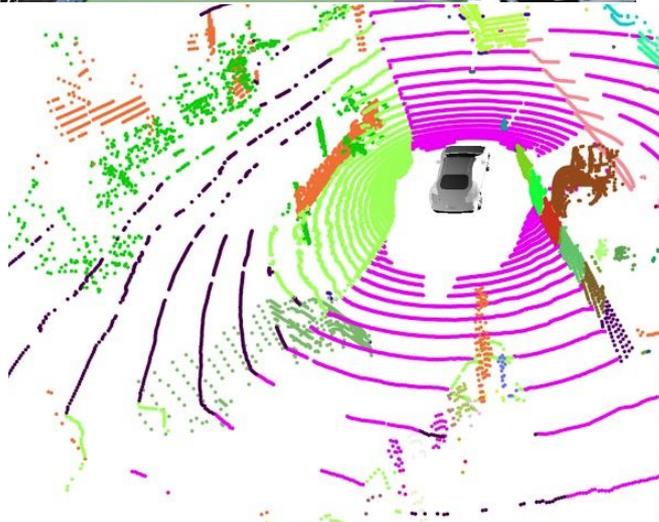


# Nuages de Points et Modélisation 3D

7 - Machine learning IV

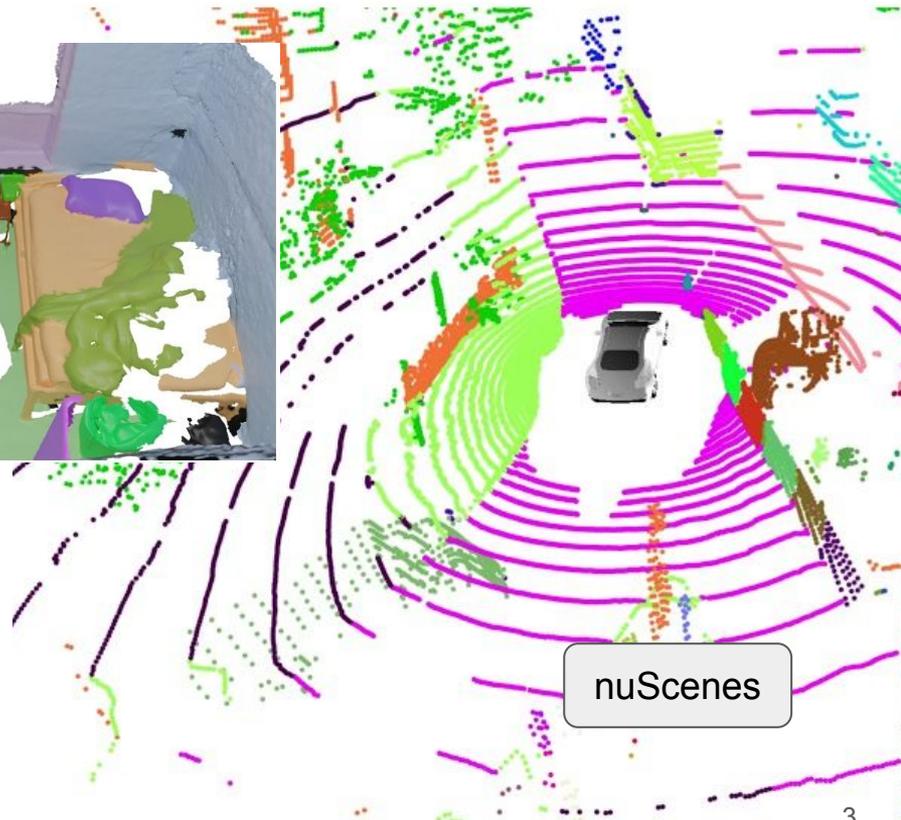
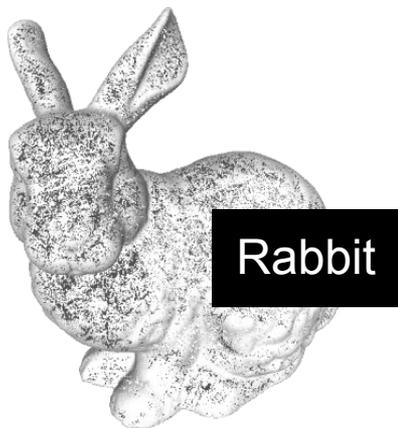
# Overview

- I. Tasks
- II. Self-supervised learning
  - A. Geometric reconstruction
  - B. Contrastive learning
  - C. Distillation
- III. Domain adaptation
- IV. OpenWorld



# Classification and semantic segmentation

## I - Tasks



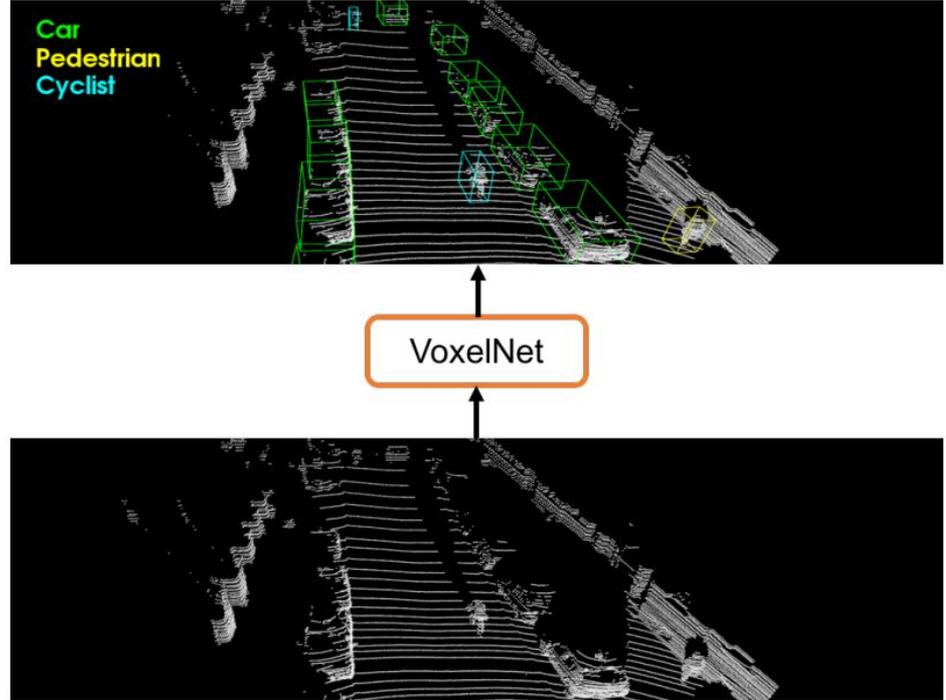
**Classification** → one label per point cloud

