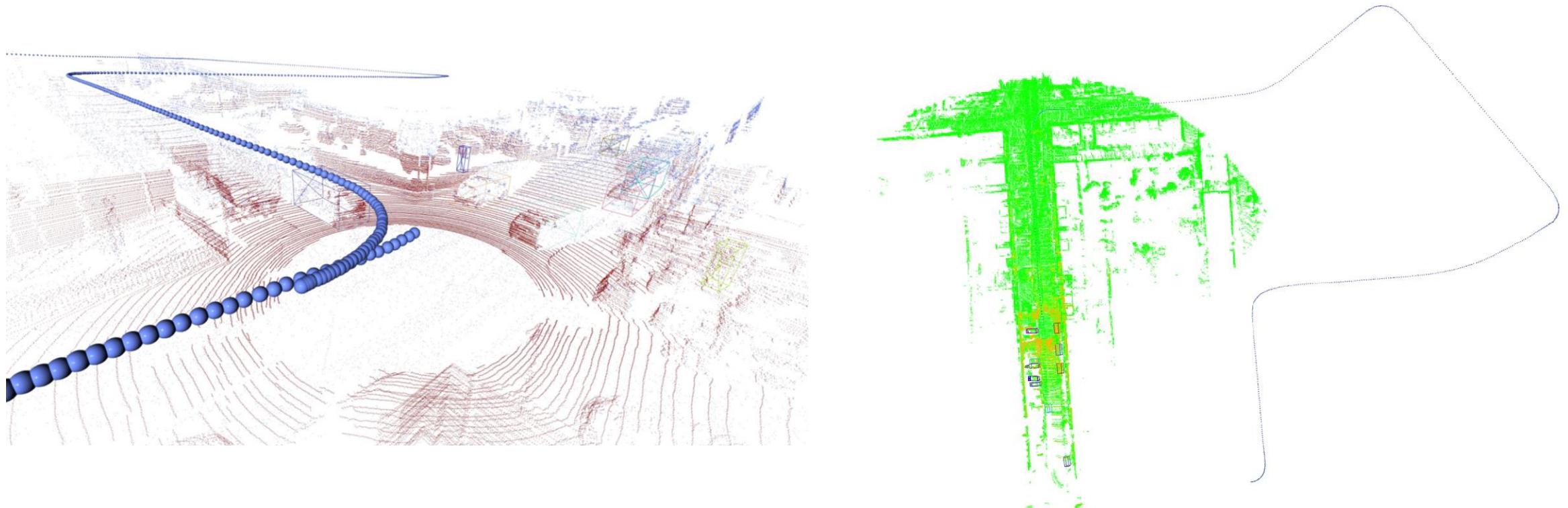
PyTorch3D IoU

Improving IoU calculation

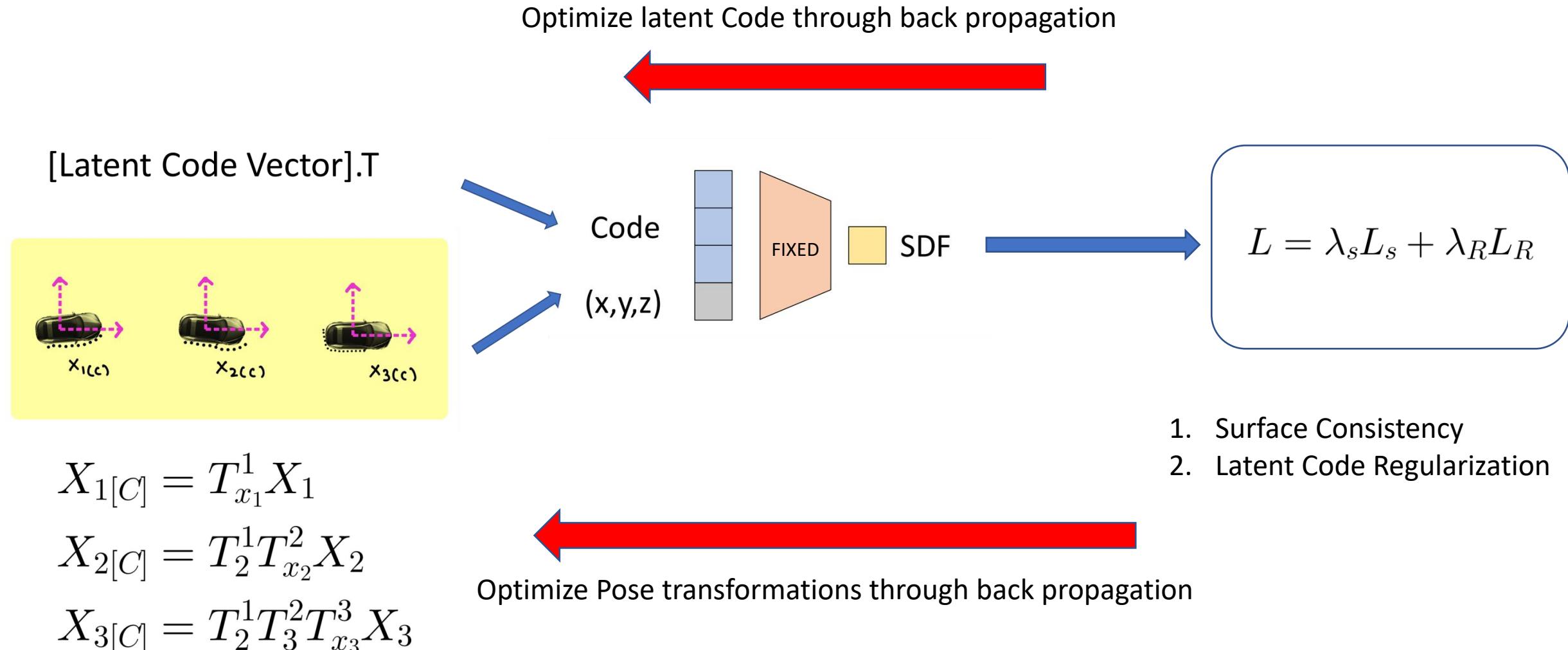


More robust instance association track cars' velocities

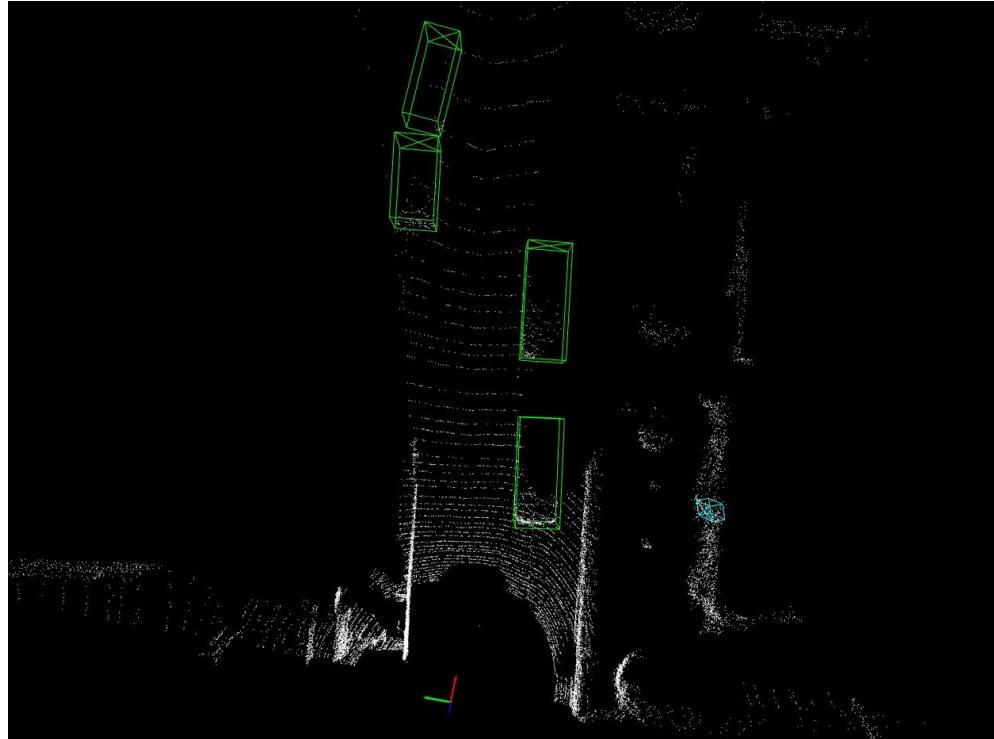
# Tracking in Global and Sensor Frame



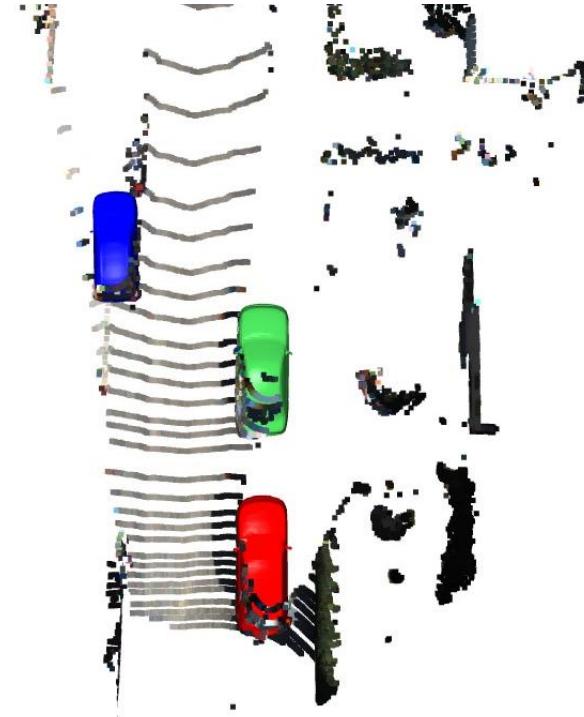
After proper instance association, we can use DeepSDF to jointly optimize shape and pose of the detected moving cars



