

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

## Team members

Panyawat (Ohm) Rattana  
Shashank (Sai) Dammalapati

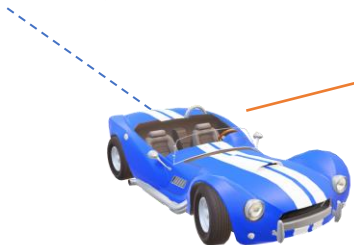
## Supervisors

Xingguang (Starry) Zhong  
Yue Pan

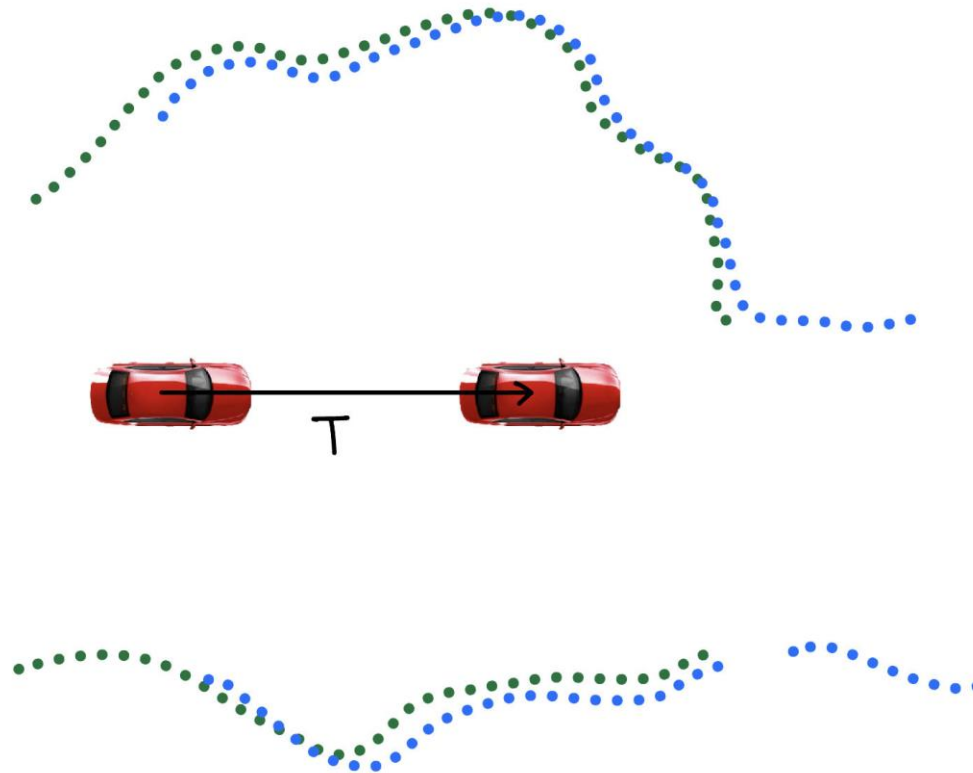
Motion Estimation



Odometry



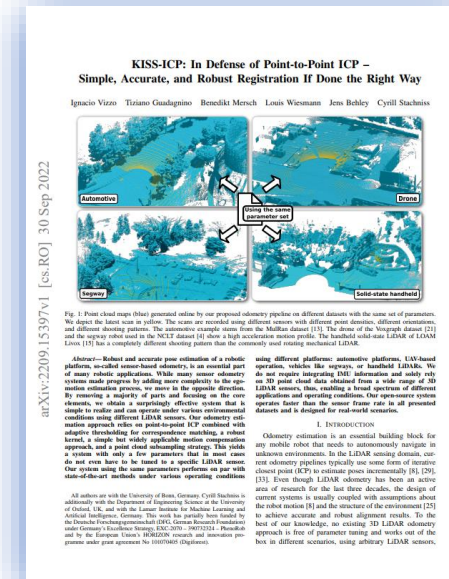
# Lidar Odometry (KISS ICP)



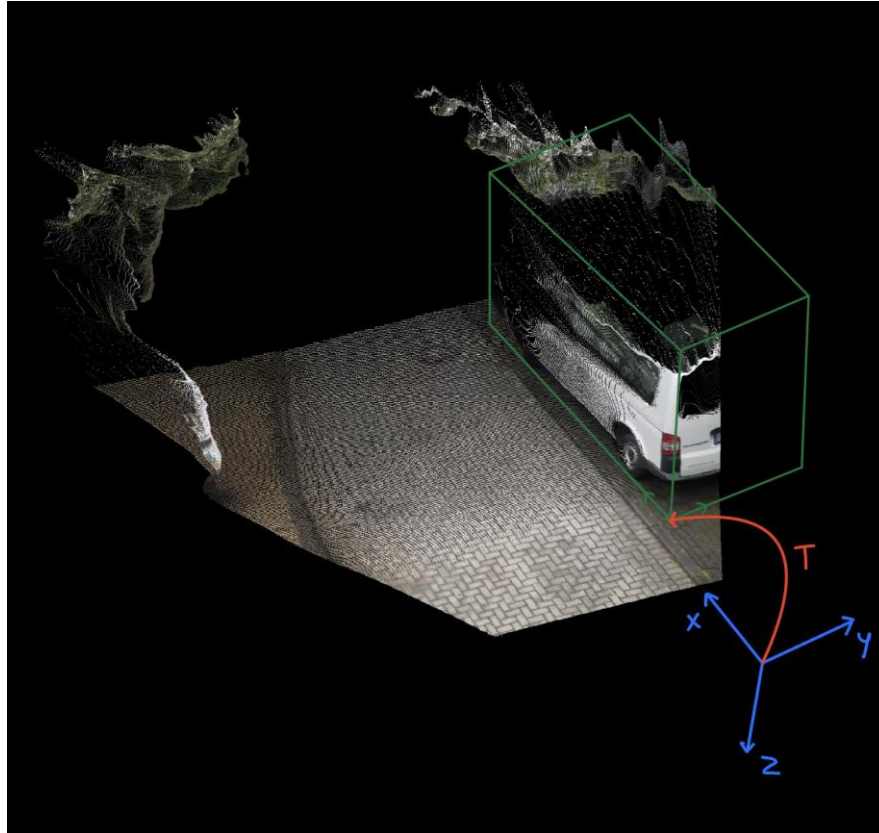
Motion Estimation of Ego Car

“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.

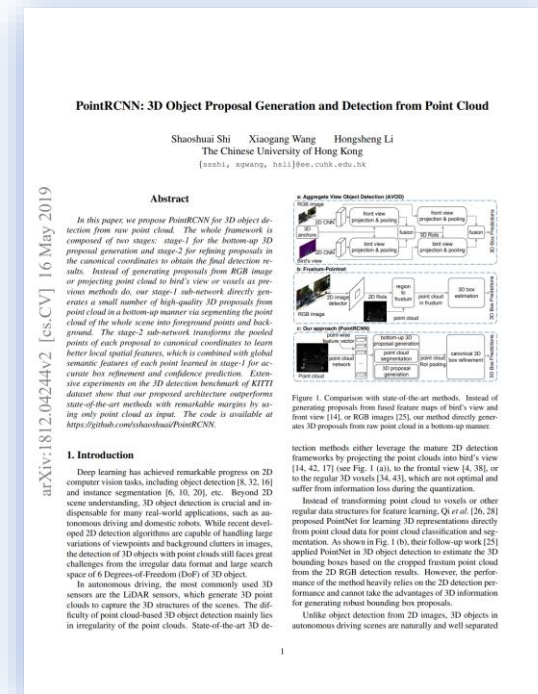


# Point RCNN

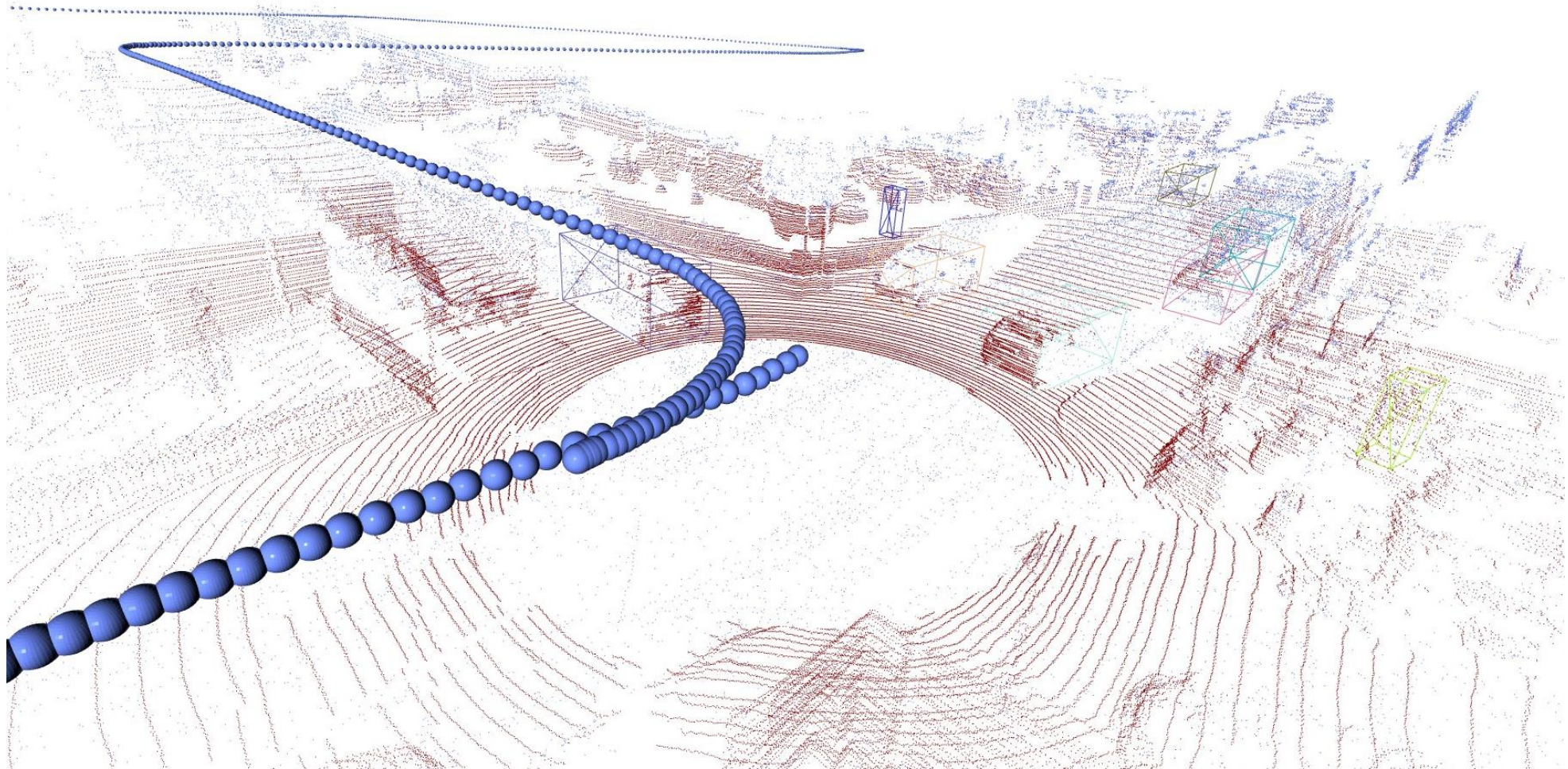


- Idea is to generate bounding boxes for cars
- Point RCNN gives us an initial guess for the relative pose of the external cars relative to ego car.

Pose Estimation of surrounding cars in 1 frame

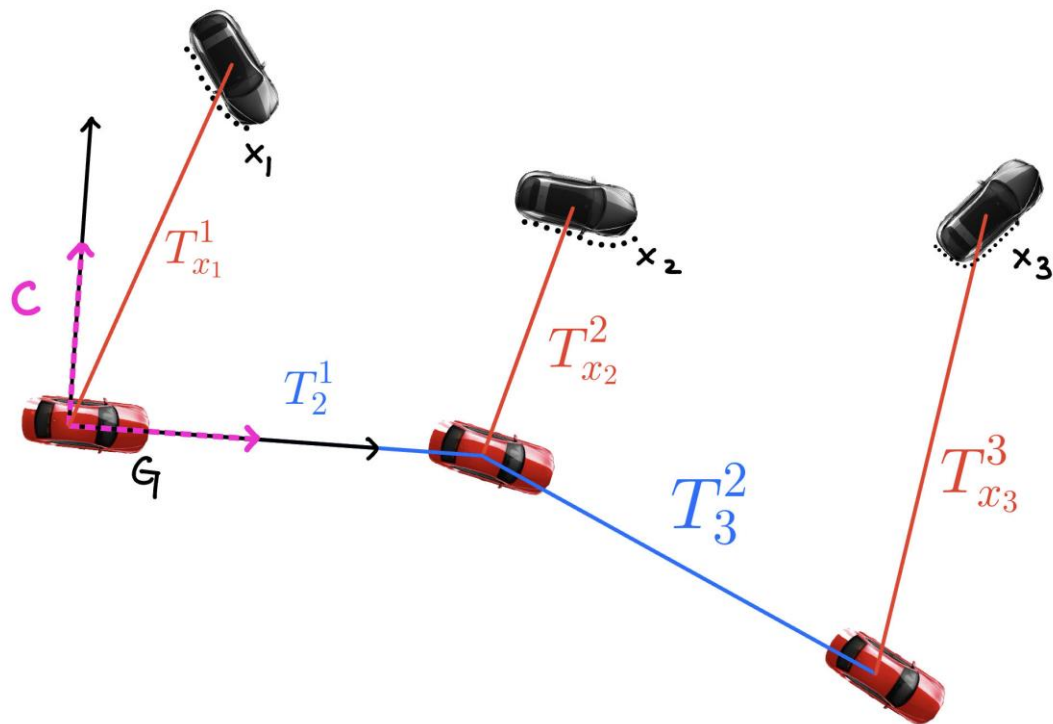






KISS ICP + Point RCNN

# Workflow: Transformation



C – Canonical Frame

G – Global Frame

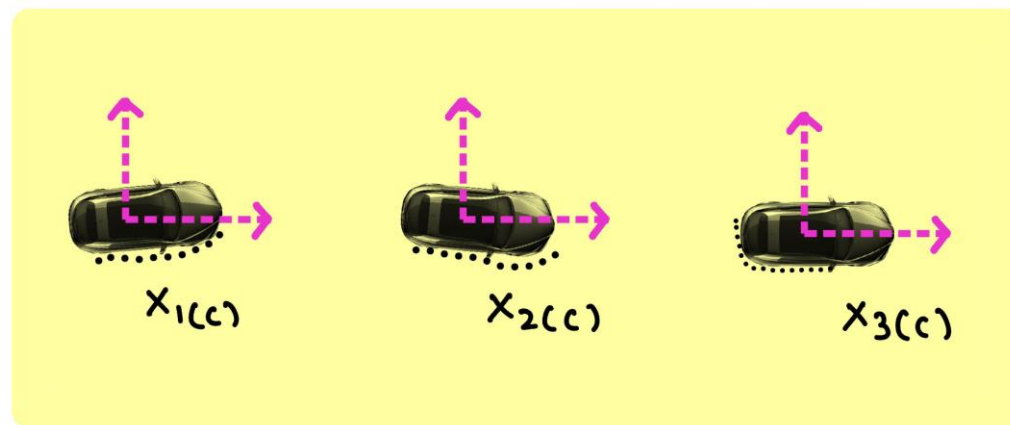
$T^N_{x_N}$  Transformation from Point RCNN