**Segmentation** → one label per point

# Detection

## I - Tasks

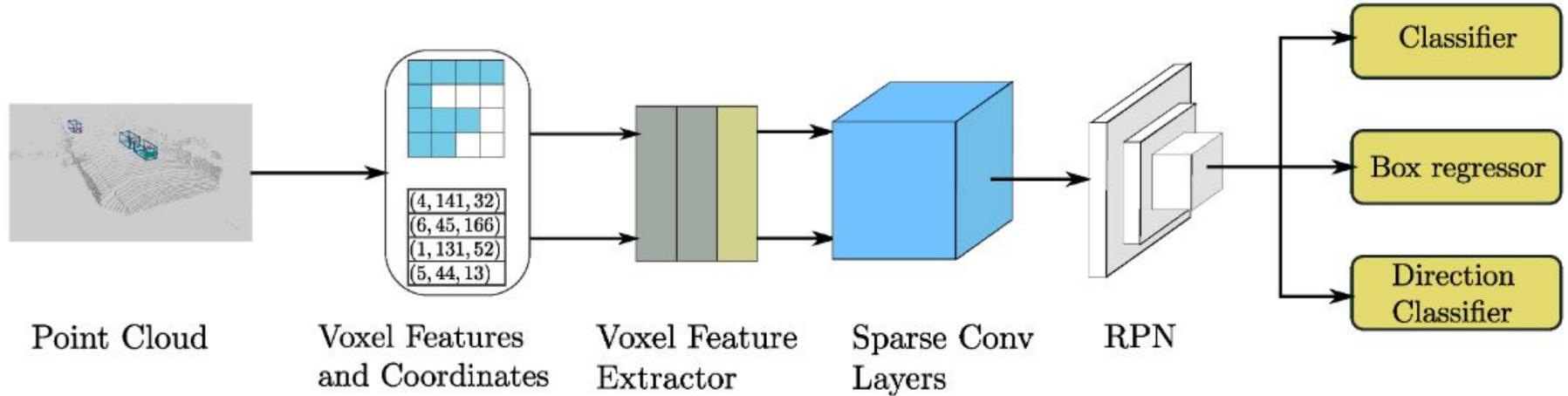
Put oriented boxes around objects of interest



# Detection

## I - Tasks

## Second

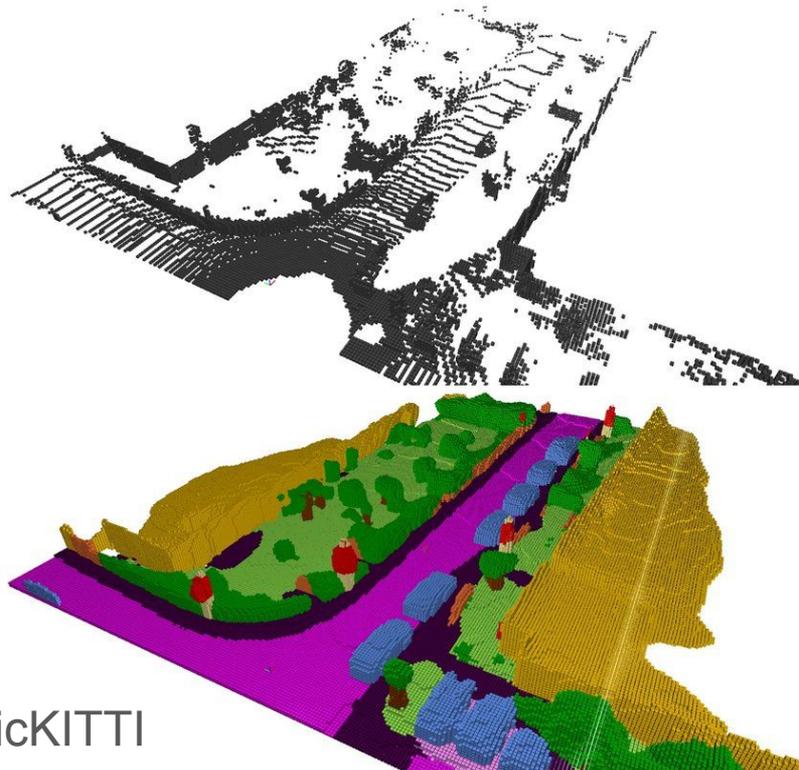


# Scene completion

I - Tasks

**Input:** lidar scan

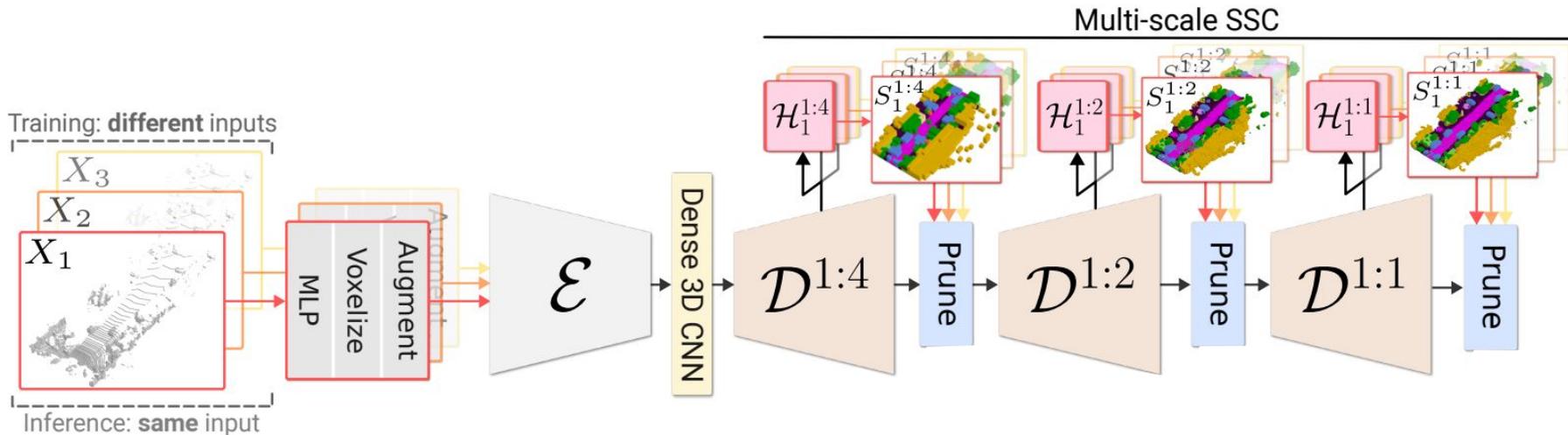
**Output:** completed voxel scene with semantics