# Optimizing Latent Code and Pose Estimation

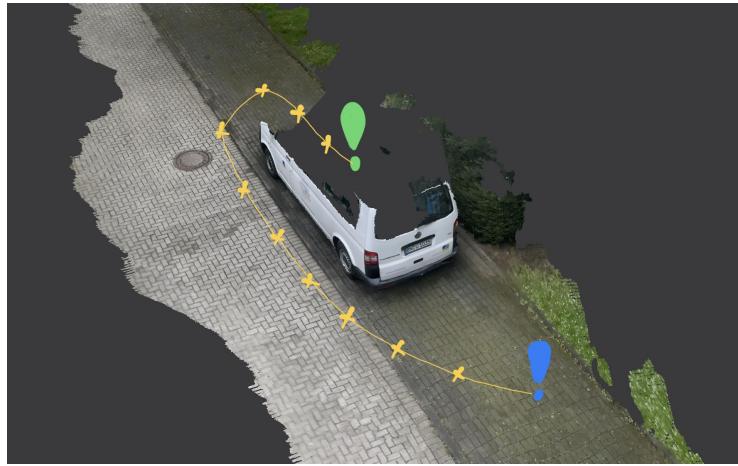


Lidar + Bounding Box

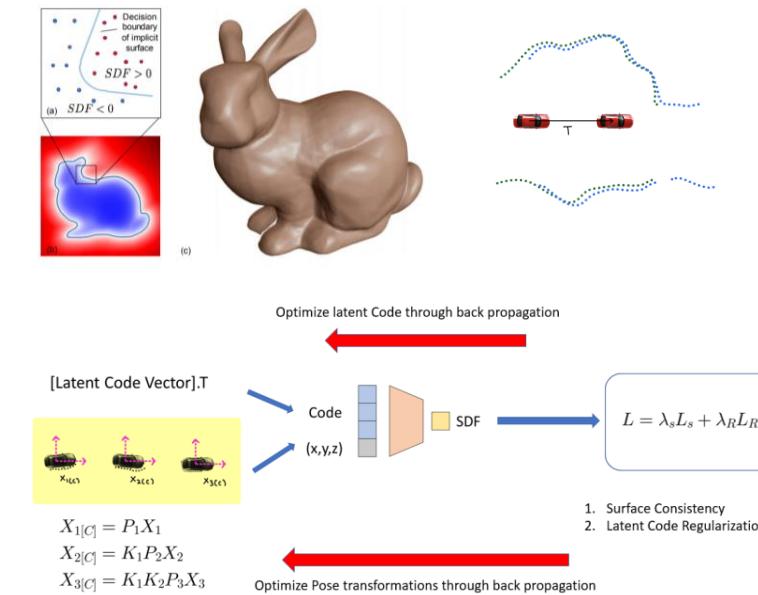


Stereo Image + Shape Completion

# Overview of the project



## Challenges

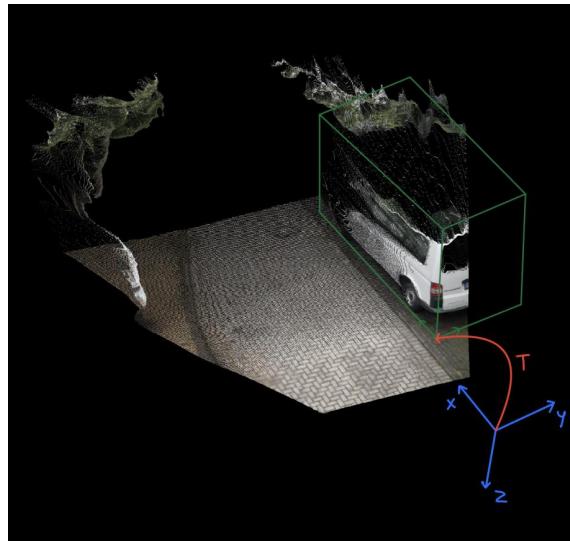


## Method

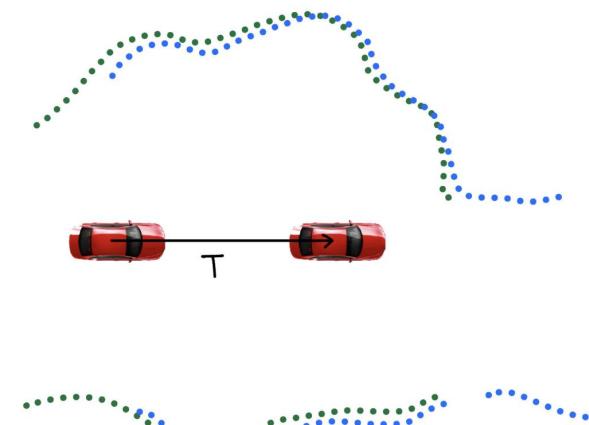


## Final Output

# Completed Tasks



Bounding Box Prediction

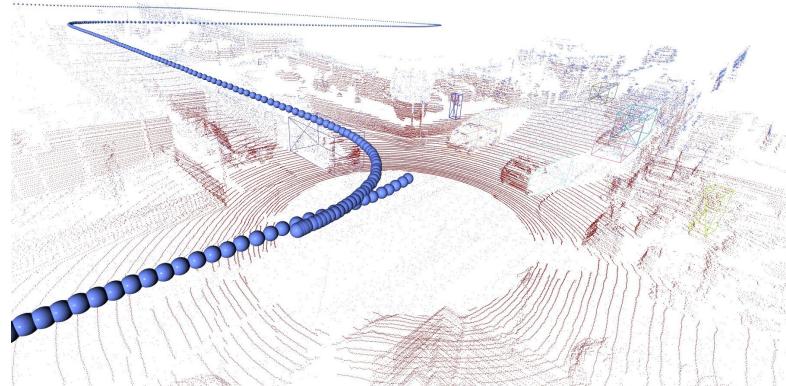


Lidar Odometry

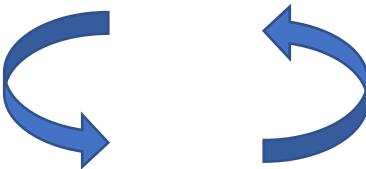


Instance Association and  
Multi Object Tracking

# Next Steps

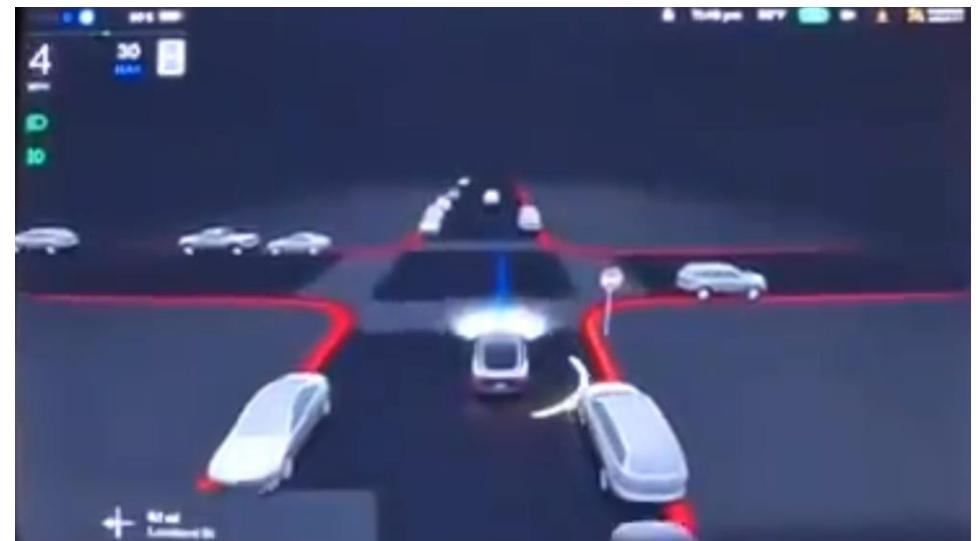
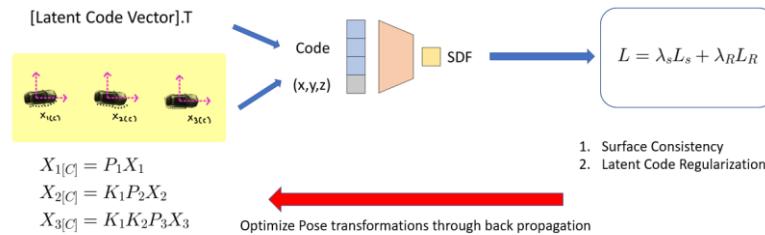