$T^{N-1}_N$  Transformation from KISS ICP

$$X_1[C] = T^1_{x_1} X_1$$

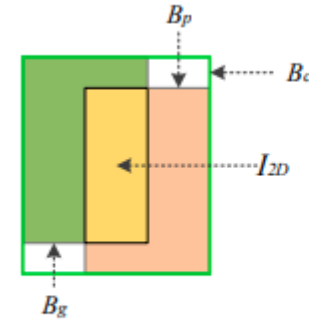
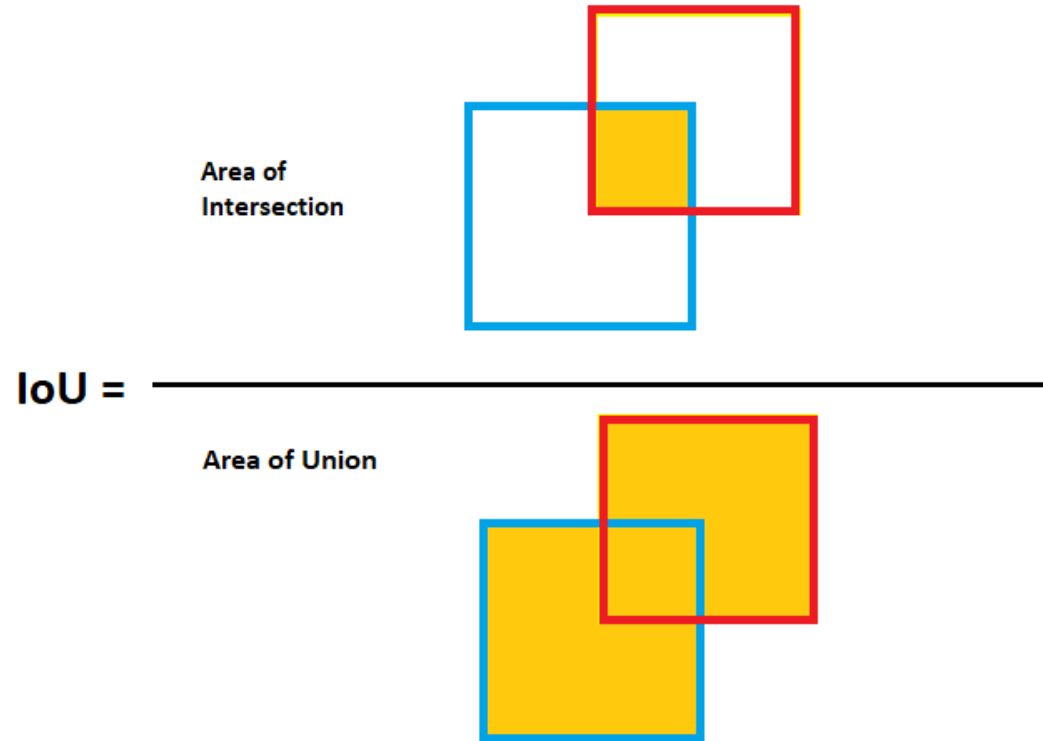
$$X_2[C] = T^1_2 T^2_{x_2} X_2$$