SemanticKITTI

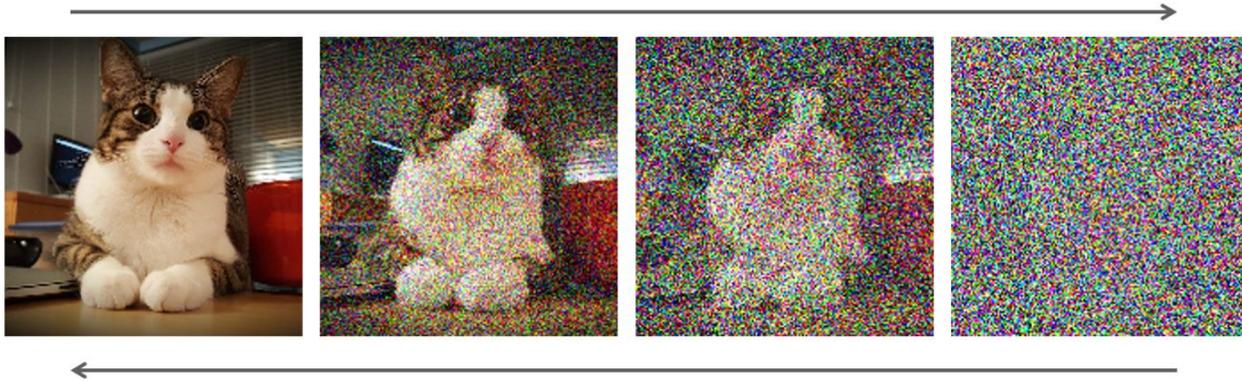
# Scene completion

## I - Tasks



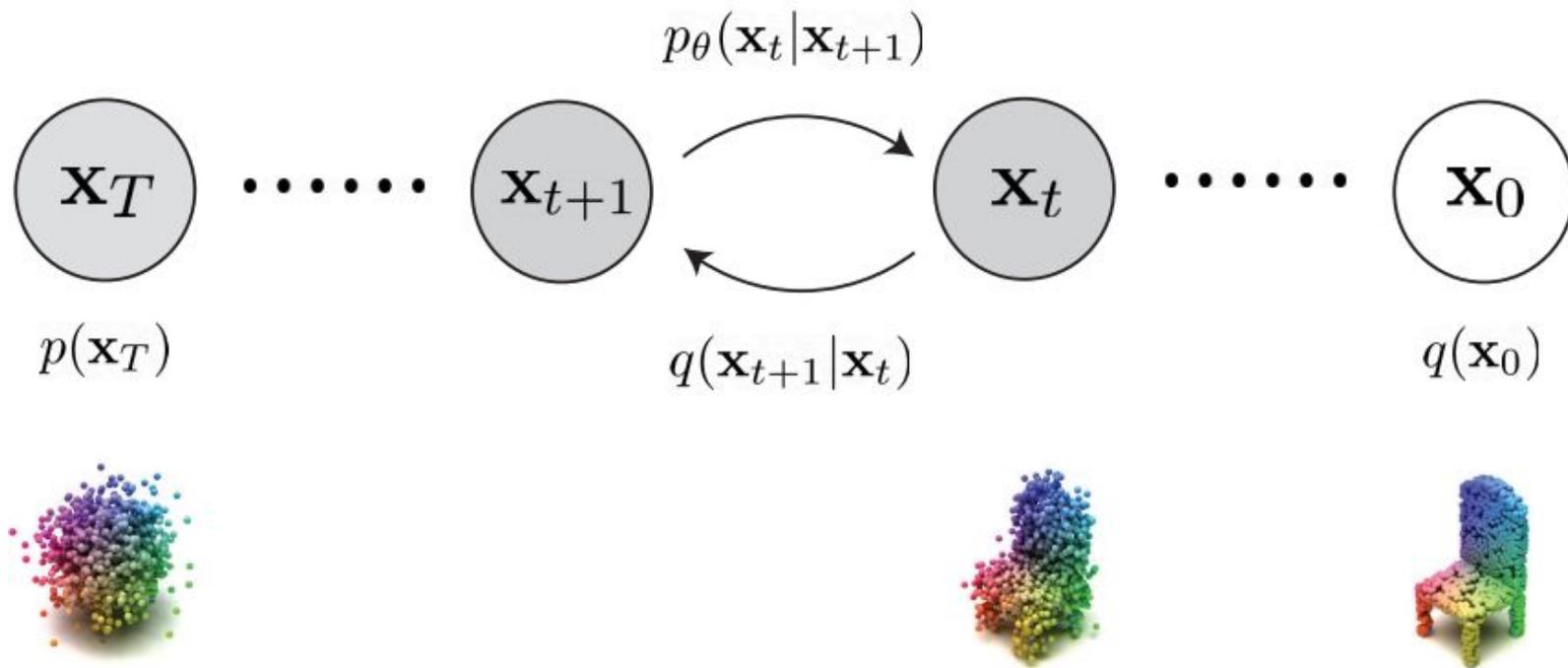
# Generation (e.g. diffusion models)