Pose and  
Point  
Cloud of  
the cars



Optimized  
Pose and  
Shape Code of  
the cars

Marching Cubes  
+  
Visualization



# Thank You

# Other Works using Shape Completion

## Panoptic Mapping with Fruit Completion and Pose Estimation for Horticultural Robots

Yue Pan  
Claus Smitt

Federico Magistri  
Chris McCool

Thomas Läbe  
Jens Behley

Elias Marks  
Cyrill Stachniss

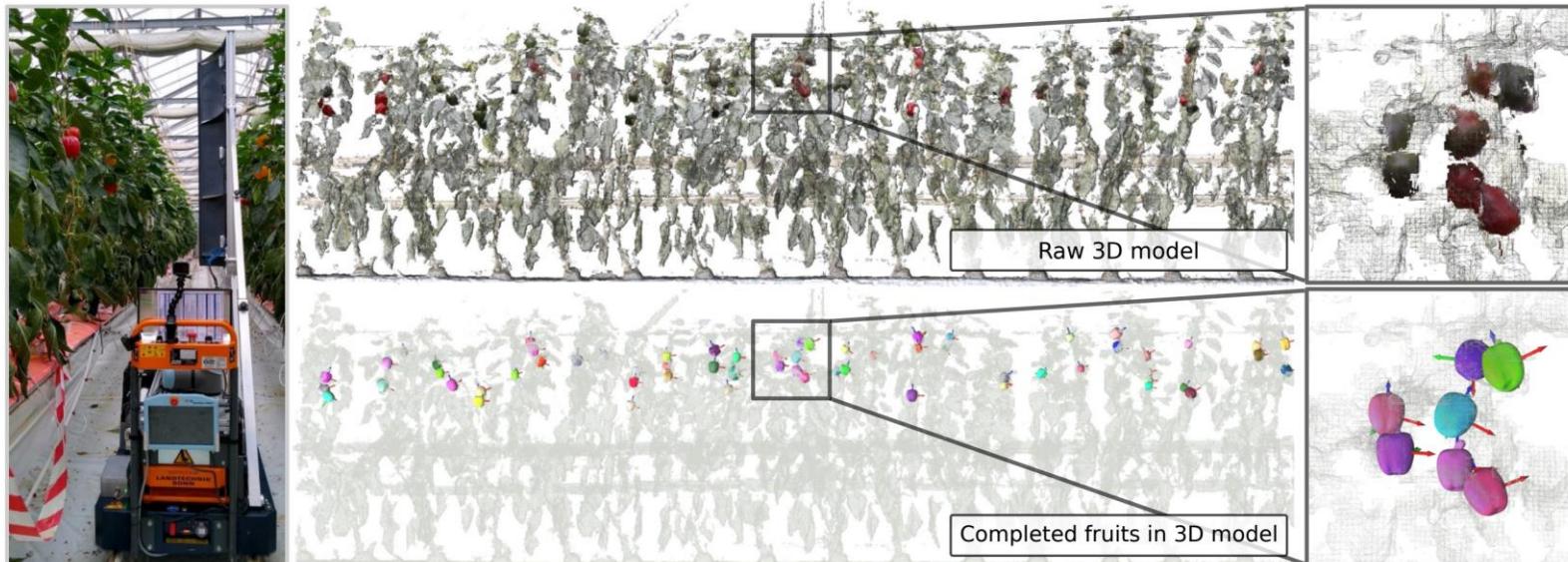
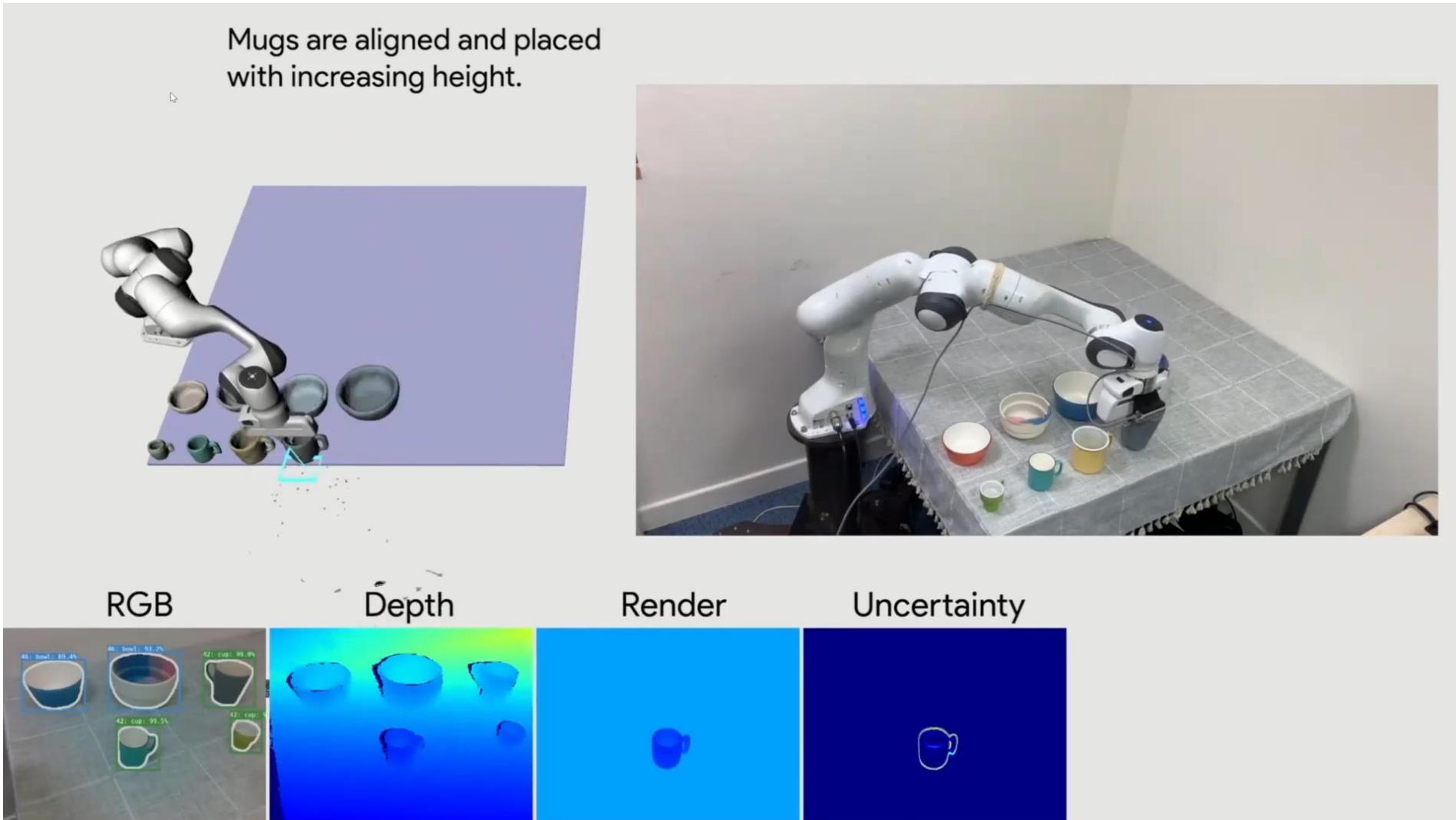


Fig. 1: Our method is able to build a multi-resolution panoptic map (top) of a challenging commercial glasshouse environment online using a mobile horticultural robot equipped with RGB-D cameras (left). Furthermore, our method manages to jointly estimate the complete shape and pose of each fruit in the map (bottom).