$$X_3[C] = T^1_2 T^2_3 T^3_{x_3} X_3$$

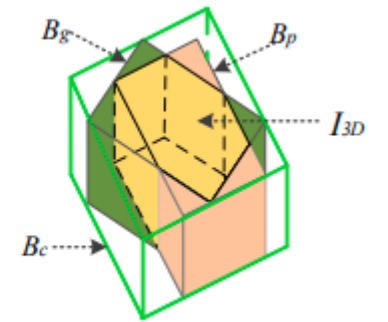


Point Clouds in Canonical Space

# IOU: Definition



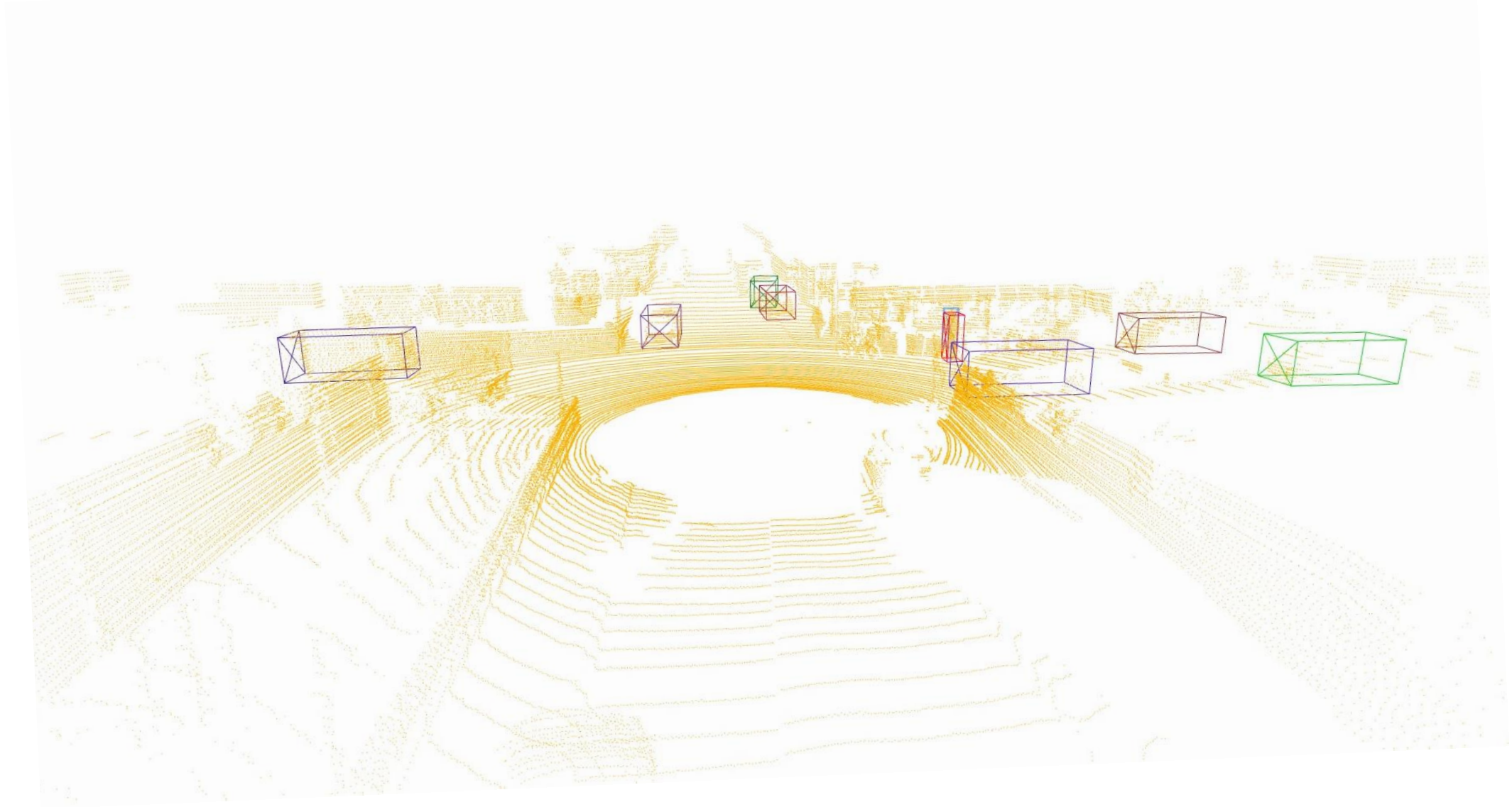
(a) 2D

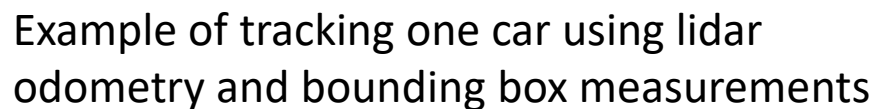


(b) 3D



# Video

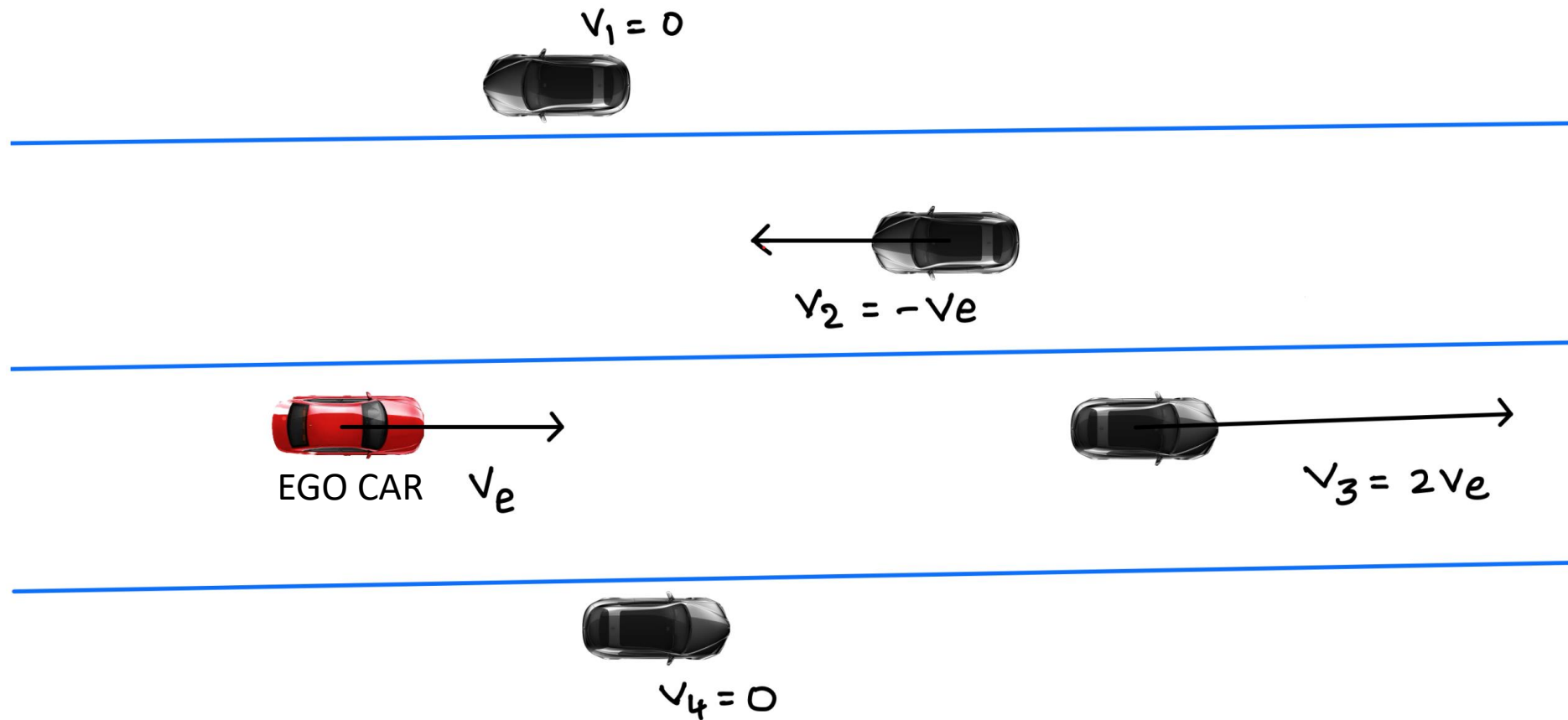






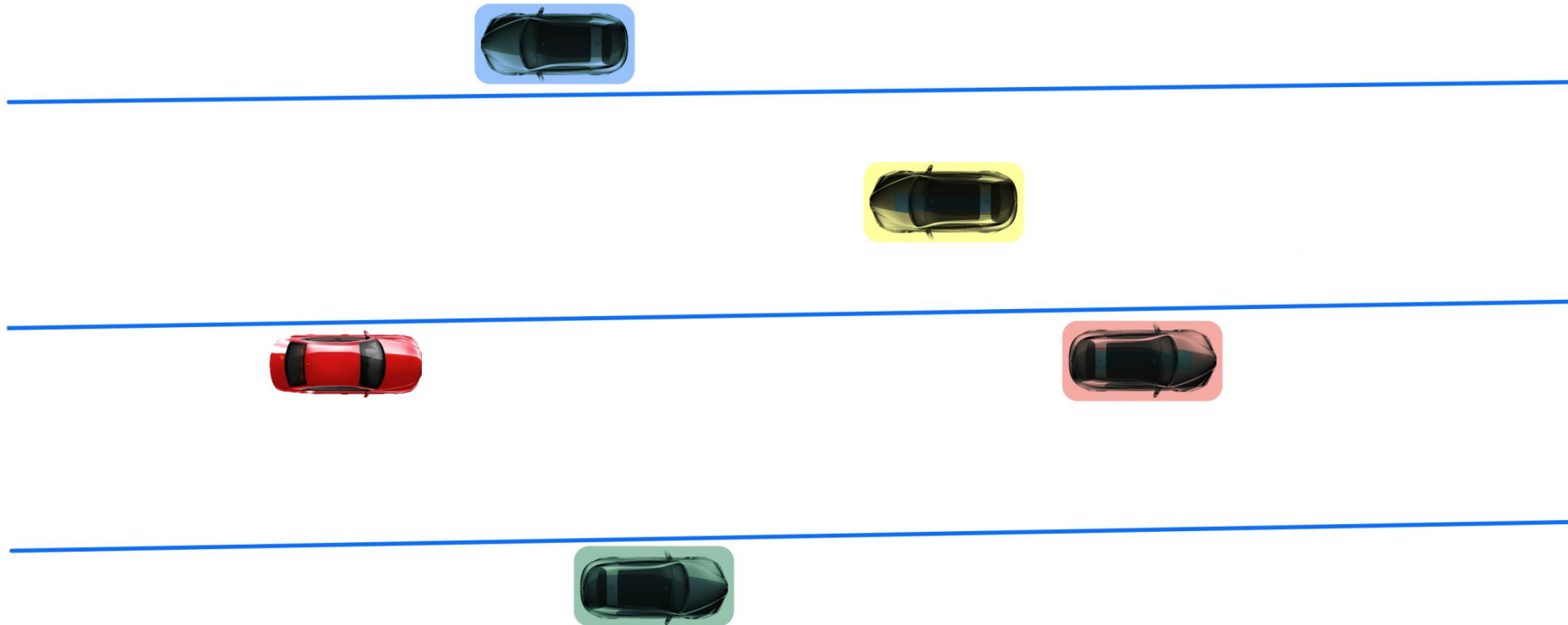
# Tracking Multiple Objects

Global frame



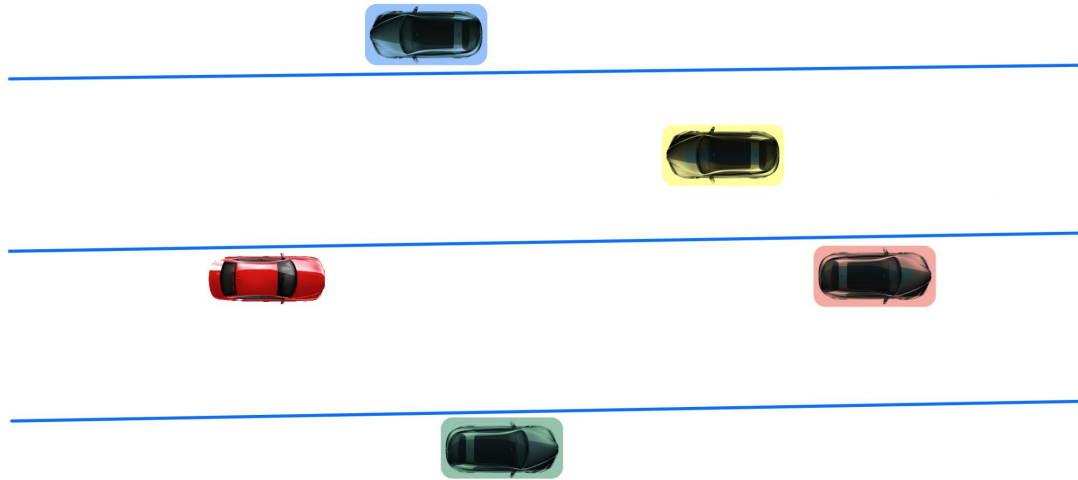
# Prediction of Bounding Boxes

Global frame

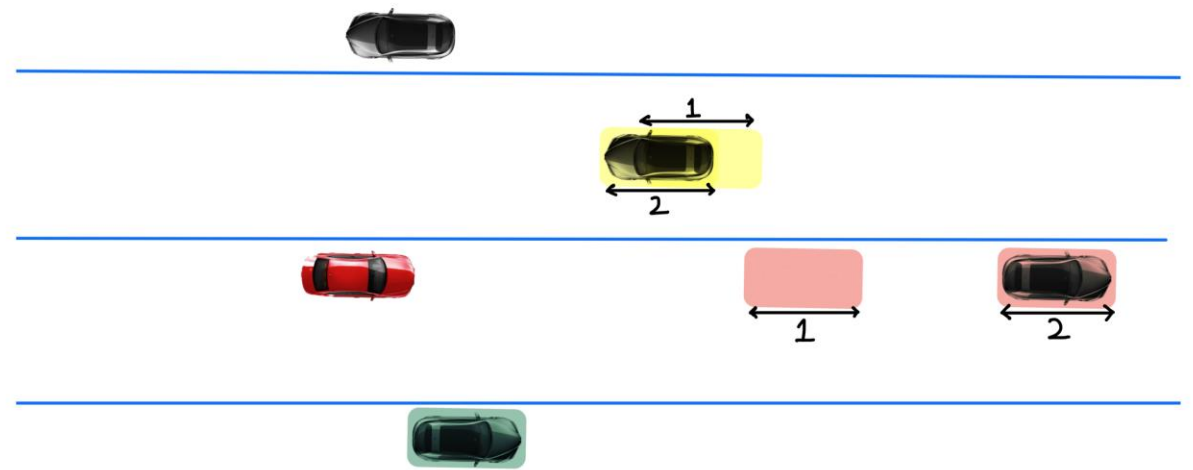


# Basic IOU problems

Global frame



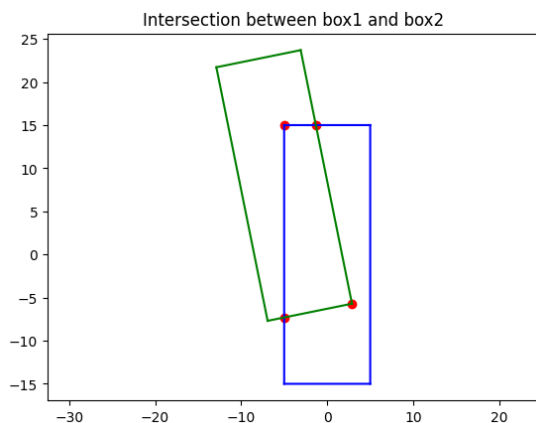
Global frame



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?



0.266 IOU



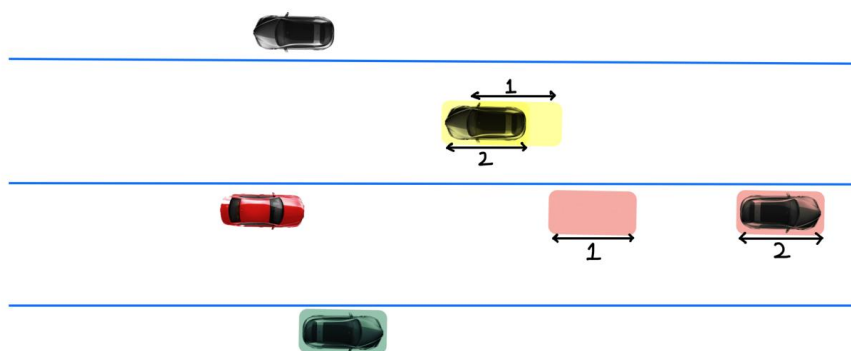
PyTorch3D IoU

Improving IoU calculation

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)

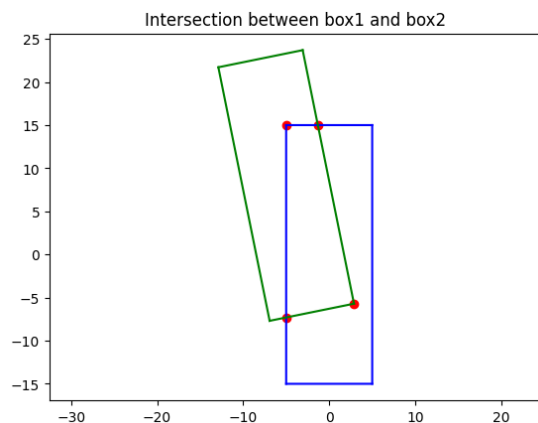
Global frame



More robust instance association using motion models



# How can we improve instance tracking?



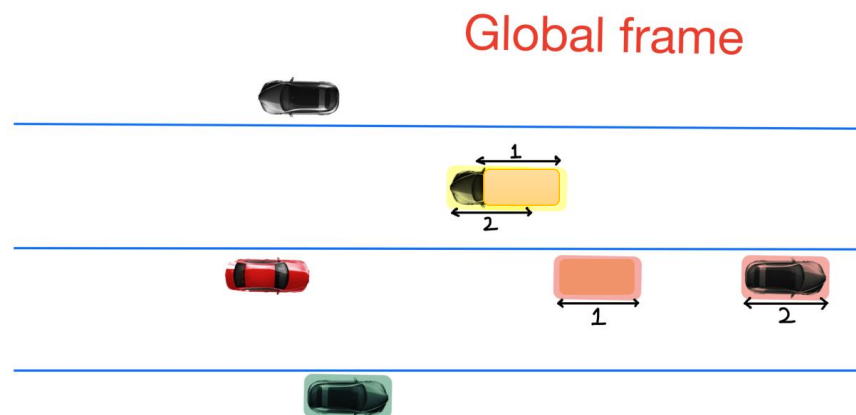
0.266 IOU



Improving IoU calculation

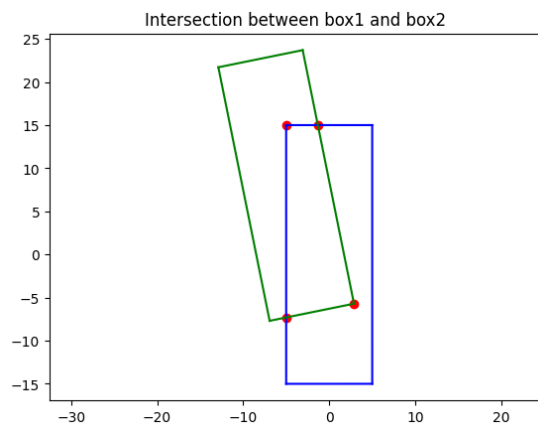
Frame	Cars in Scene	Cars Detected	Buffer
1	1 2 3	1 2 3	1 2 3
2	1 2 3 4	1 3 4	1 3 4
3	2 3 4	2 3 4	2 3 4

If (Buffer = Detections)



More robust instance association track cars' velocities

# How can we improve instance tracking?



0.266 IOU



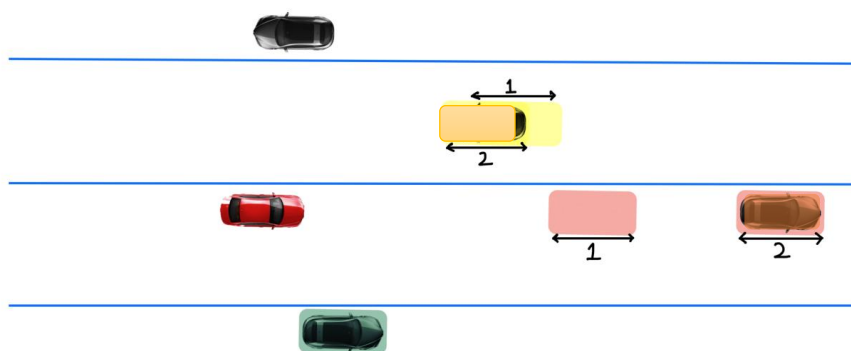
PyTorch3D IoU

Improving IoU calculation

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

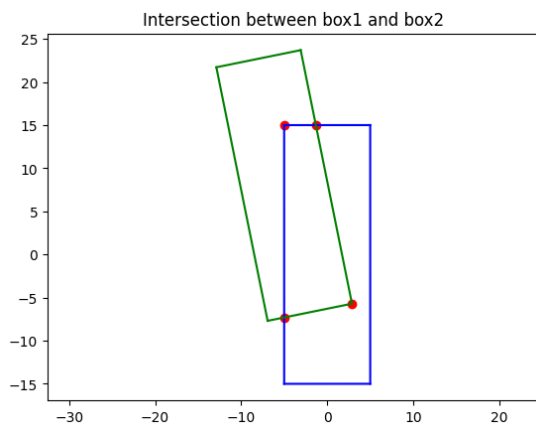
If (Buffer = Detections)