## I - Tasks



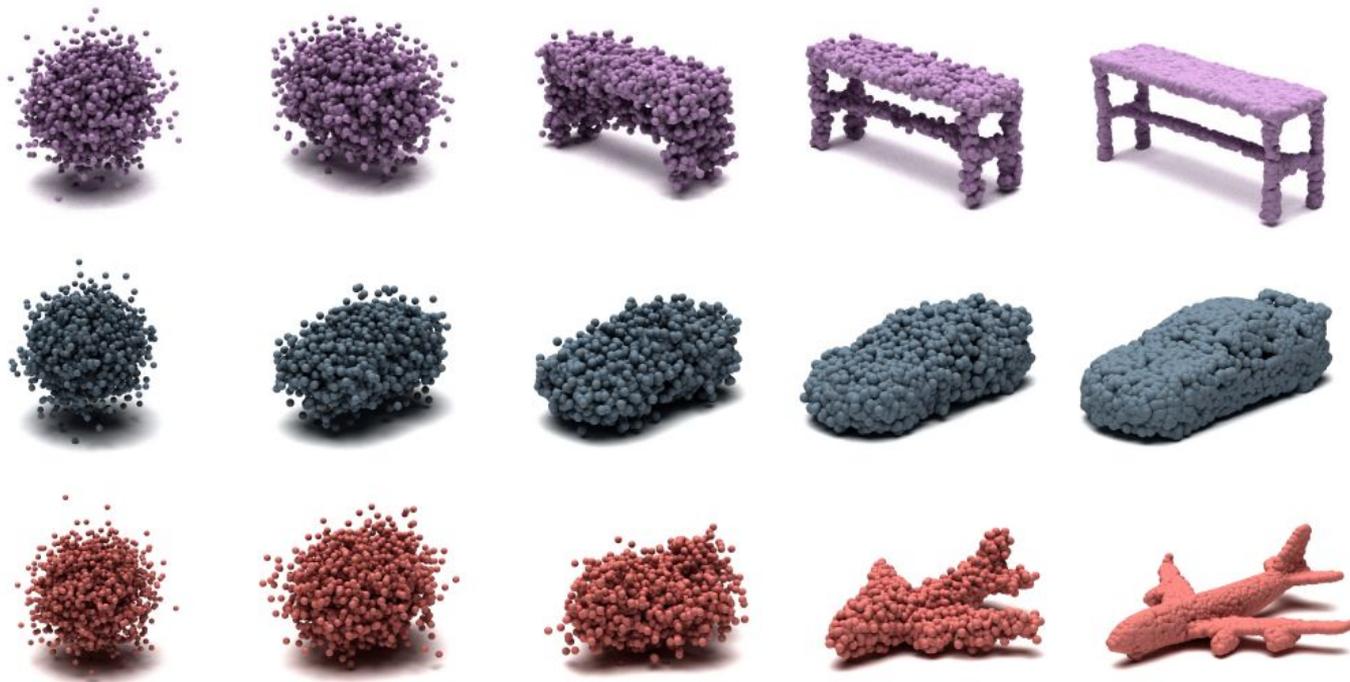
# Generation (e.g. diffusion models)

## I - Tasks



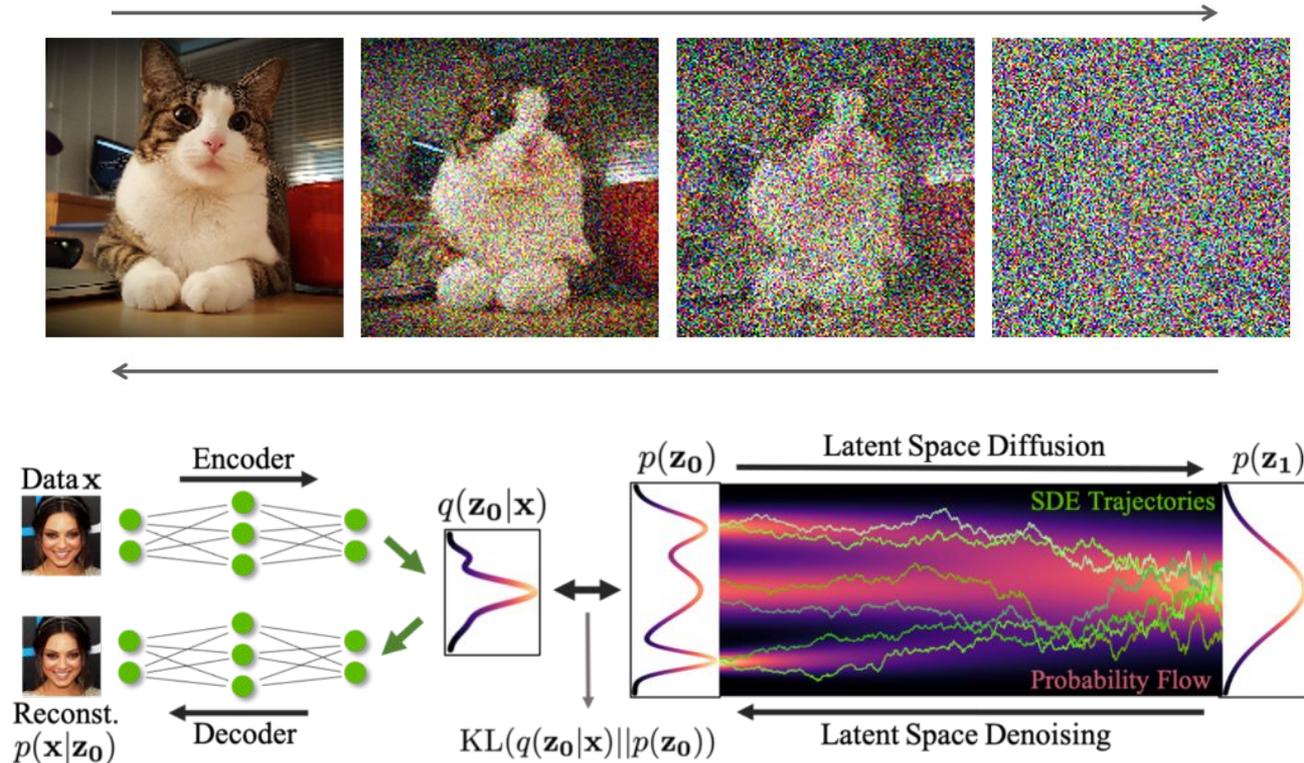
# Generation (e.g. diffusion models)

## I - Tasks



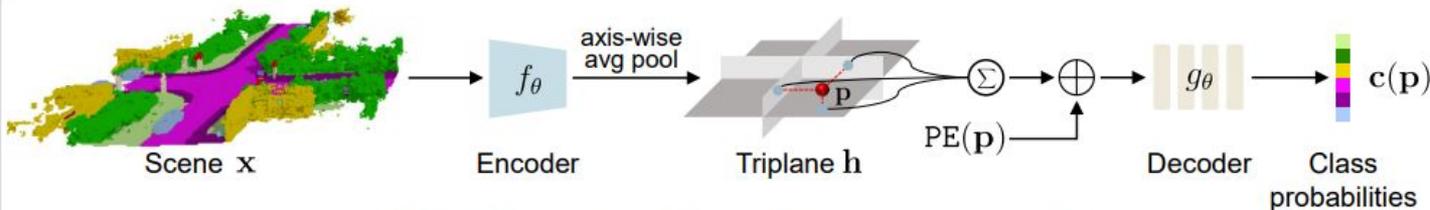
# Generation (e.g. diffusion models)