# Node SLAM

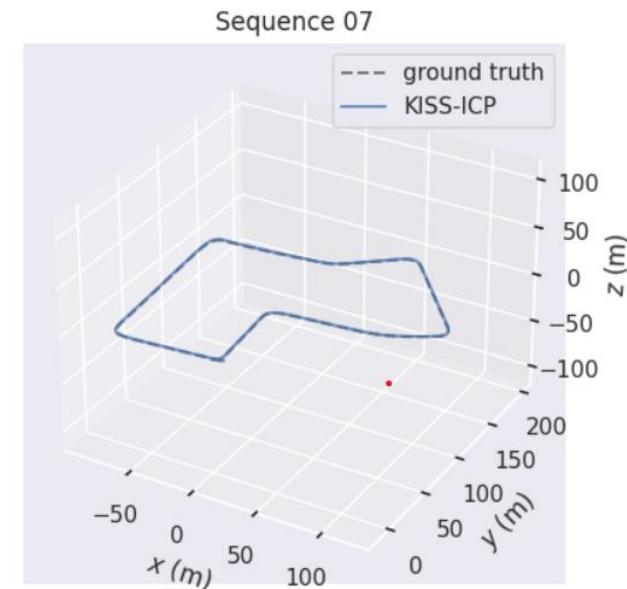


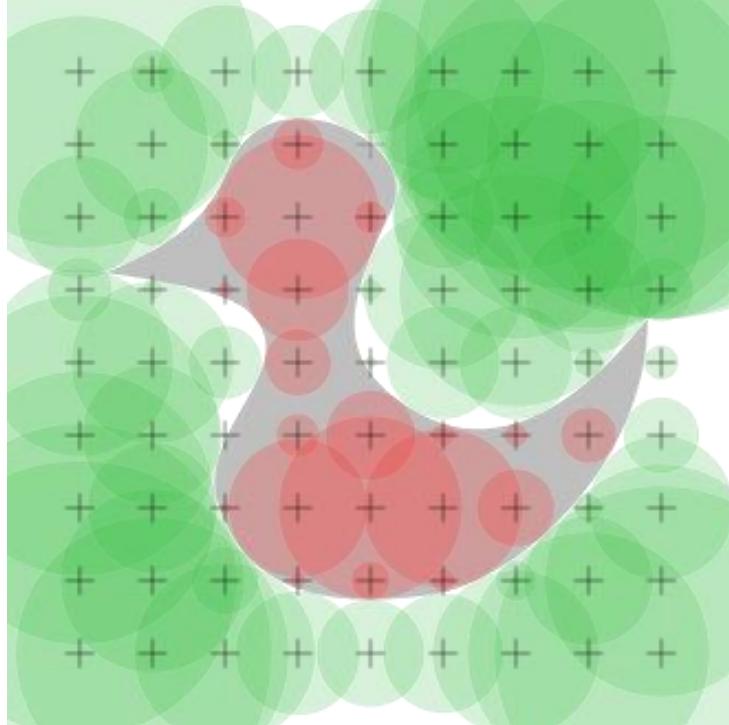
# Metrics

- PointRCNN - 8 fps
- Kiss ICP – 60 fps
- Shape Completion – 5 fps (avg)
- Kiss ICP + Bounding Box (Stored) – 30 fps
- DSP SLAM – 5 fps

Now evaluating sequence 07  
0% |

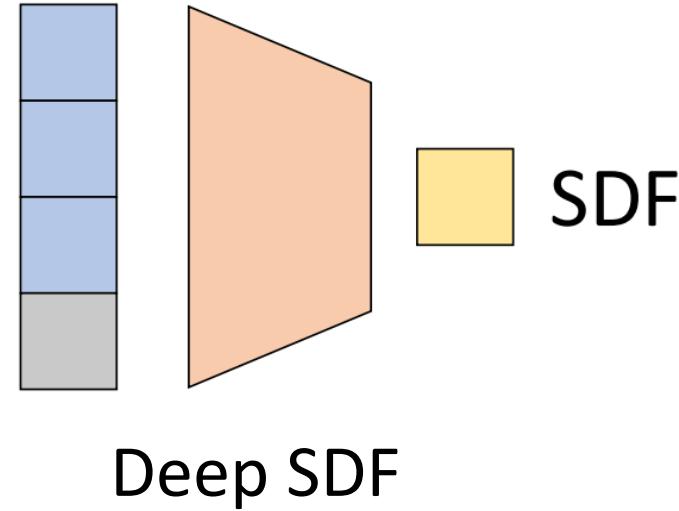
Metric	Value	Units
Average Translation Error	0.328	%
Average Rotational Error	0.164	deg/m
Absoulte Trajectory Error (ATE)	0.776	m
Absoulte Rotational Error (ARE)	0.007	rad
Average Frequency	69	Hz
Average Runtime	14	ms





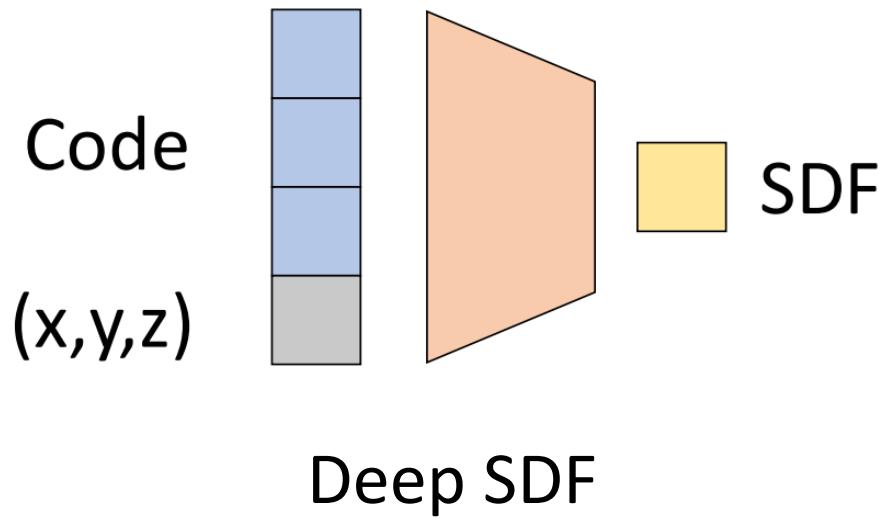
SDF  $f_{\theta}(x, y, z) \approx SDF(x, y, z)$

Code  
( $x, y, z$ )



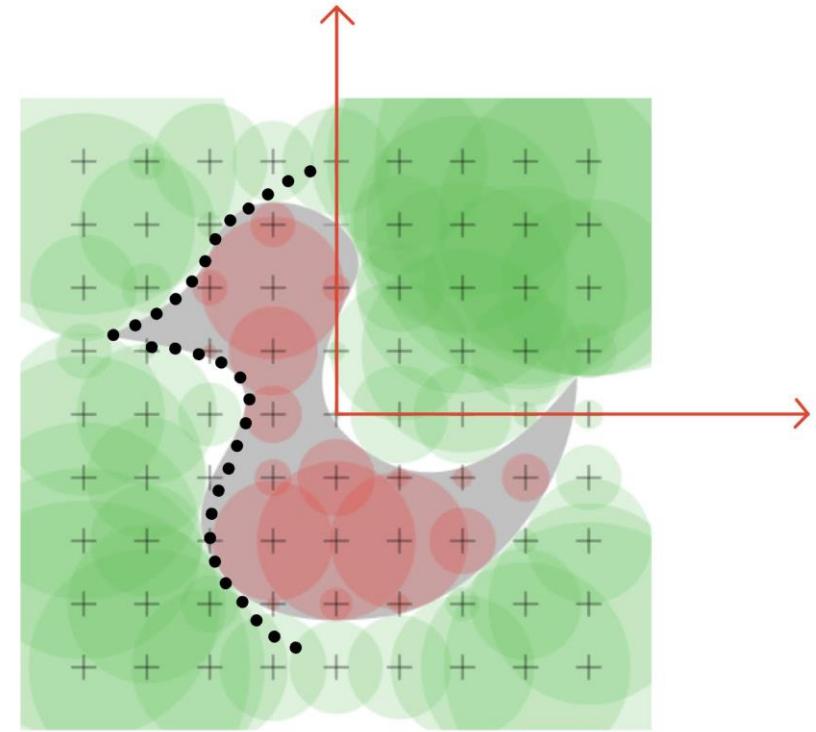
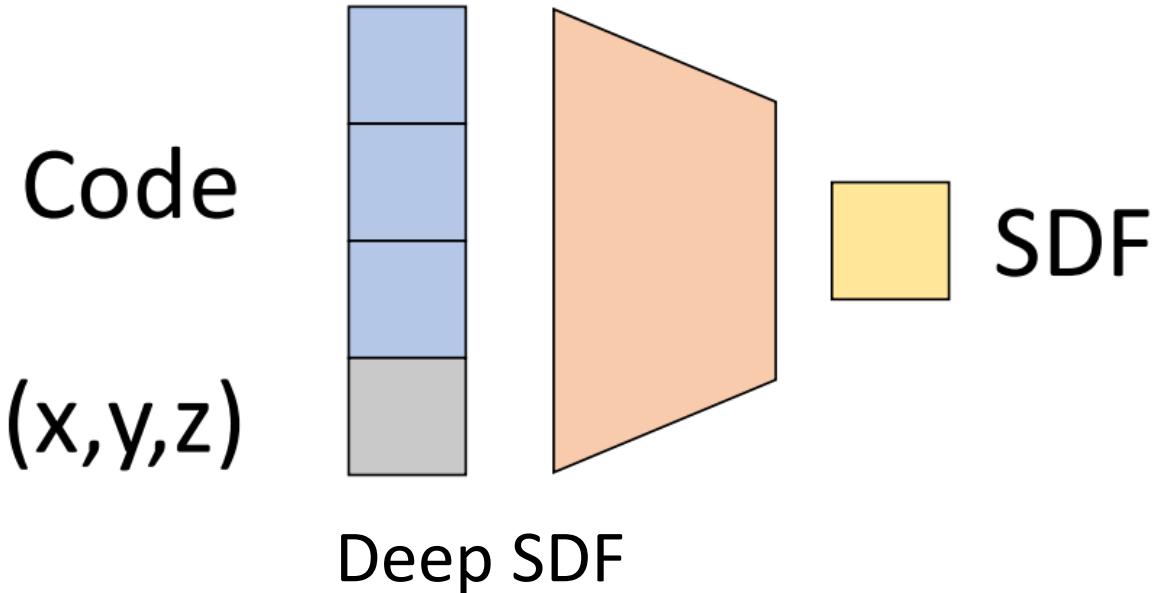
DeepSDF can implicitly model a class of objects