Global frame



More robust instance association track cars' velocities

# How can we improve instance tracking?

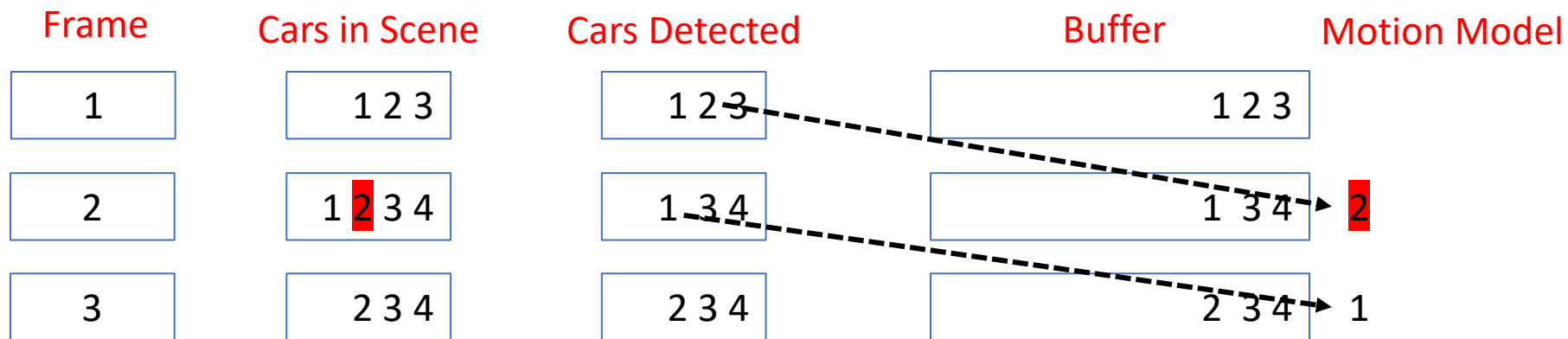


0.266 IOU

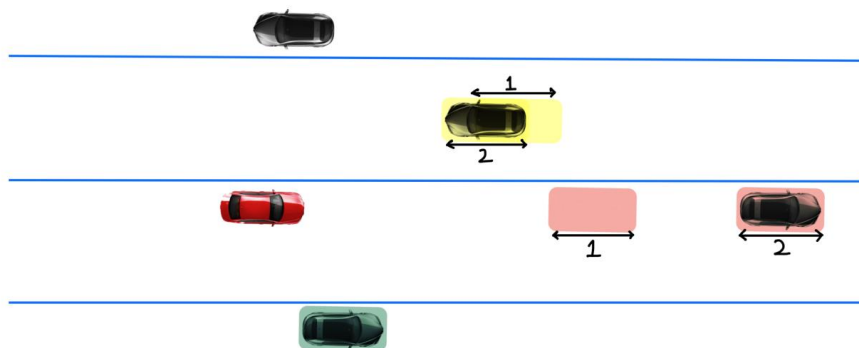


PyTorch3D IoU

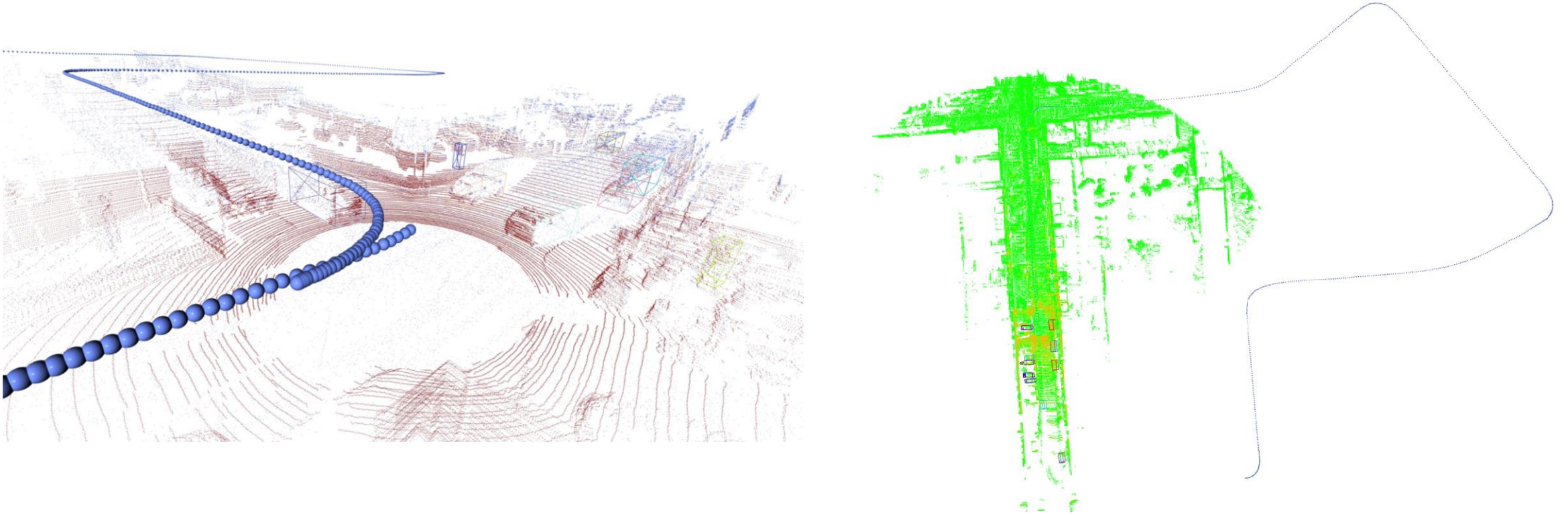
Improving IoU calculation



Global frame



More robust instance association track cars' velocities



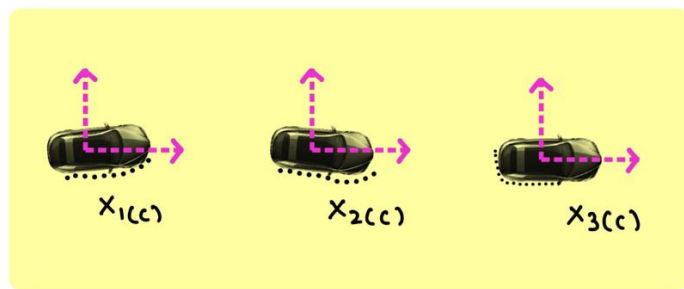
After proper instance association, we can use DeepSDF to jointly optimize shape and pose of the detected moving cars



Optimize latent Code through back propagation



[Latent Code Vector].T



1. Surface Consistency
2. Latent Code Regularization

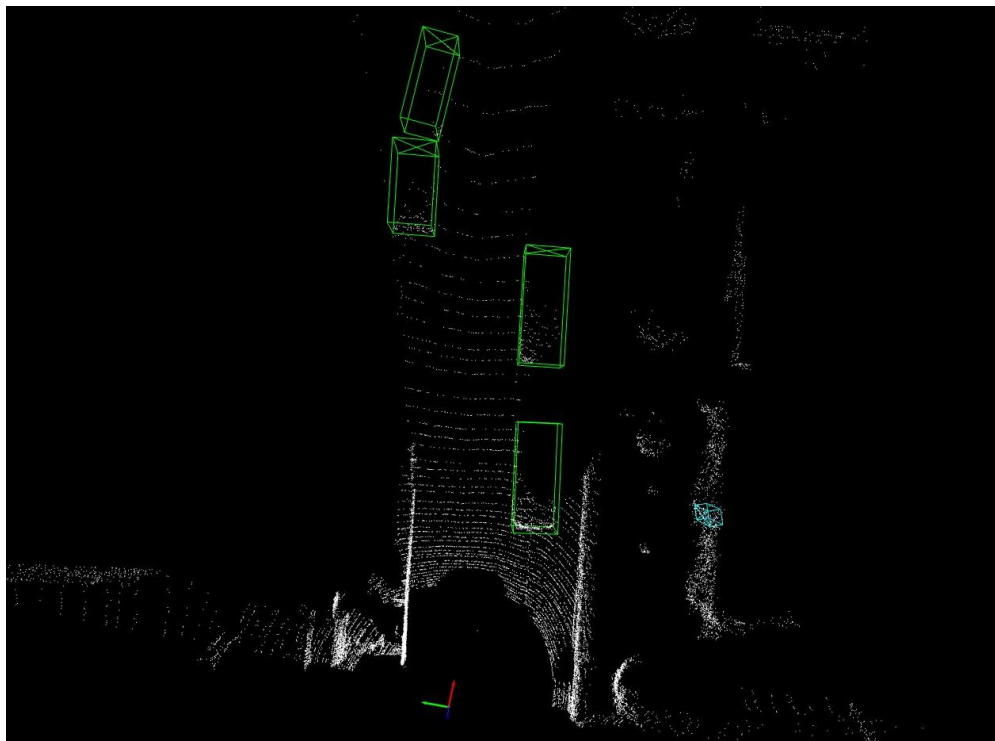
$$X_{1[C]} = T_{x_1}^1 X_1$$

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

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

Optimize Pose transformations through back propagation



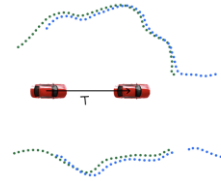
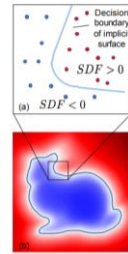
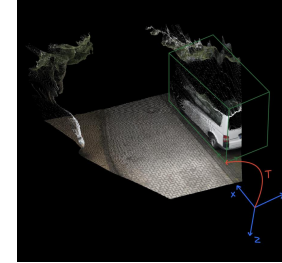
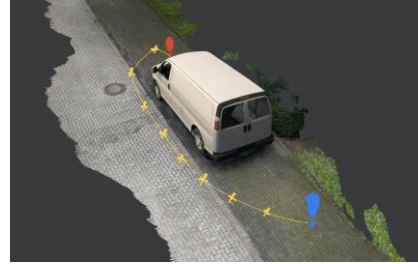


Lidar + Bounding Box

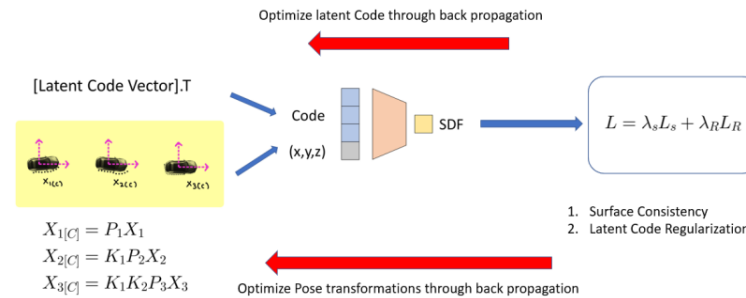


Stereo Image + Shape Completion

# Overview of the project



Challenges

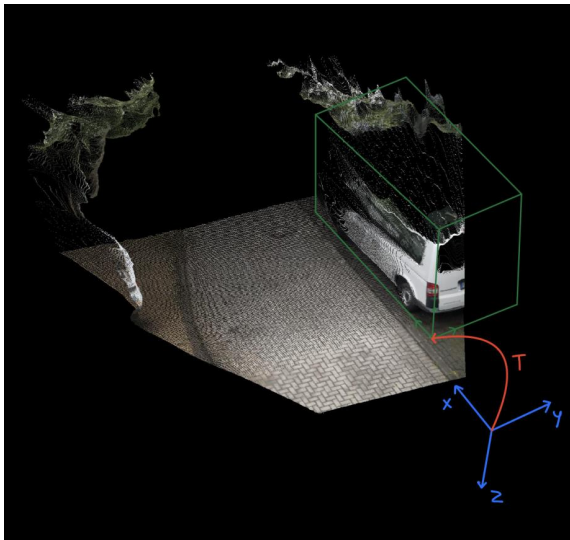


Method

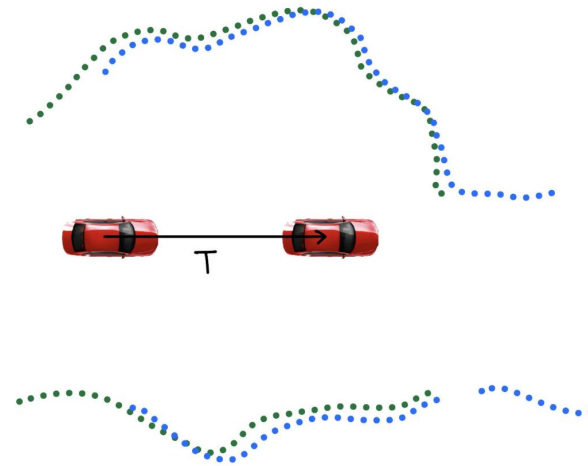


Final Output

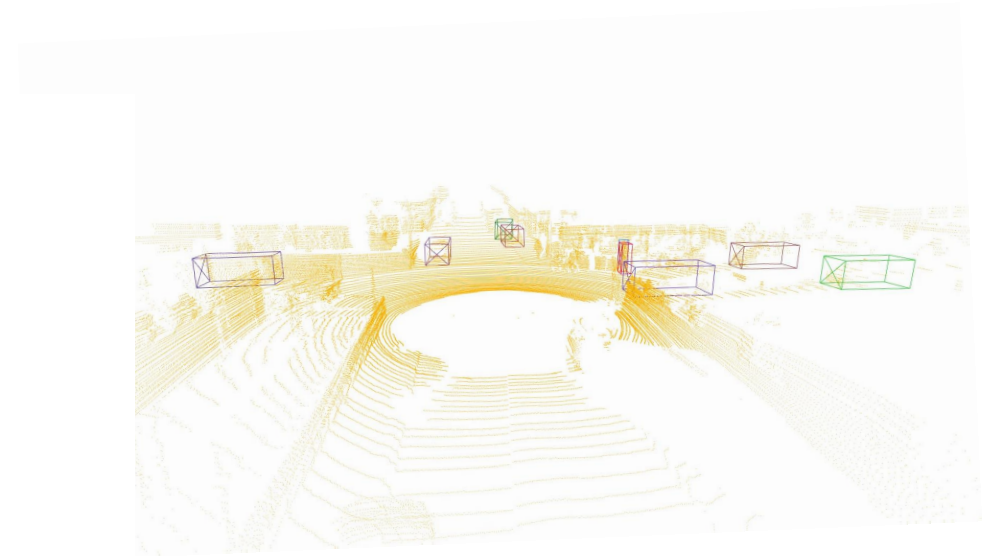
# Completed Tasks



Bounding Box Prediction



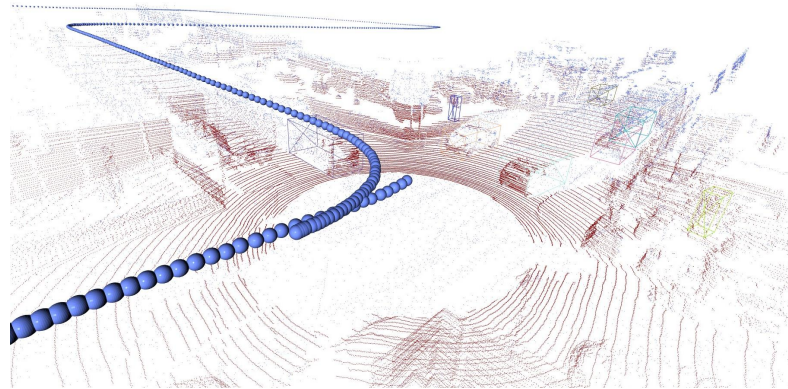
Lidar Odometry



Instance Association and  
Multi Object Tracking



# Next Steps



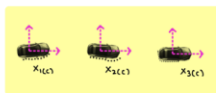
Pose and  
Point  
Cloud of  
the cars



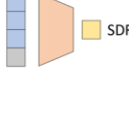
Optimized  
Pose and  
Shape Code  
of the cars