## I - Tasks

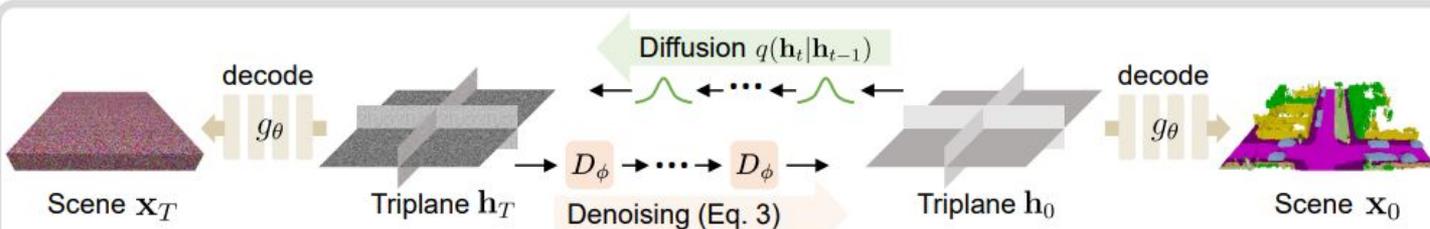


# Generation (e.g. diffusion models)

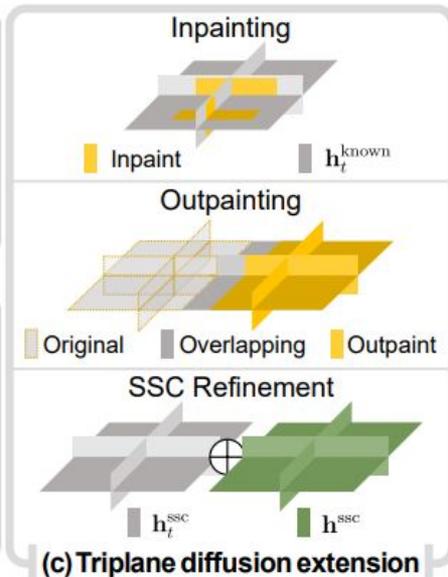
## I - Tasks



(a) Triplane learning for efficient outdoor scene compression



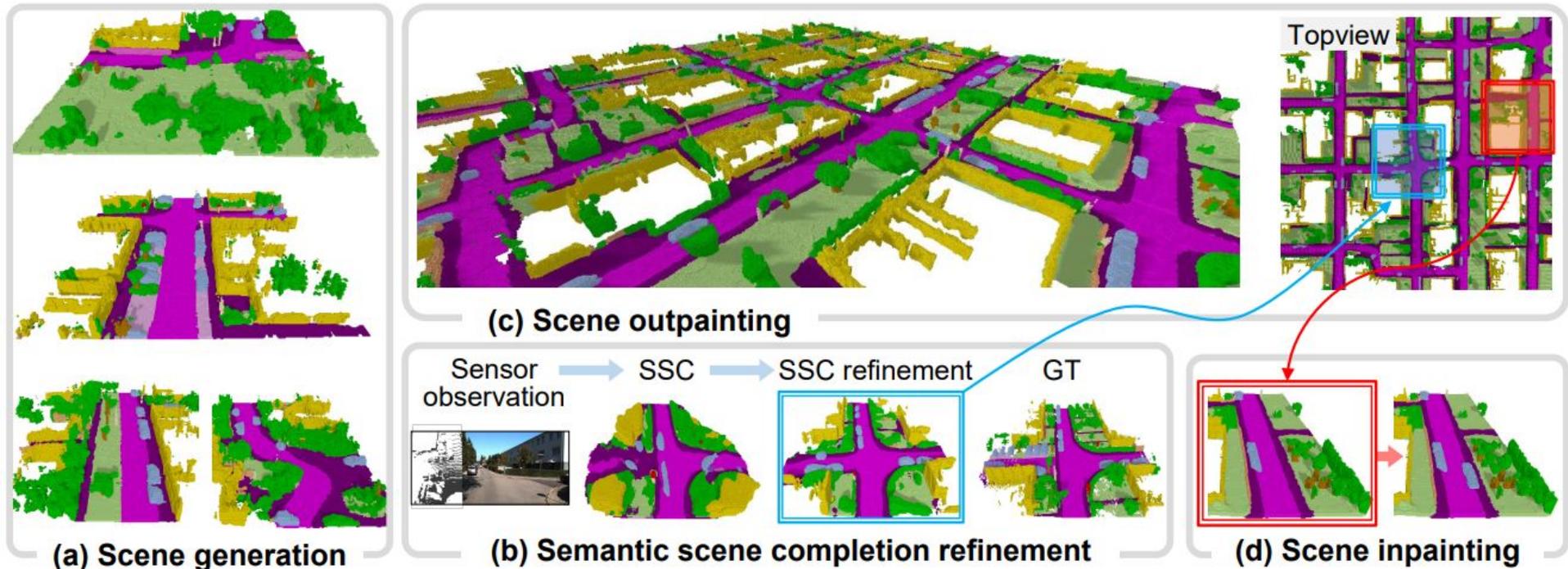
(b) Triplane diffusion for outdoor scene generation



(c) Triplane diffusion extension

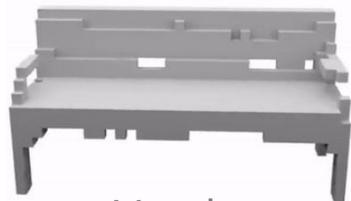
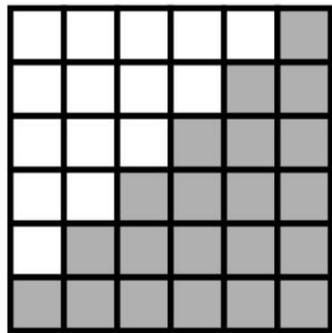
# Generation (e.g. diffusion models)

## I - Tasks

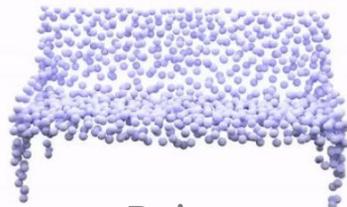
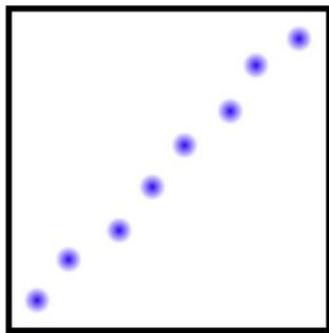