# Deep SDF



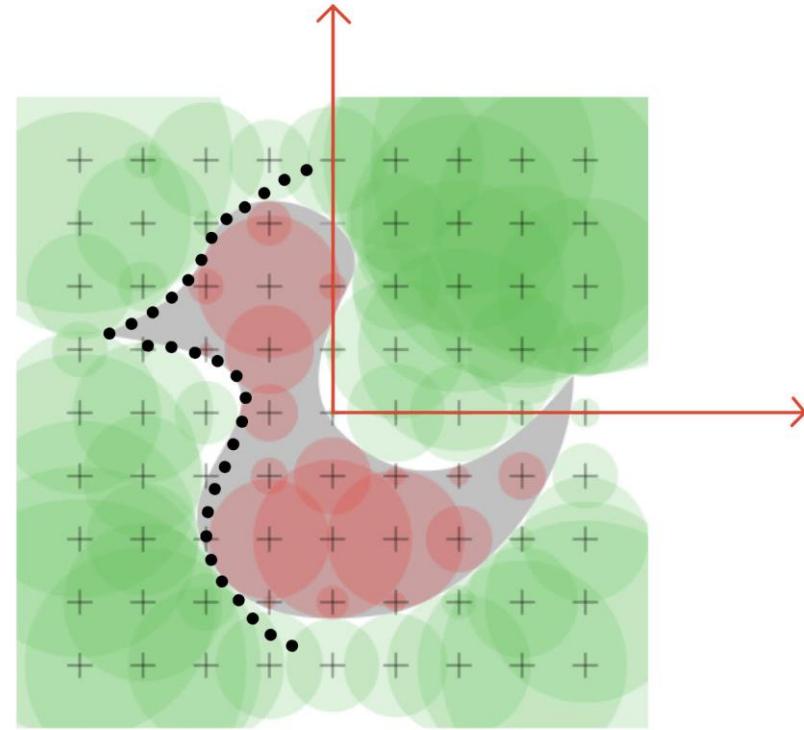
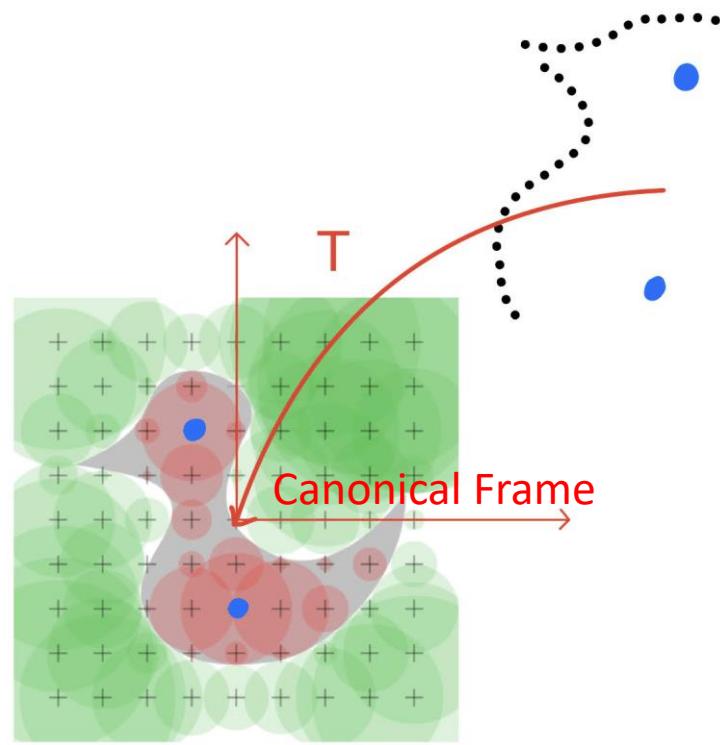
A latent code collapses a general representation into a single shape

# 0 SDF



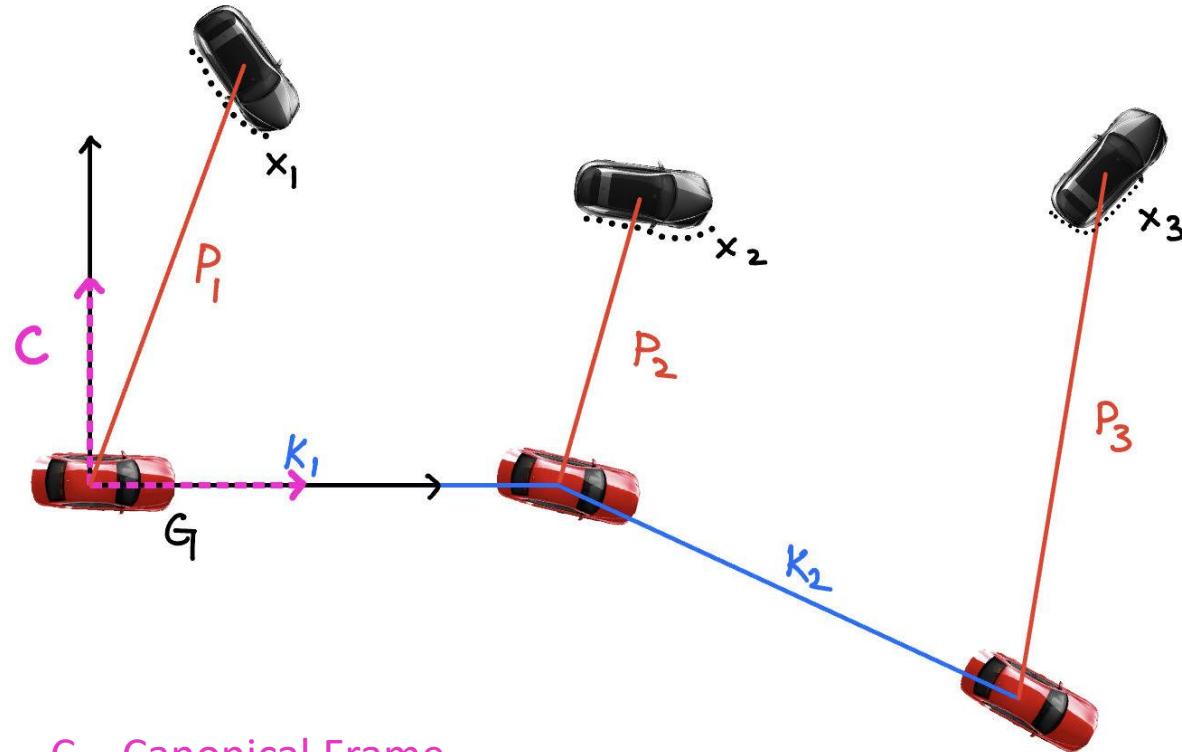
SDF value for a point on surface modelled by SDF is ZERO