Optimize latent Code through back propagation

[Latent Code Vector].T



Code  
(x,y,z)



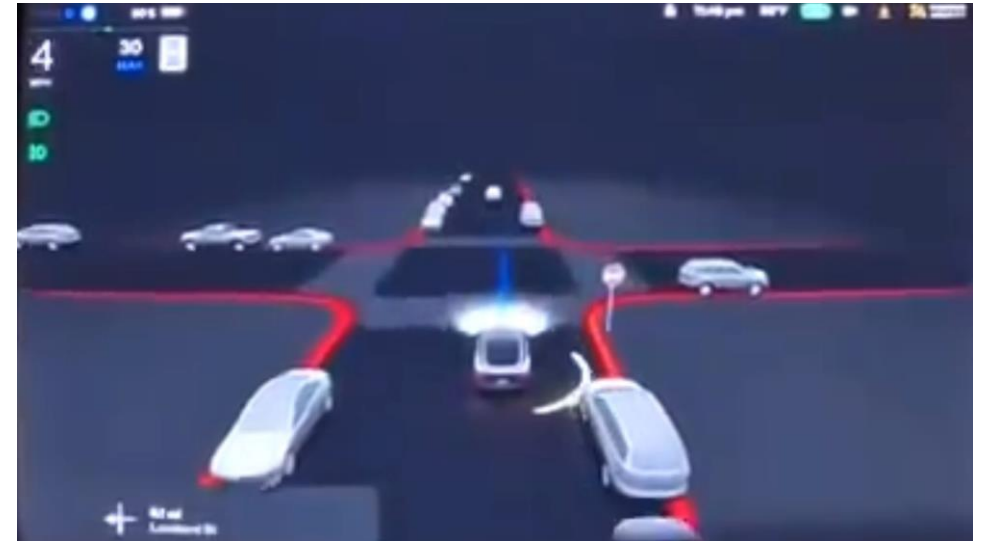
$$L = \lambda_s L_s + \lambda_R L_R$$

$$\begin{aligned} X_{1[C]} &= P_1 X_1 \\ X_{2[C]} &= K_1 P_2 X_2 \\ X_{3[C]} &= K_1 K_2 P_3 X_3 \end{aligned}$$

Optimize Pose transformations through back propagation

1. Surface Consistency
2. Latent Code Regularization

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

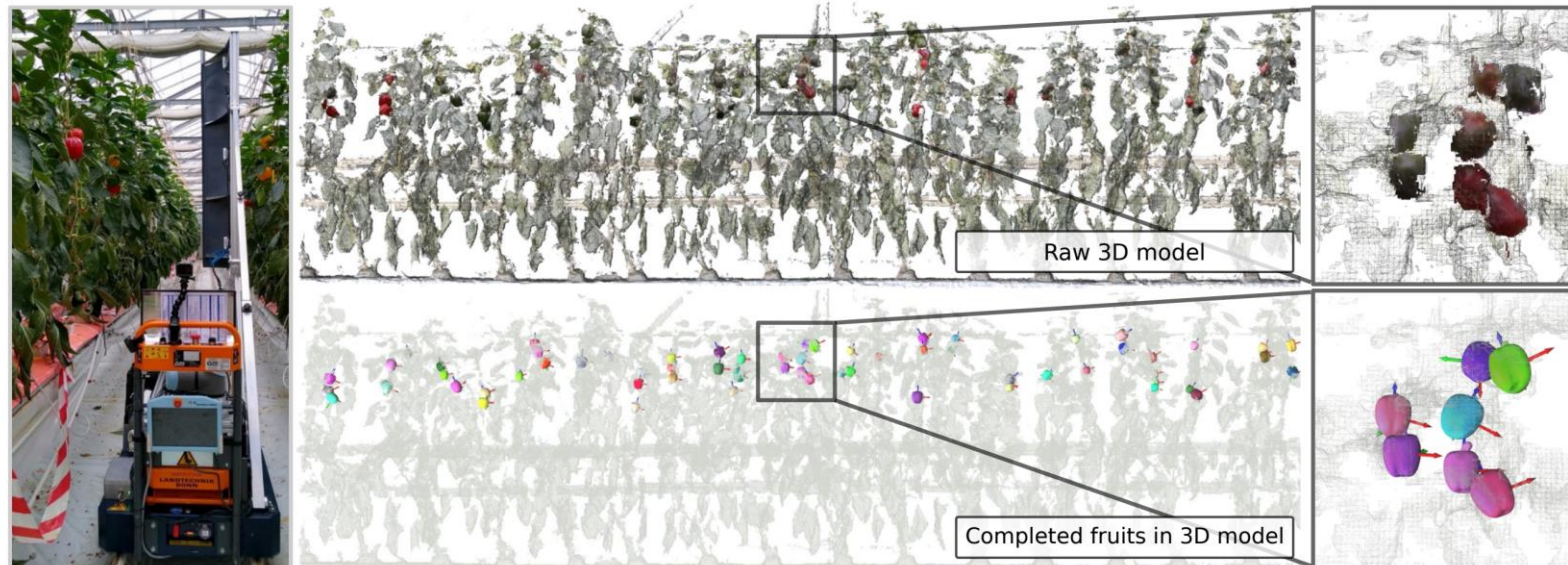
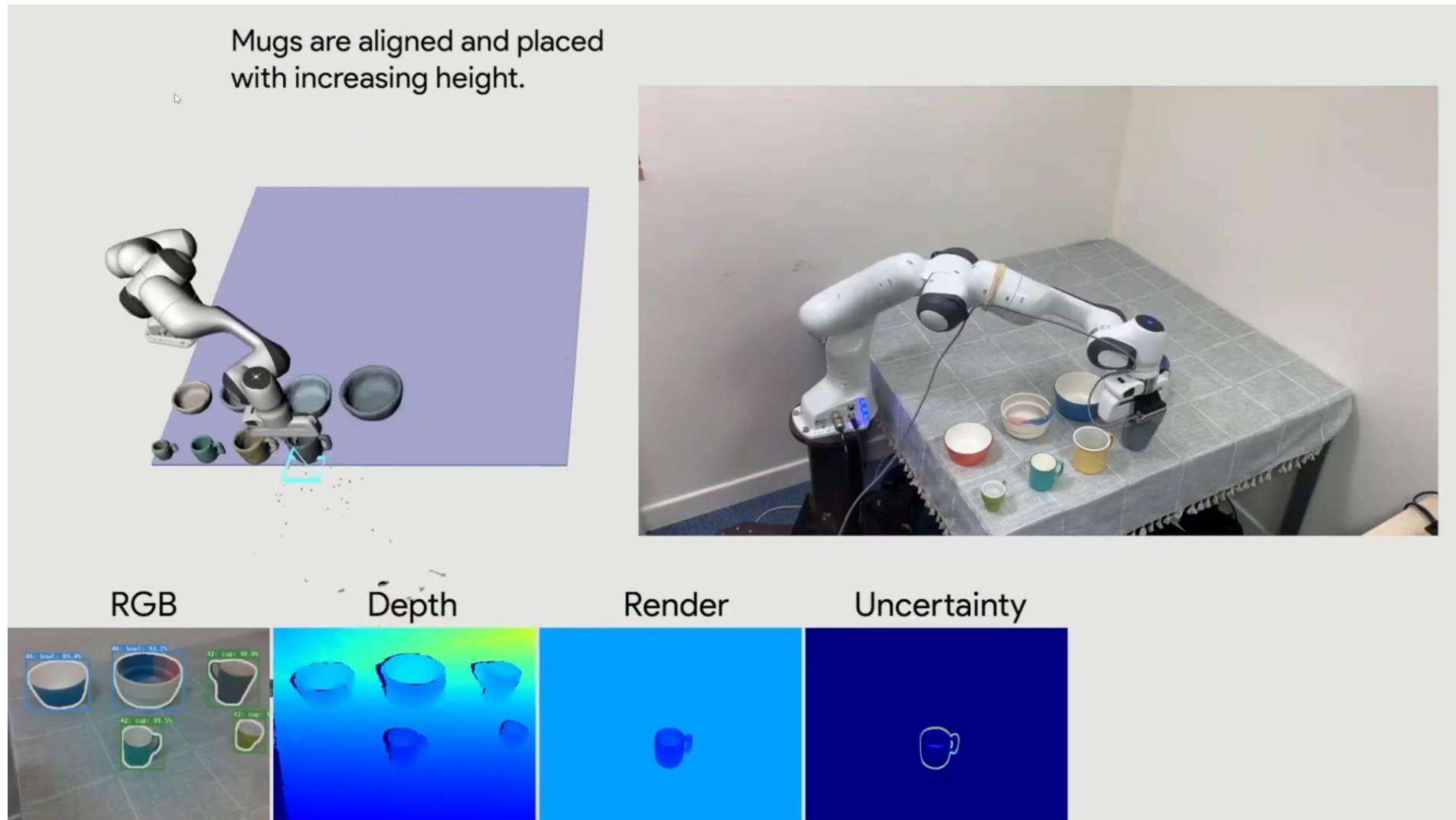


Fig. 1: Our method is able to build a multi-resolution panoptic map (top) of a challenging commercial glasshouse environment online using a mobile horticultural robot equipped with RGB-D cameras (left). Furthermore, our method manages to jointly estimate the complete shape and pose of each fruit in the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

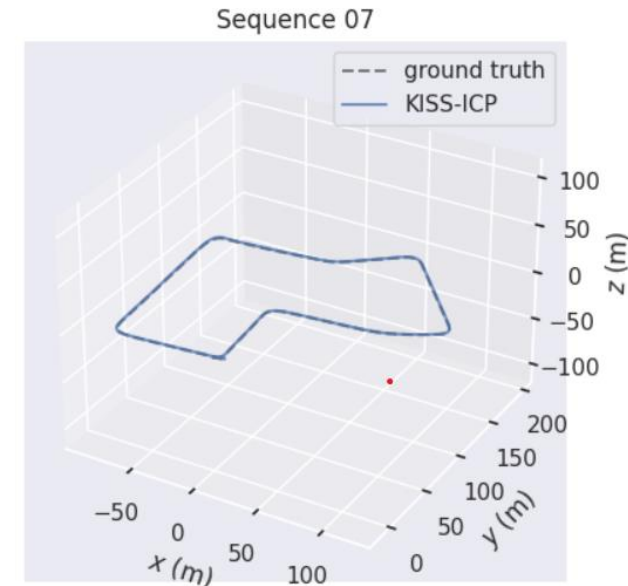


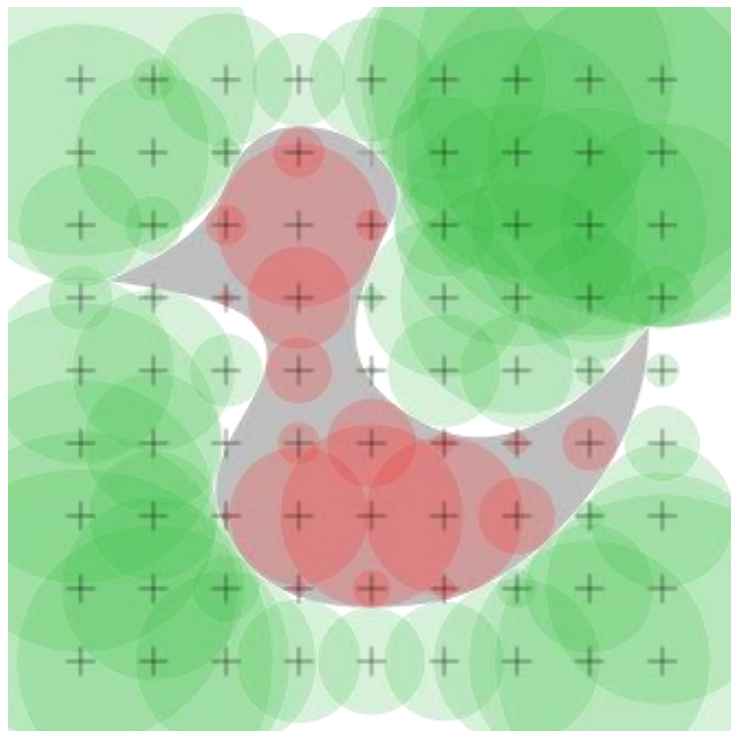


- PointRCNN - 8 fps
- Kiss ICP – 60 fps
- Shape Completion – 5 fps (avg)
- Kiss ICP + Bounding Box (Stored) – 30 fps
- DSP SLAM – 5 fps