# Surface reconstruction

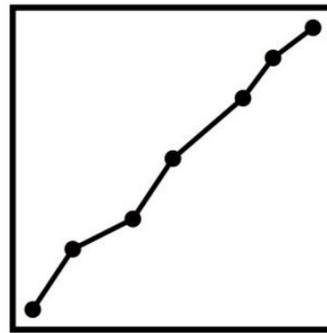
## I - Tasks



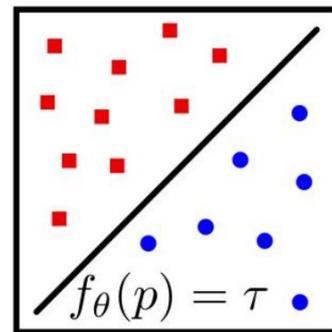
Voxels



Points



Mesh



Implicit

**Occupancy Networks: Learning 3D Reconstruction in Function Space**

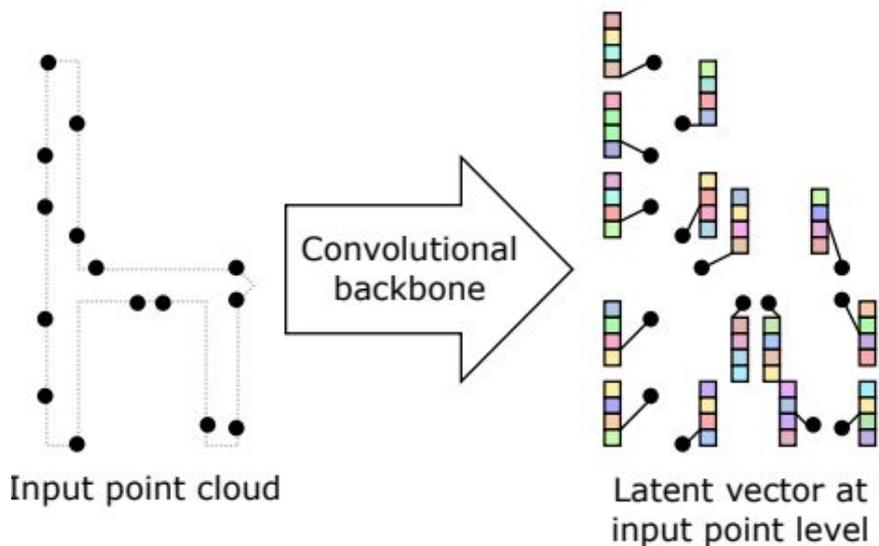
*Mescheder, Lars and Oechsle, Michael and Niemeyer, Michael and Nowozin, Sebastian and Geiger, Andreas*

Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR), 2019

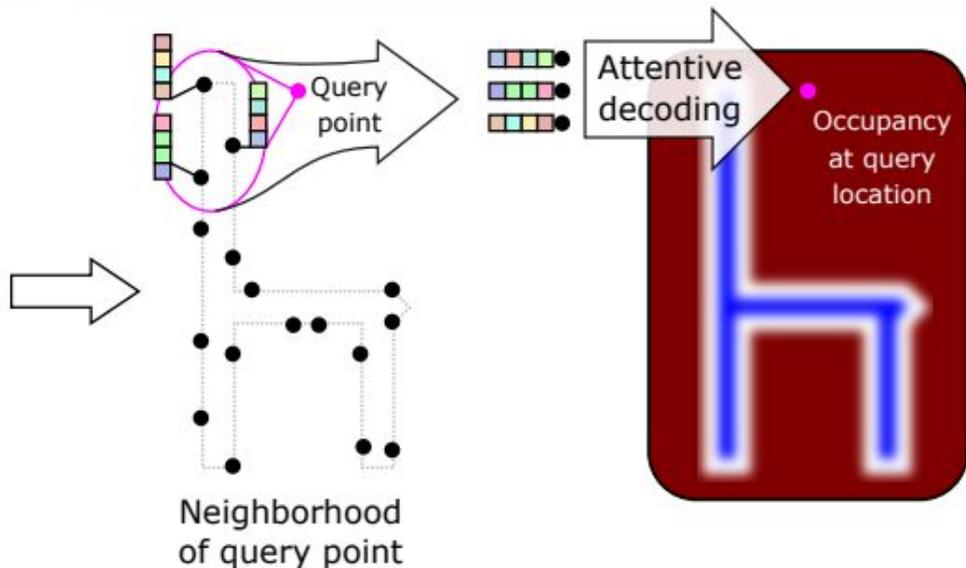
# Surface reconstruction

## I - Tasks

### Shape encoding at point level

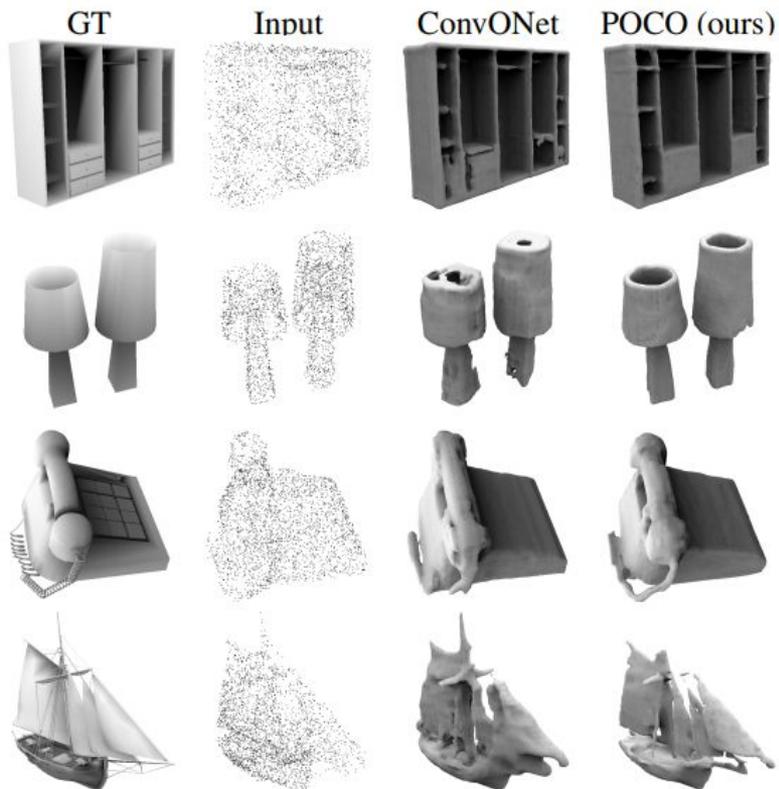


### Local decoding

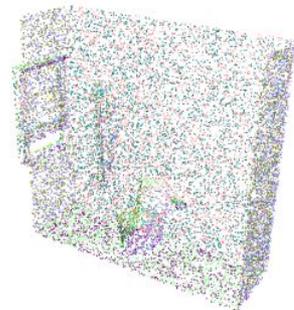


# Surface reconstruction

## I - Tasks

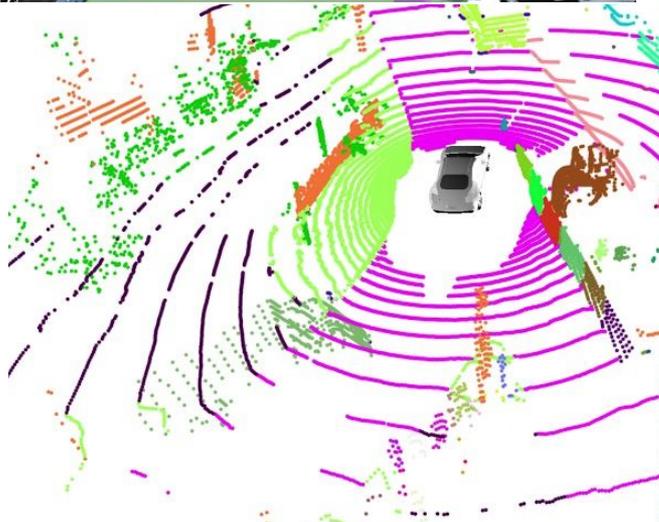


500 pts/m<sup>2</sup>



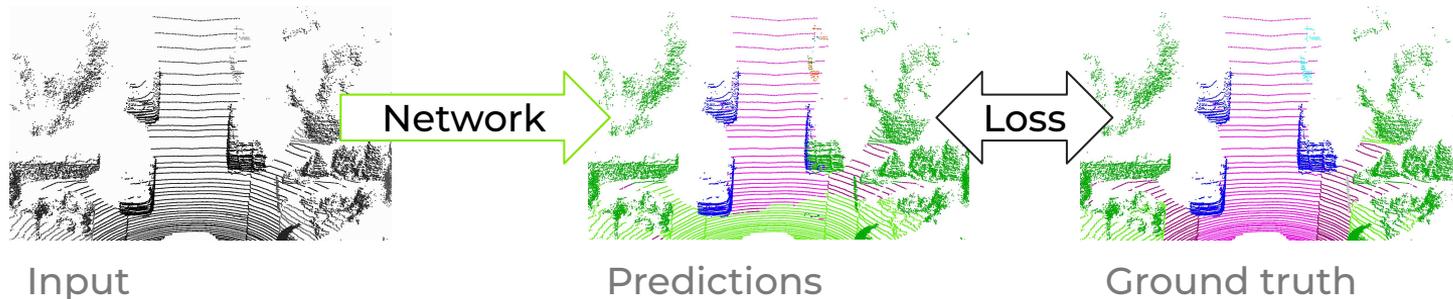
# Overview

- I. Tasks
- II. Self-supervised learning
  - A. Geometric reconstruction
  - B. Contrastive learning
  - C. Distillation
- III. Domain adaptation
- IV. OpenWorld