# Canonical Frame



SDF presumes that the input is in Canonical Frame

# Workflow: Transformation



C – Canonical Frame

G – Global Frame

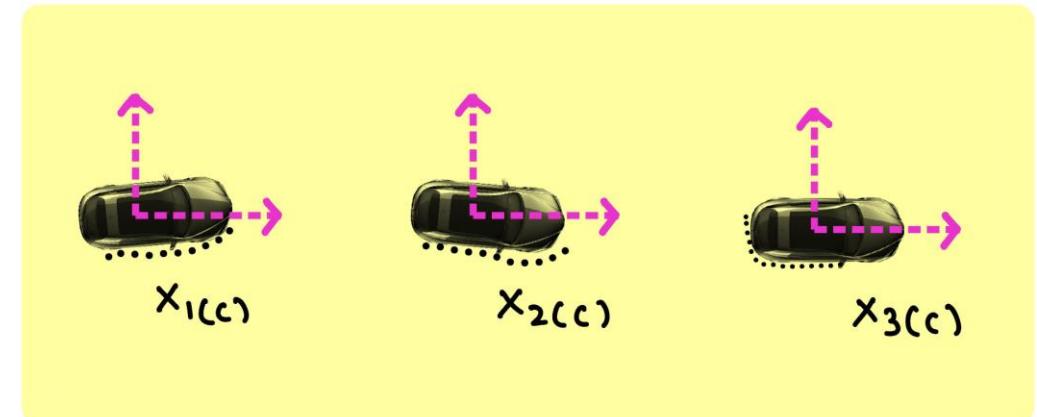
P – Transformation from Point RCNN

K – Transformation from KISS ICP

$$X_1[C] = P_1 X_1$$

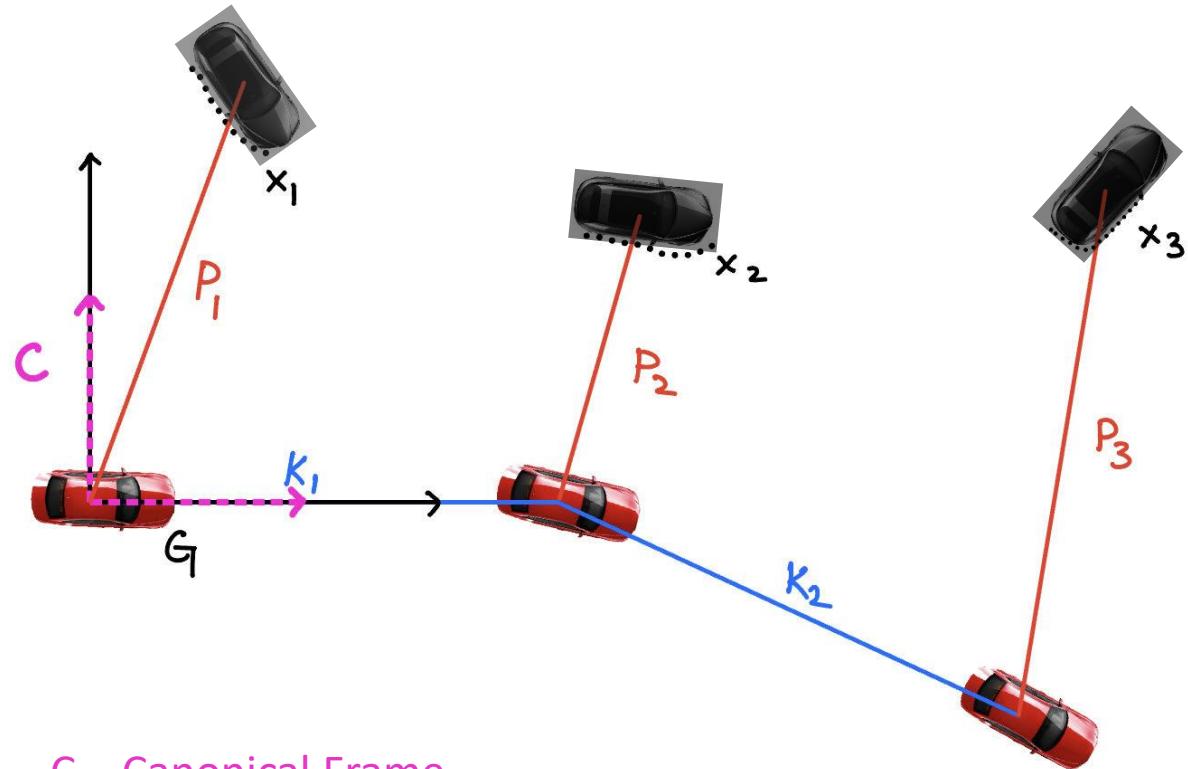
$$X_2[C] = K_1 P_2 X_2$$

$$X_3[C] = K_1 K_2 P_3 X_3$$



Canonical Space

# Workflow: Transformation



C – Canonical Frame

G – Global Frame

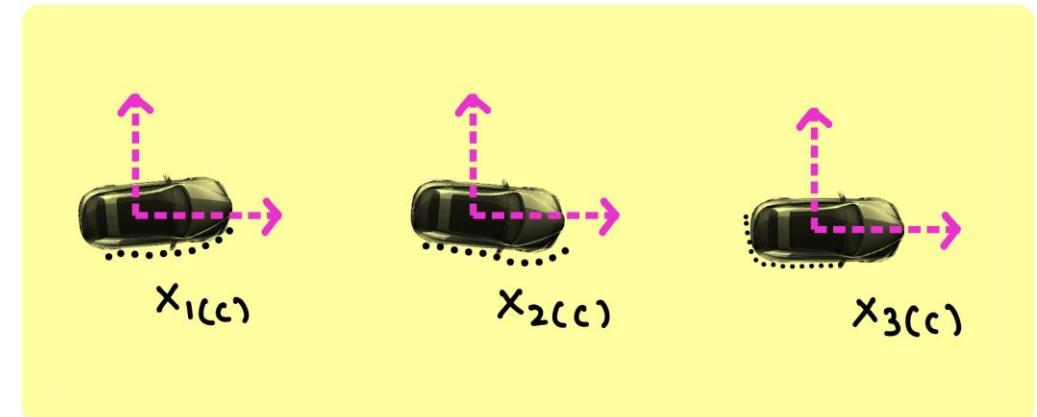
P – Transformation from Point RCNN

K – Transformation from KISS ICP

$$X_1[C] = P_1 X_1$$

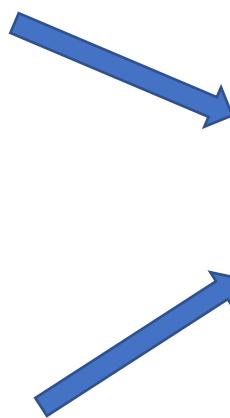
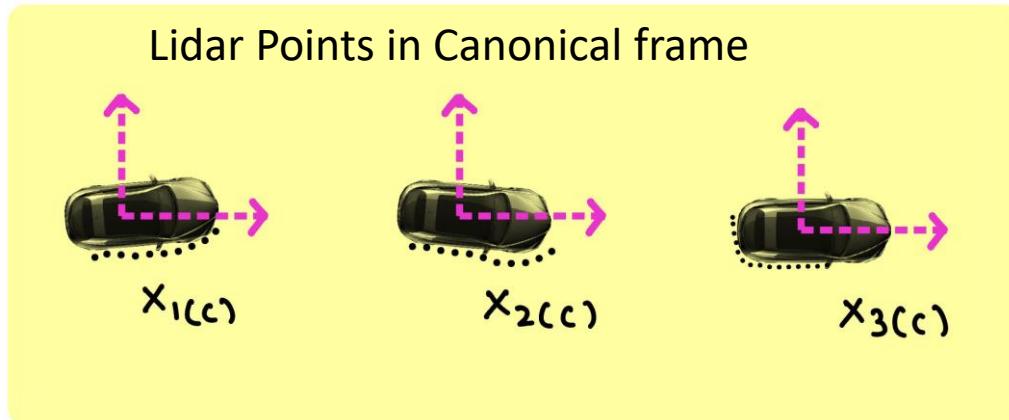
$$X_2[C] = K_1 P_2 X_2$$