Now evaluating sequence 07  
0%|

	Metric	Value	Units
Average Translation Error		0.328	%
Average Rotational Error		0.164	deg/m
Absolute Trajectory Error (ATE)		0.776	m
Absolute 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

$(x, y, z)$

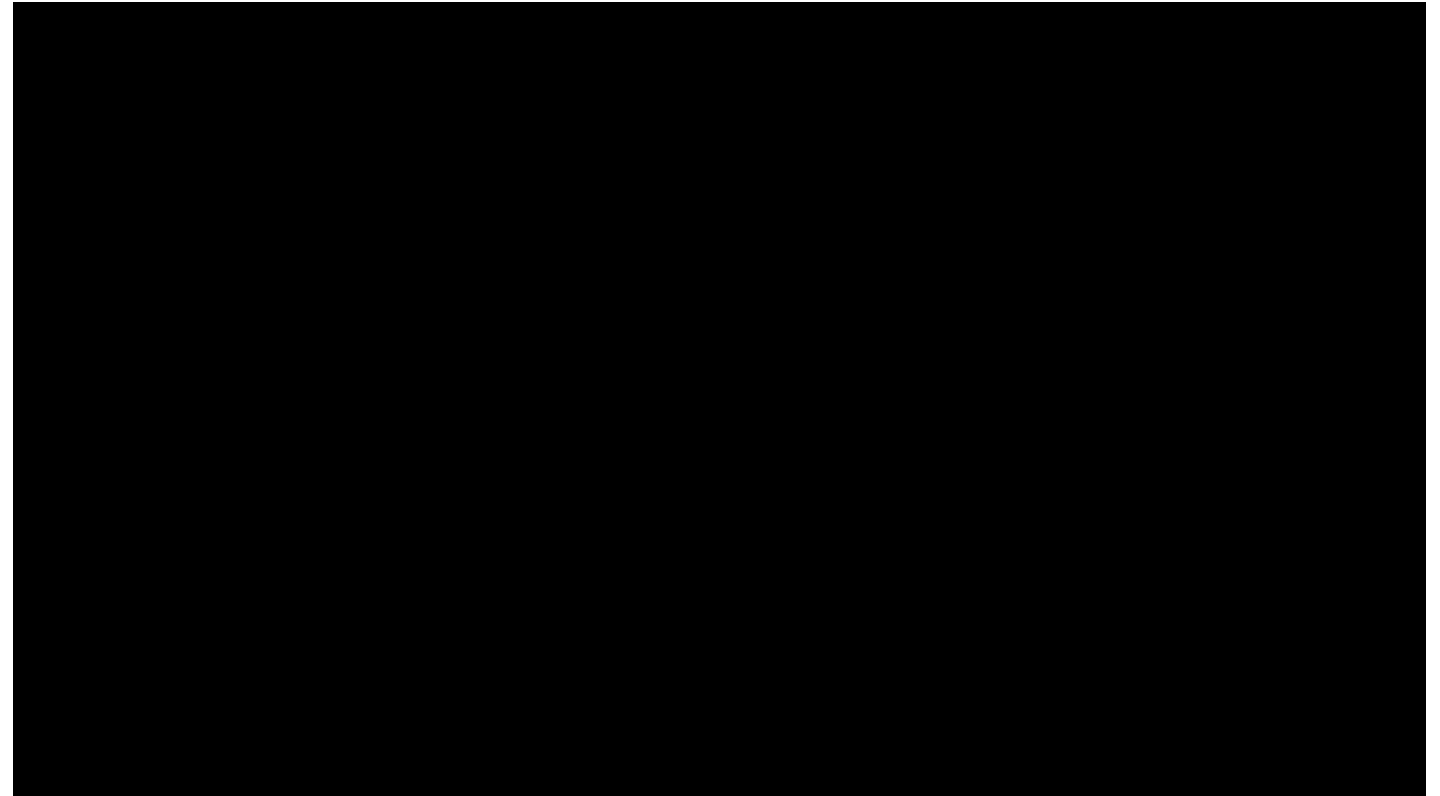
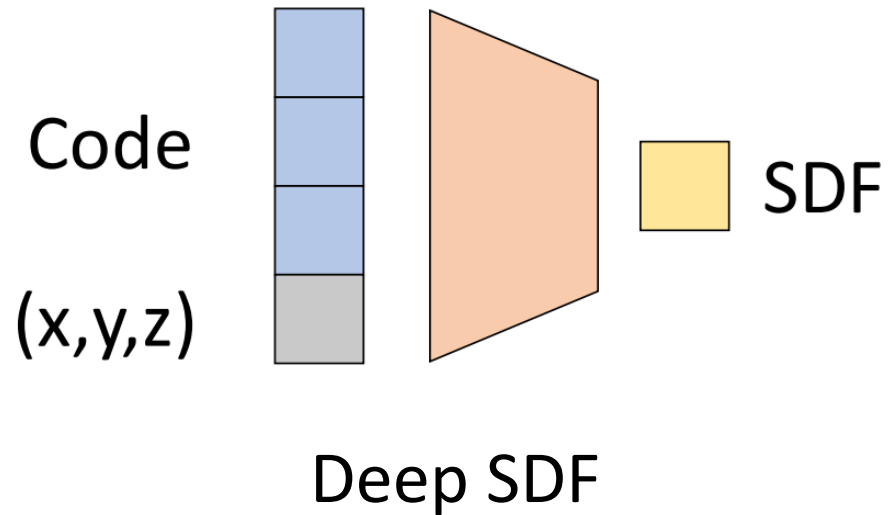


SDF

Deep SDF

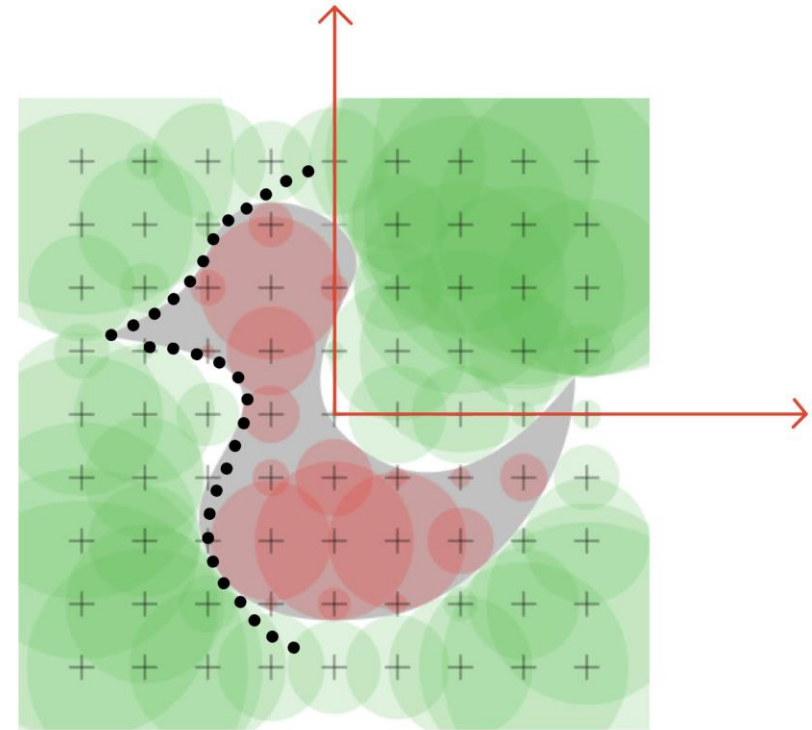
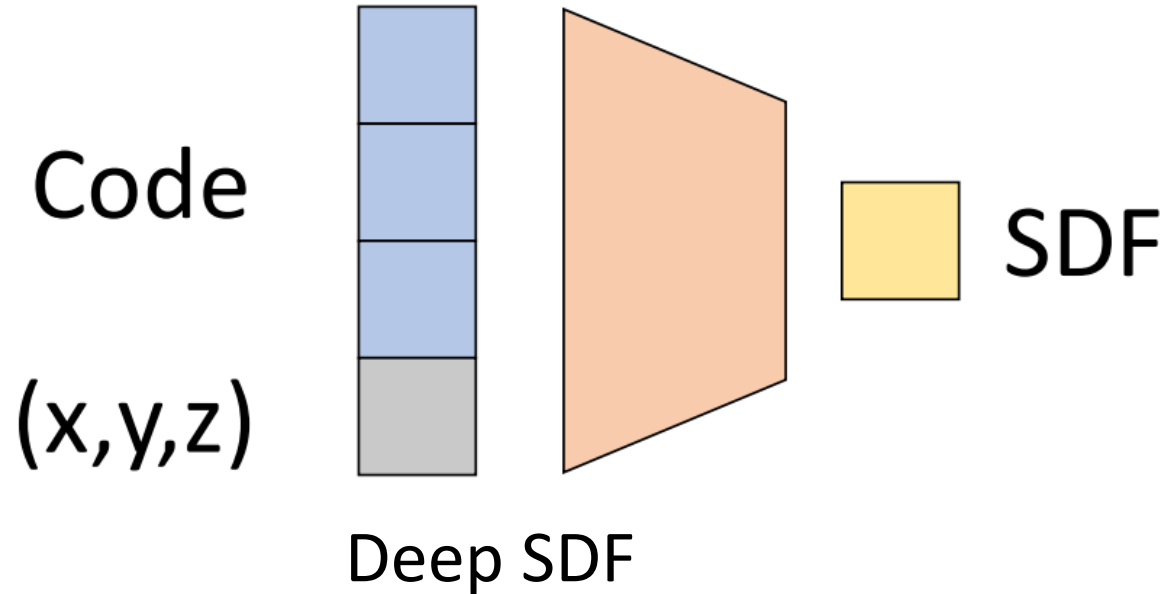
DeepSDF can implicitly model a class of objects

# Deep SDF



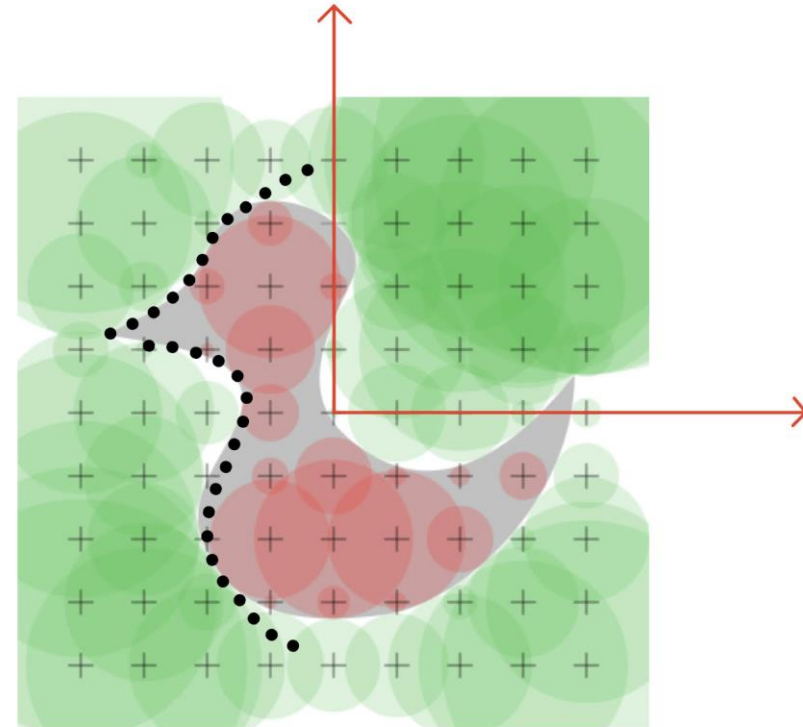
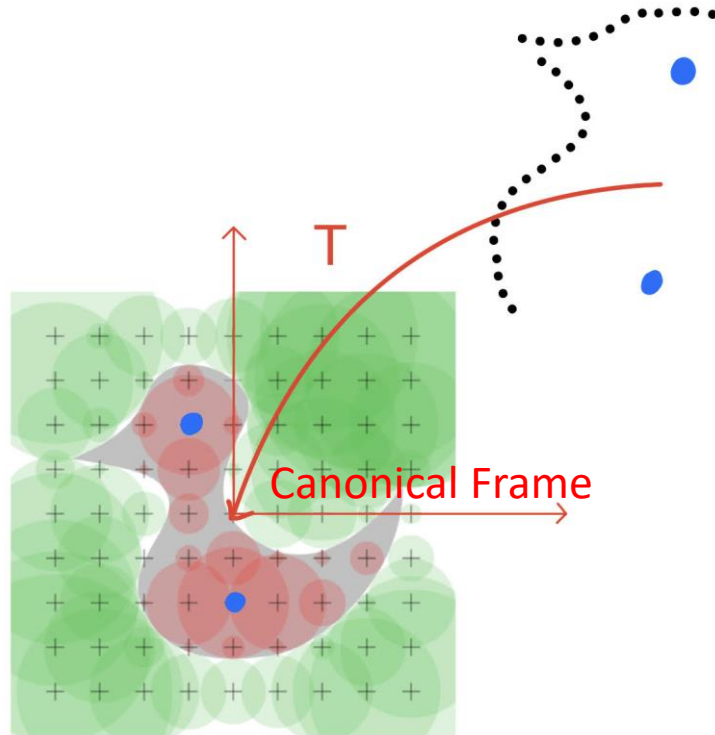
A latent code collapses a general representation into in a single shape

# 0 SDF



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

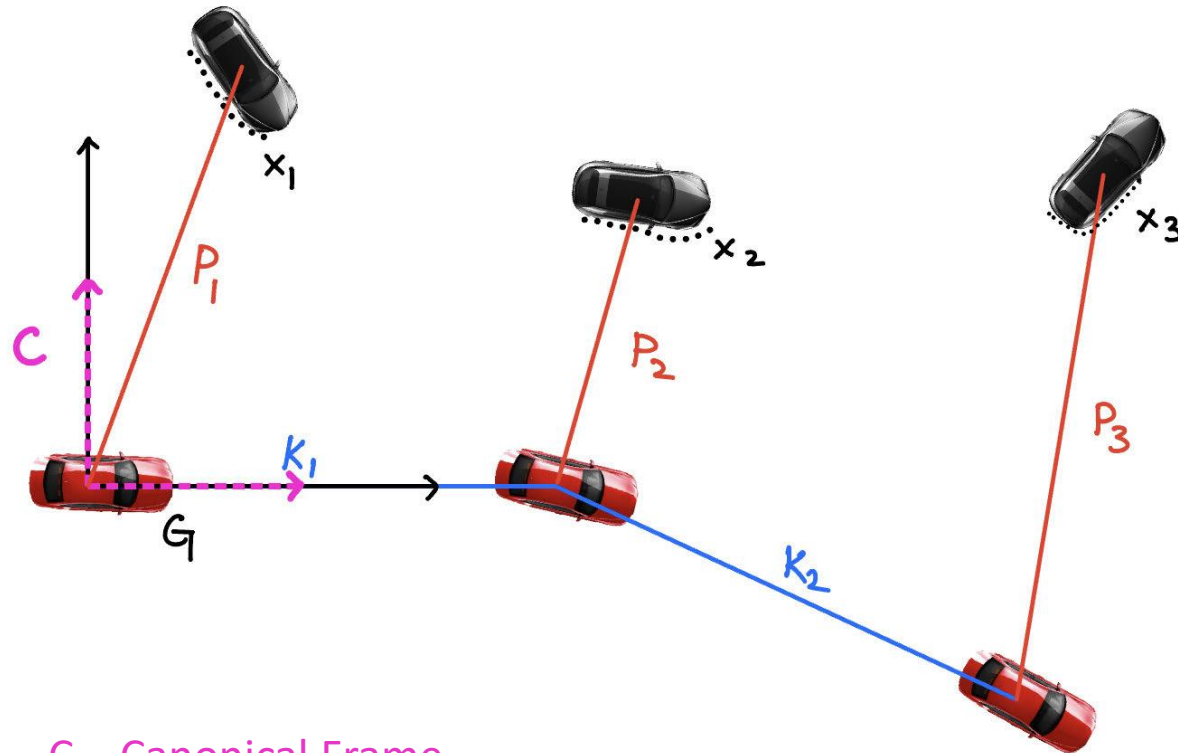
# Canonical Frame



SDF presumes that the input is in Canonical Frame



# Workflow: Transformation



$C$  – Canonical Frame

$G$  – Global Frame

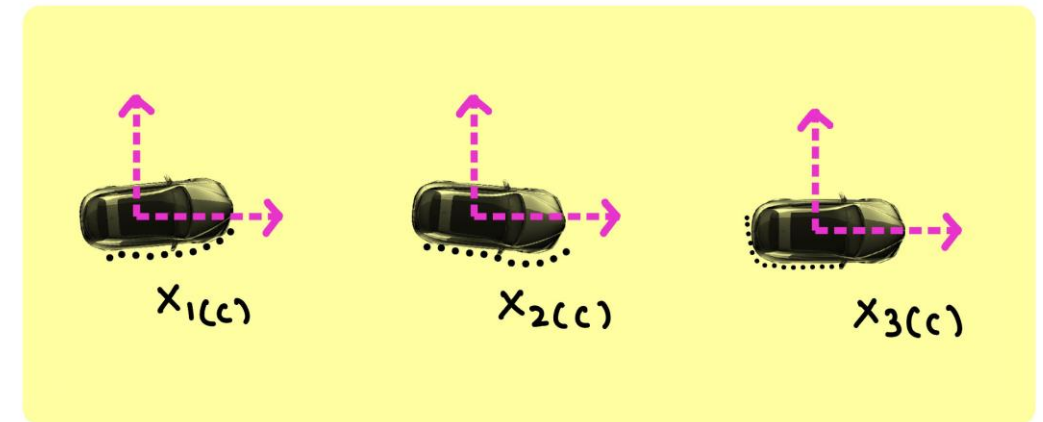
$P$  – Transformation from Point RCNN

$K$  – Transformation from KISS ICP

$$X_1[C] = P_1 X_1$$

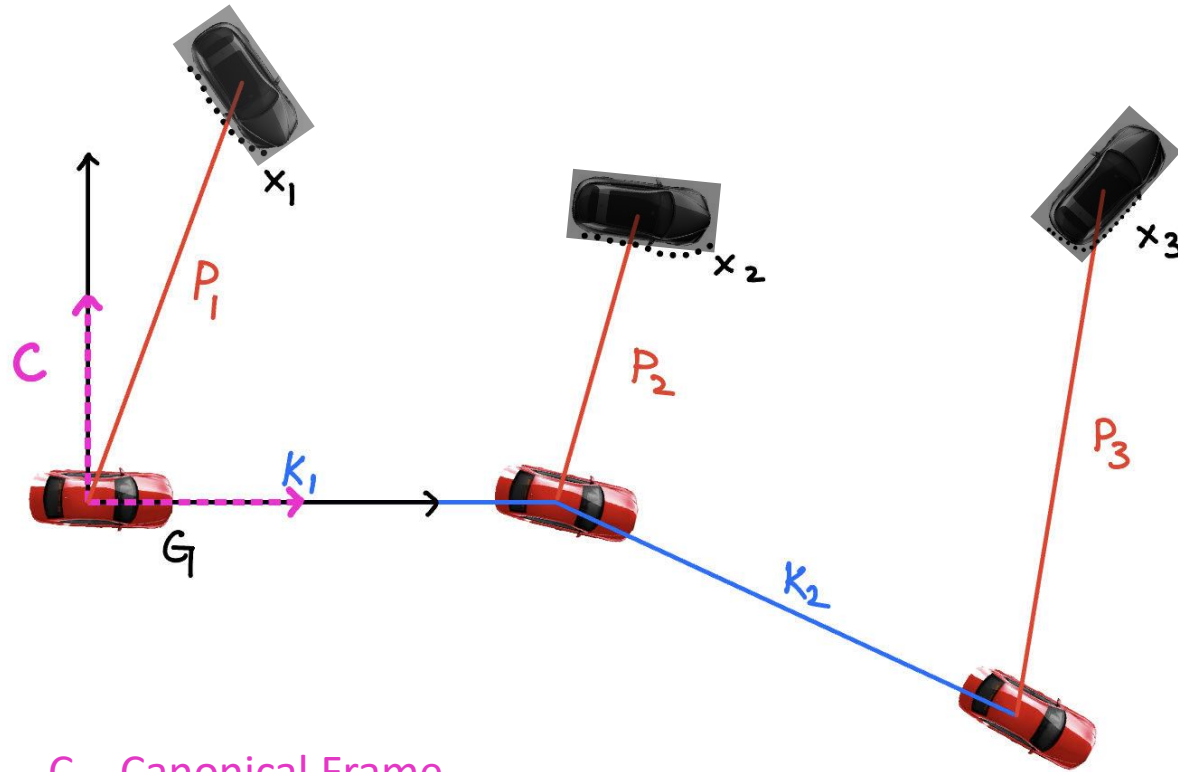
$$X_2[C] = K_1 P_2 X_2$$

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



Canonical Space

# Workflow: Transformation



C – Canonical Frame

G – Global Frame

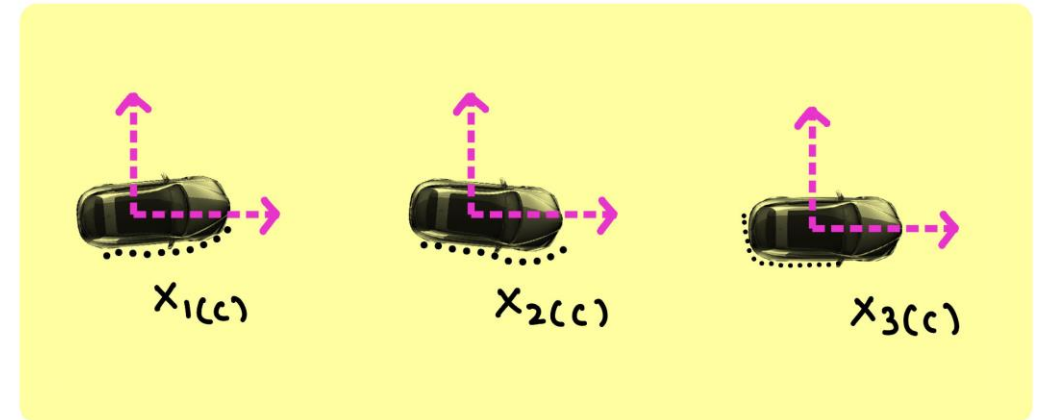
P – Transformation from Point RCNN

K – Transformation from KISS ICP

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

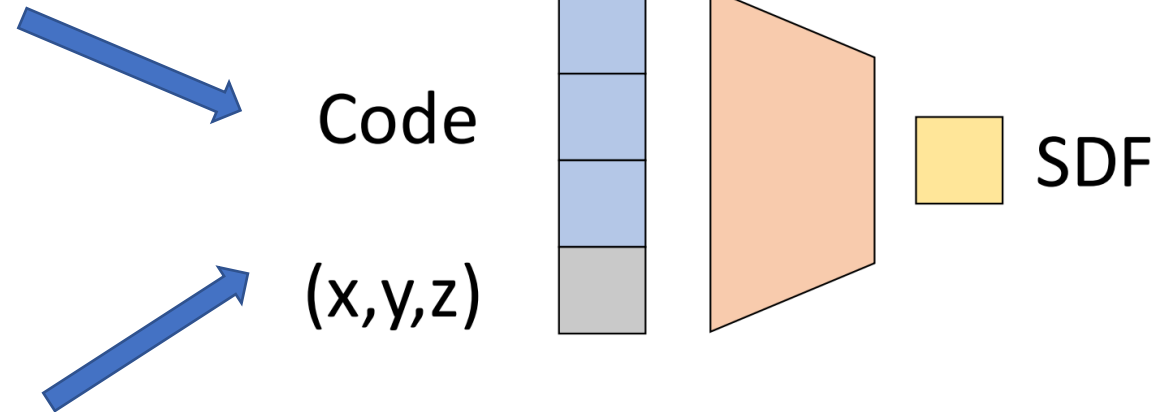
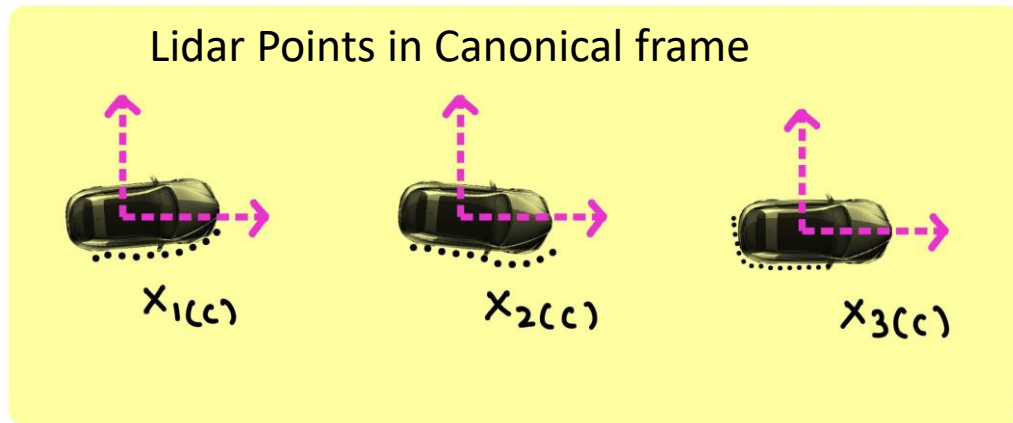
$$X_{2[C]} = K_1 P_2 X_2$$