$$X_3[C] = K_1 K_2 P_3 X_3$$

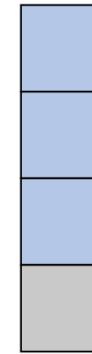


Canonical Space

[Latent Code Vector].T

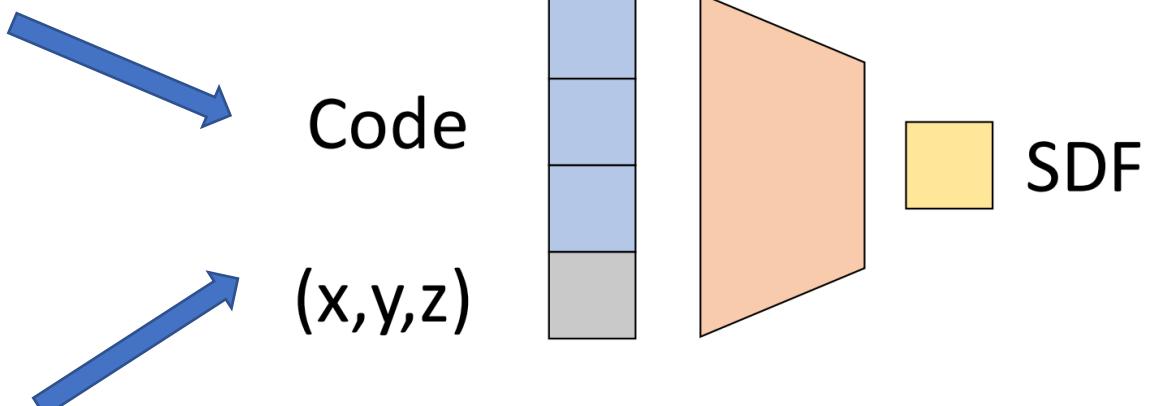
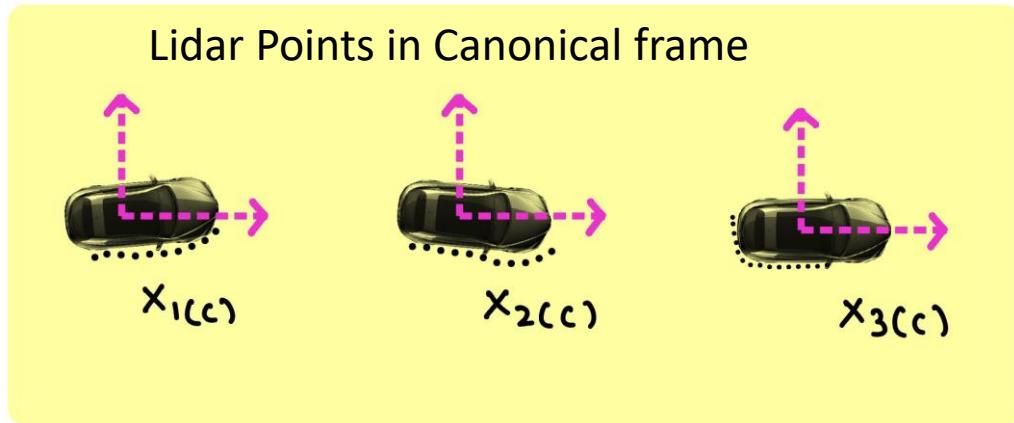


Code  
 $(x,y,z)$



SDF

[Latent Code Vector].T



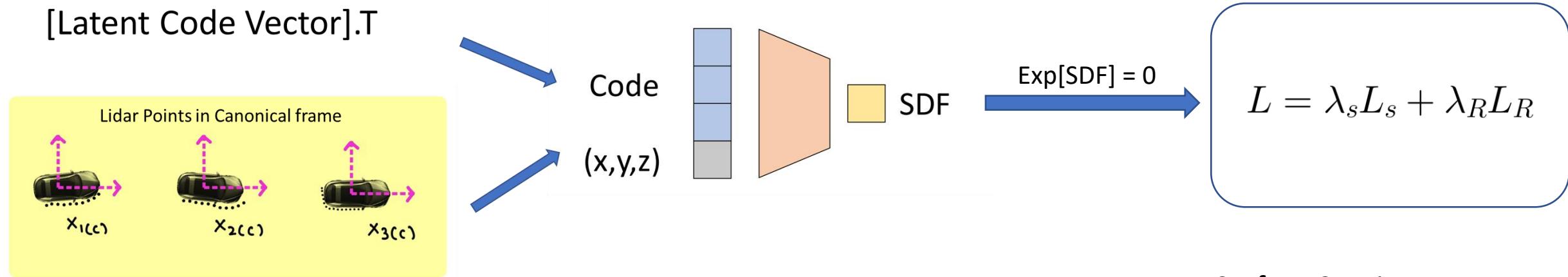
$$X_1[C] = P_1 X_1$$

$$X_2[C] = K_1 P_2 X_2$$

$$X_3[C] = K_1 K_2 P_3 X_3$$



Remember, points in canonical frame are a function of transformation predicted by PointRCNN



$$X_{1[C]} = P_1 X_1$$

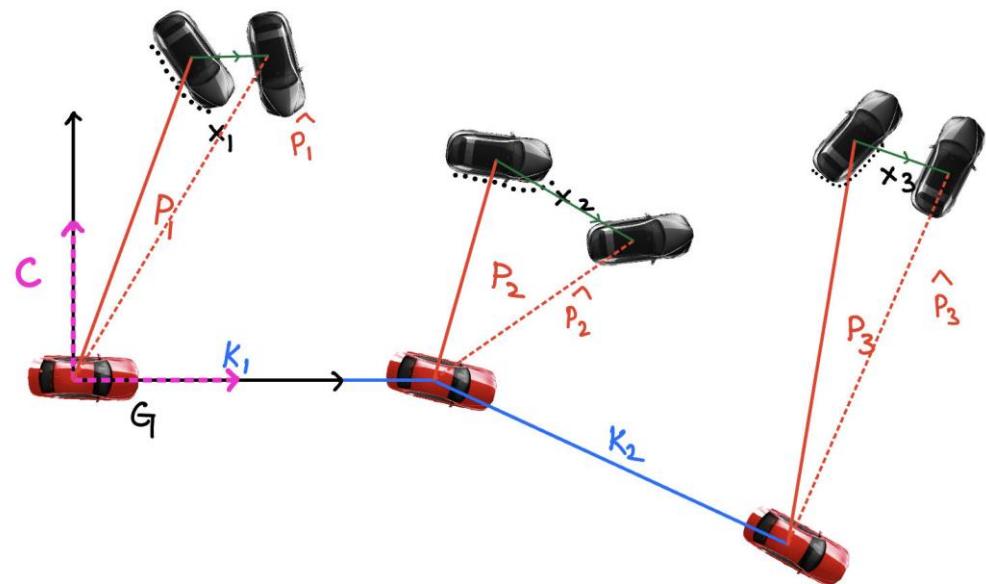
$$X_{2[C]} = K_1 P_2 X_2$$

$$X_{3[C]} = K_1 K_2 P_3 X_3$$

1. Surface Consistency
2. Latent Code Regularization

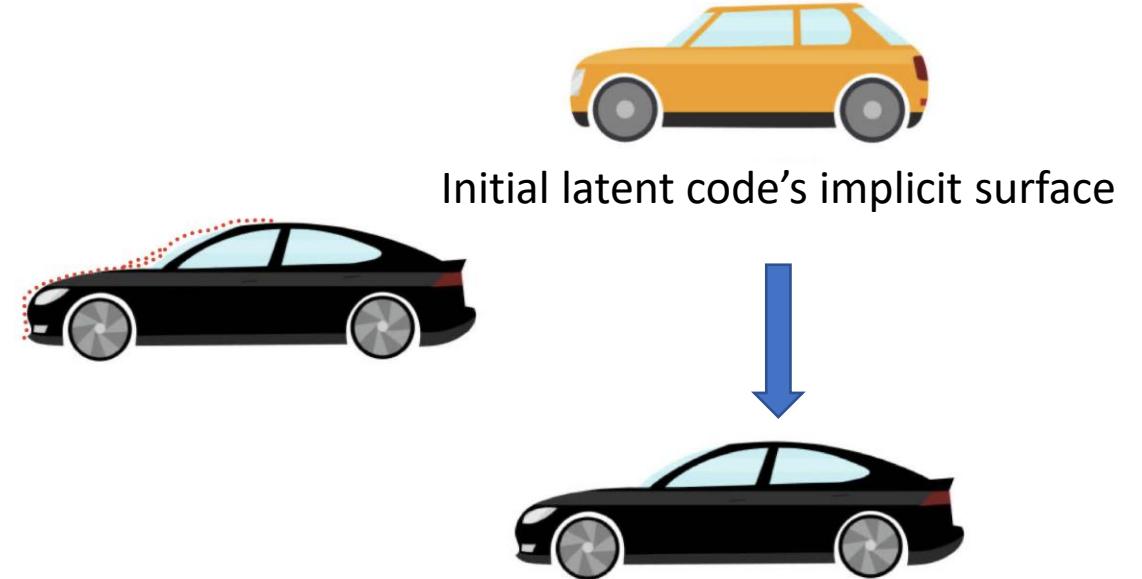
# Optimizing Latent Code and Pose Estimation