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

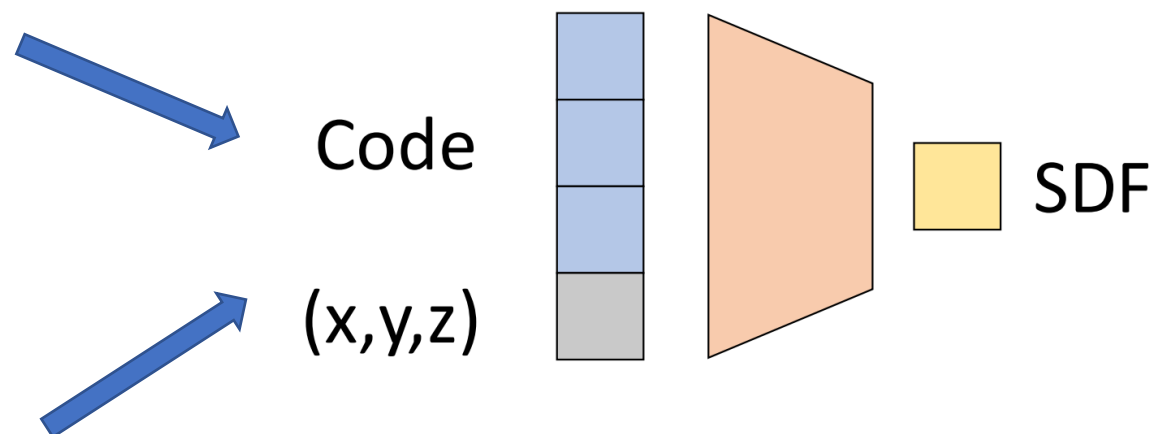
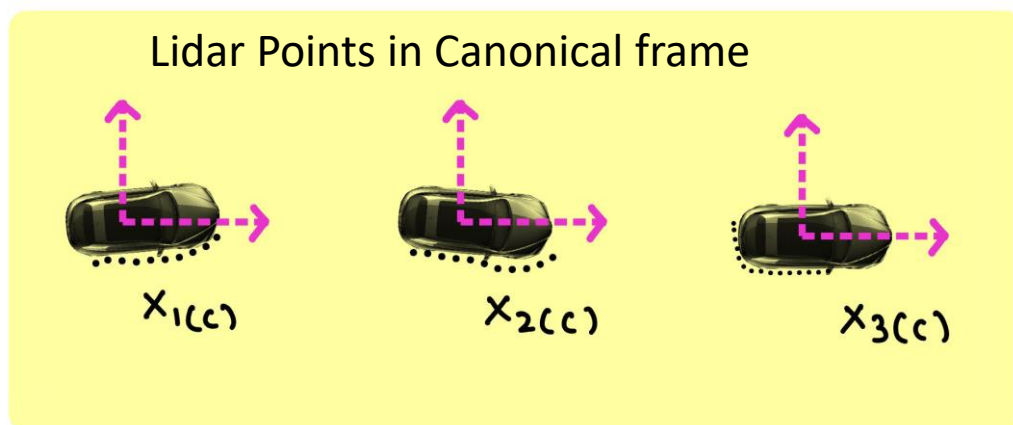


Canonical Space

[Latent Code Vector].T



[Latent Code Vector].T

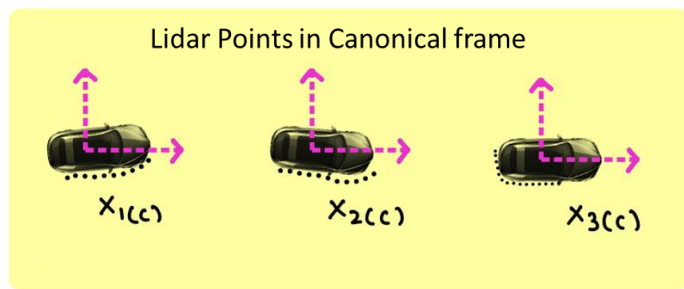


$$\begin{aligned} X_{1[C]} &= P_1 X_1 \\ X_{2[C]} &= K_1 P_2 X_2 \\ X_{3[C]} &= K_1 K_2 P_3 X_3 \end{aligned}$$



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

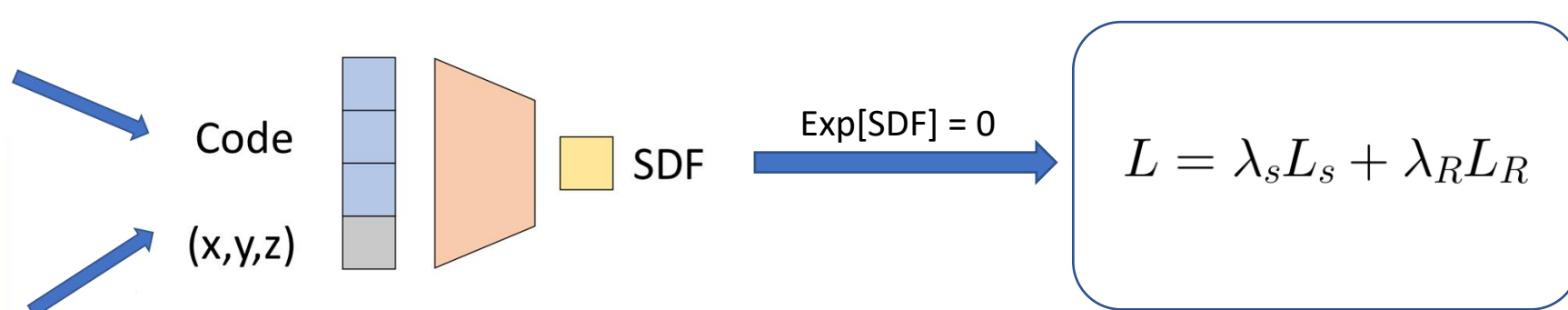
[Latent Code Vector].T



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

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

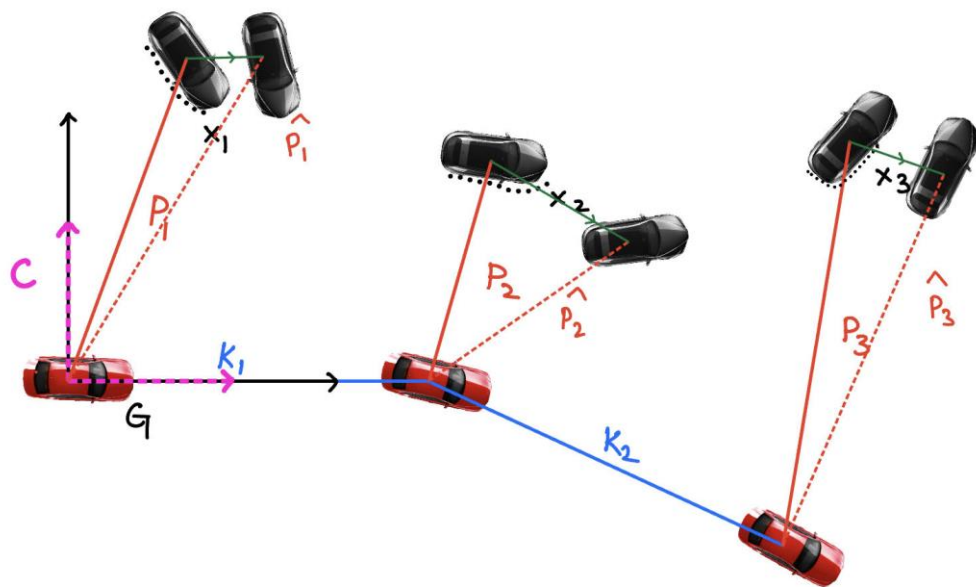
$$X_{3[C]} = K_1 K_2 P_3 X_3$$



1. Surface Consistency
2. Latent Code Regularization

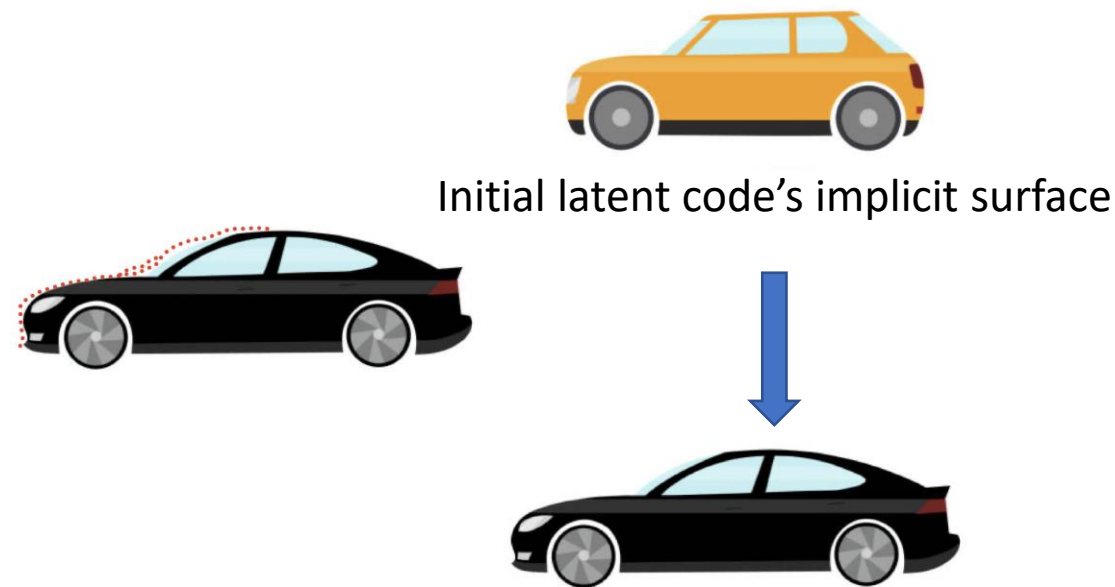


## Optimized Pose Estimation



Optimized poses of dynamic obstacles

## Optimized latent Code



Optimized latent code's implicit surface

# DataSets and Evaluation

## Dynamic Object Pose Estimation

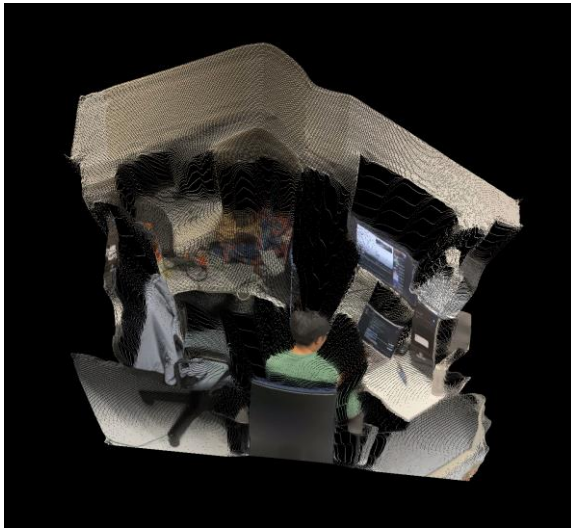




# 3D Reconstruction



Multiple Images



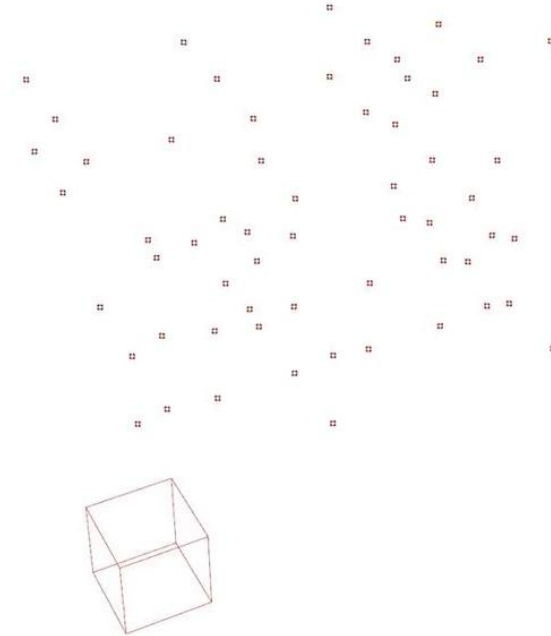
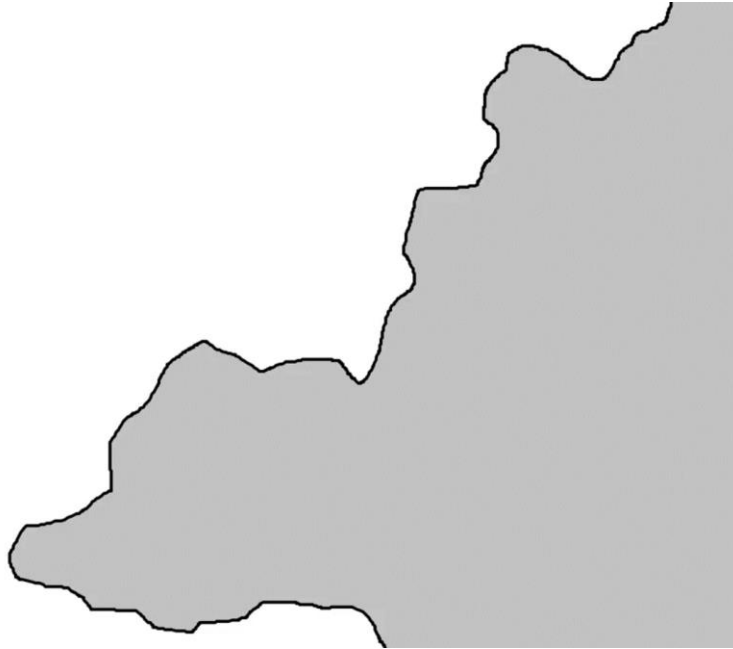
Partial Point Clouds



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).

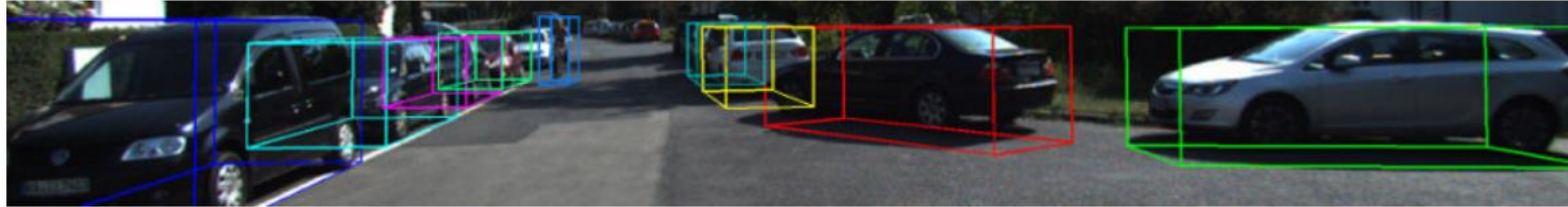
## M4: Make some contribution (~6 weeks)

- 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.

## 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.