Optimized Pose Estimation



Optimized poses of dynamic obstacles

Optimized latent Code



Optimized latent code's implicit surface

# DataSets and Evaluation

Dynamic Object Pose Estimation

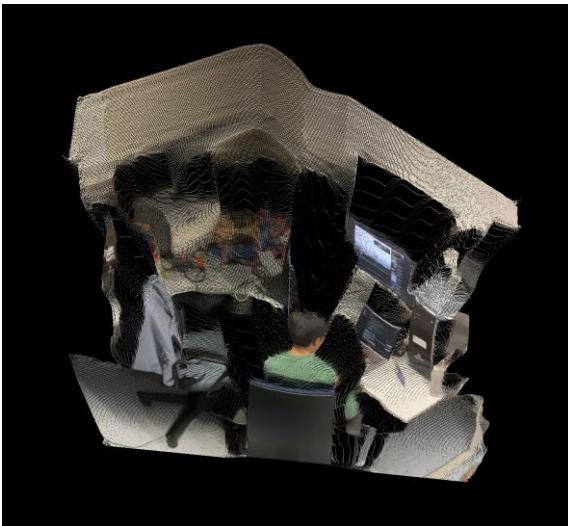




# 3D Reconstruction



Multiple Images

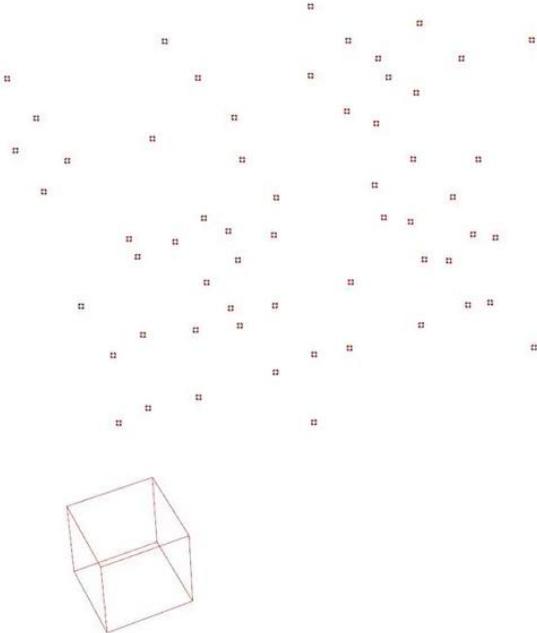
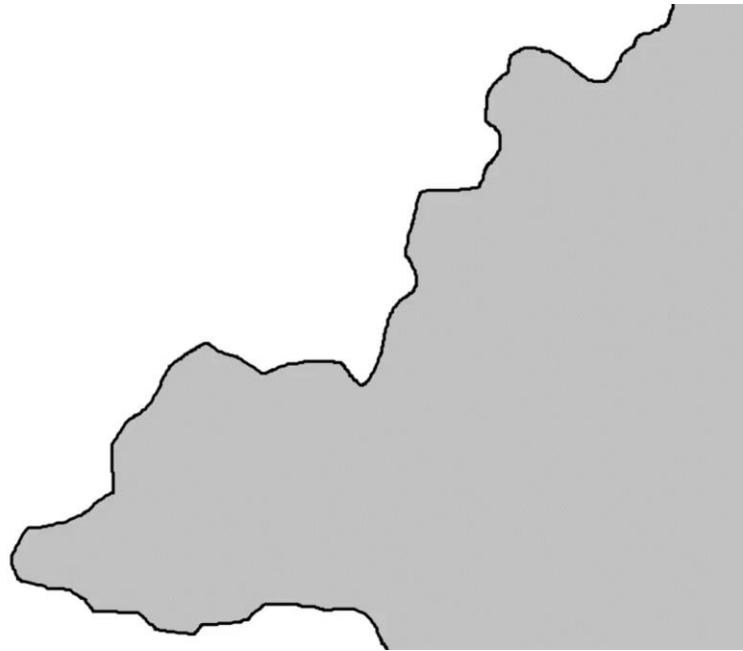


Partial Point Clouds



3D Scene

Read reasons for why shape completion is important from  
papers



## + :: M1: Basic knowledge learning and paper reading (~4 weeks)

- Pytorch learning: finish [\[A1\]](#) [\[A2\]](#) [\[A3\]](#) Assignments of eecs498. (Only if you know nothing about Pytorch before)
- Read and understand the following papers:
  - Cluster-VO [\[code\]](#)[\[paper\]](#)[\[video\]](#)
  - + :: ◦ DeepSDF [\[code\]](#)[\[paper\]](#)
  - Nerf [\[paper\]](#)
  - DSP-SLAM [\[code\]](#)[\[paper\]](#)[\[video\]](#)
  - Weakly Supervised Learning of Rigid 3D Scene Flow [\[code\]](#)[\[paper\]](#)
- Get familiar with the following datasets we will use:
  1. KITTI
  2. Nuscenes
  3. Waymo

EECS 498.008 / 598.008  
Deep Learning for Computer Vision  
Winter 2022

## Course Description

Computer Vision has become ubiquitous in our society, with applications in search, image understanding, apps, mapping, medicine, drones, and self-driving cars. Core to many of these applications are visual recognition tasks such as image classification and object detection. Recent developments in neural network approaches have greatly advanced the performance of these state-of-the-art visual recognition systems. This course is a deep dive into details of neural-network based deep learning methods for computer vision. During this course, students will learn to implement, train and debug their own neural networks and gain a detailed understanding of cutting-edge research in computer vision. We will cover learning algorithms, neural network architectures, and practical engineering tricks for training and fine-tuning networks for visual recognition tasks.

## + :: M2: Code Reading and demo test (~2 weeks)

- Read the code of DSP-SLAM and test it in different datasets, especially in dynamic scenes, and observe failure situations.
- Read the instance association part of Cluster-VO.

## M3: Build the main program framework (~4 weeks)

- Finish the data preprocessing part, including detection network and lidar odometry.
- Finish 2 in Methodology based on DSP-SLAM.
- Finish 3 in Methodology with the simplest method (bounding box overlap).

## M4: Make some contribution (~6 weeks)

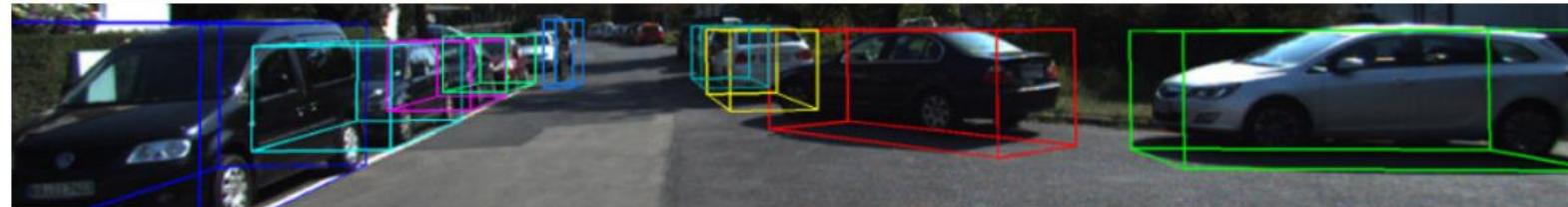
- Finish 4 in Methodology by designing a proper loss function.
- Explore more robust methods to achieve instance association.
- Explore the complete differentiable process from instance association to motion estimation.

## M5: Final Result: Evaluation (~4 weeks)

- Use scene flow dataset and metric to evaluate motion estimation.
- Design another metric to evaluate instance association with panoptic segmentation dataset.

# The KITTI Vision Benchmark Suite

## 3D Object Detection Evaluation 2017



The 3D object detection benchmark consists of 7481 training images and 7518 test images as well as the corresponding point clouds, comprising a total of 80.256 labeled objects. For evaluation, we compute precision-recall curves. To rank the methods we compute average precision. We require that all methods use the same parameter set for all test pairs. Our development kit provides details about the data format as well as MATLAB / C++ utility functions for reading and writing the label files.

- [Download left color images of object data set \(12 GB\)](#)
- [Download right color images, if you want to use stereo information \(12 GB\)](#)
- [Download the 3 temporally preceding frames \(left color\) \(36 GB\)](#)
- [Download the 3 temporally preceding frames \(right color\) \(36 GB\)](#)
- [Download Velodyne point clouds, if you want to use laser information \(29 GB\)](#)
- [Download camera calibration matrices of object data set \(16 MB\)](#)
- [Download training labels of object data set \(5 MB\)](#)
- [Download object development kit \(1 MB\) \(including 3D object detection and bird's eye view evaluation code\)](#)
- [Download pre-trained LSVM baseline models \(5 MB\) used in Joint 3D Estimation of Objects and Scene Layout \(NIPS 2011\). These models are referred to as LSVM-MDPM-sv \(supervised version\) and LSVM-MDPM-us \(unsupervised version\) in the tables below.](#)
- [Download reference detections \(L-SVM\) for training and test set \(800 MB\)](#)
- Qianli Liao (NYU) has put together [code to convert from KITTI to PASCAL VOC file format](#) (documentation included, requires Emacs).
- Karl Rosaen (U.Mich) has released [code to convert between KITTI, KITTI tracking, Pascal VOC, Udacity, CrowdAI and AUTTI formats](#).
- Jonas Heylen (TRACE vzw) has released [pixel accurate instance segmentations](#) for all 7481 training images.
- We thank [David Stutz](#) and [Bo Li](#) for developing the 3D object detection benchmark.