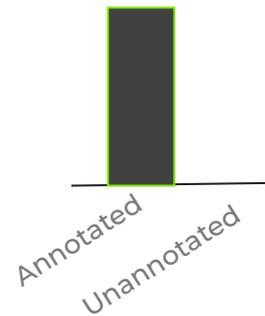
# Supervised learning

## Self-supervised learning



## Train with annotated data

- Annotations are costly
- Limitation of the database size

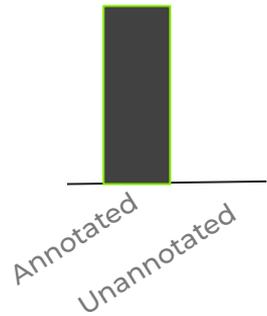


# Frugal learning

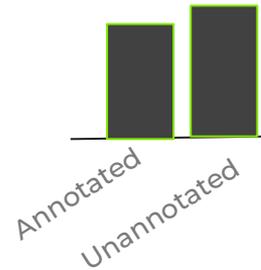
Self-supervised learning

Objective: learning with less annotations

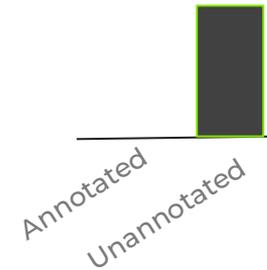
Supervised



Semi -  
Supervised



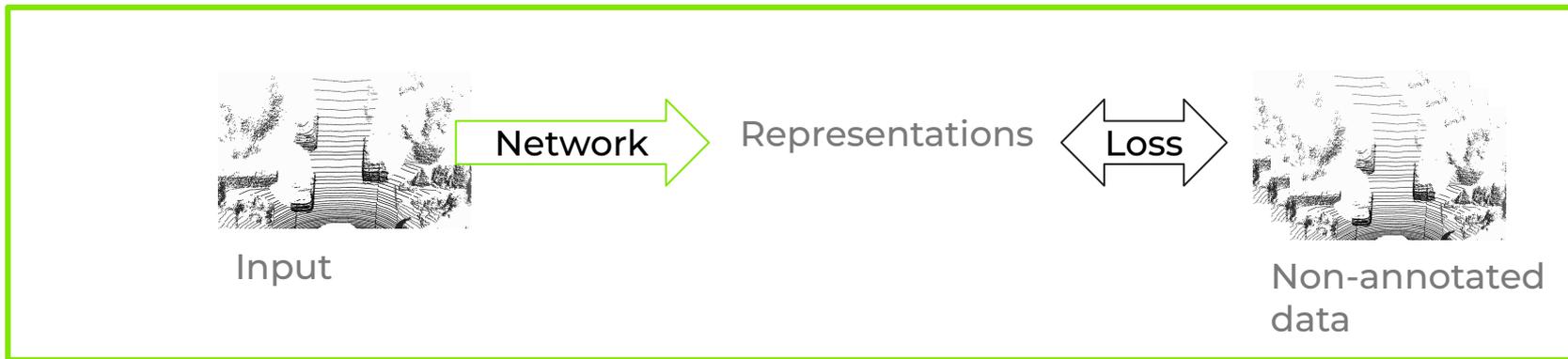
Un / self - Supervised



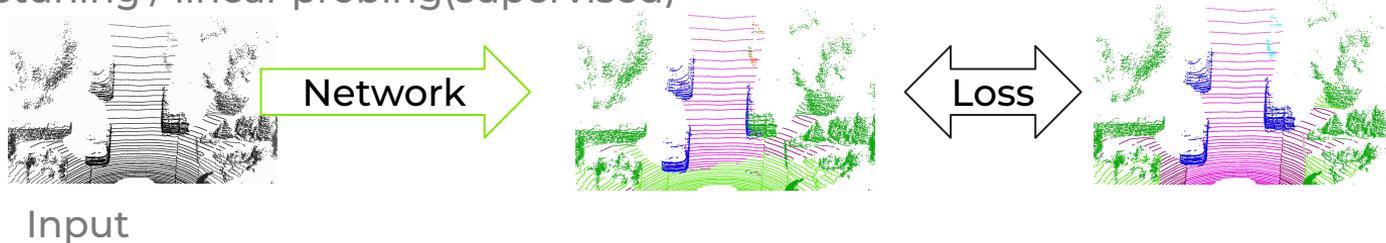
# Self-supervision

**What?** learn useful representations without annotations

**Why?** better performance when finetuning / data-efficiency



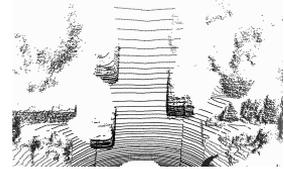
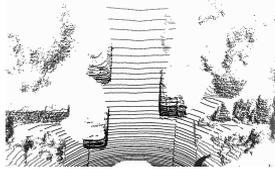
Step 2: finetuning / linear probing(supervised)



# Self-supervision

## Step 1

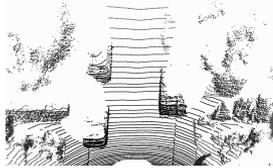
Pretraining



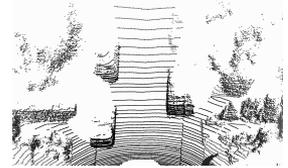
# Self-supervision

## Step 1

Pretraining

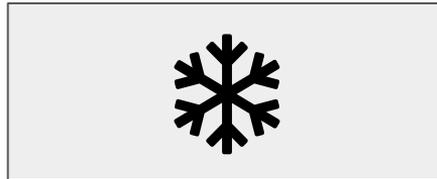
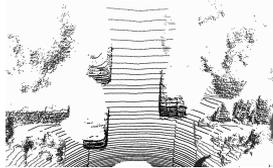


Self-sup.  
head

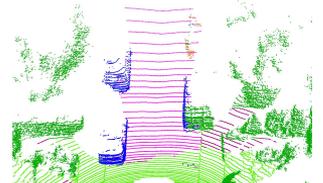


## Step 2

Probing



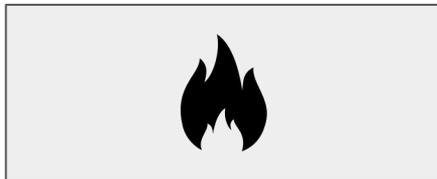
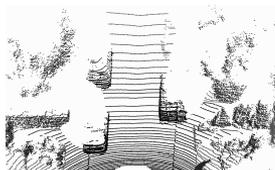
Taks  
head



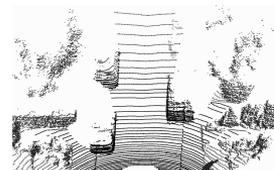
# Self-supervision

## Step 1

Pretraining

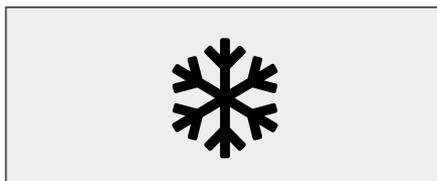
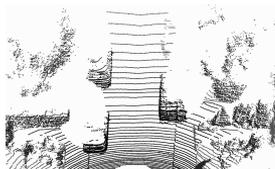


Self-sup.  
head

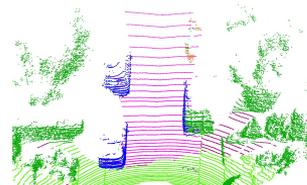


## Step 2

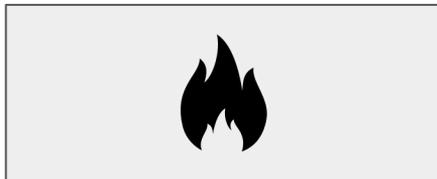
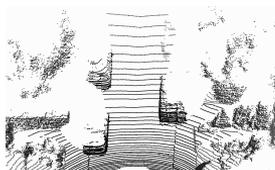
Probing



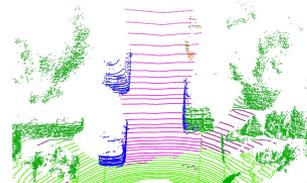
Task  
head



Finetuning

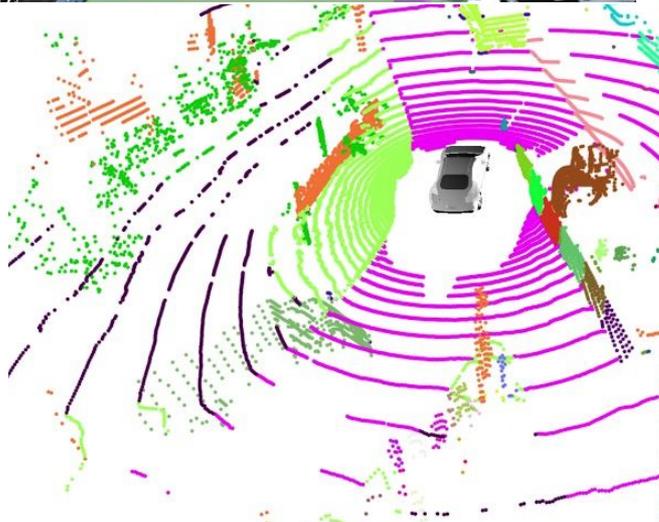


Task  
head



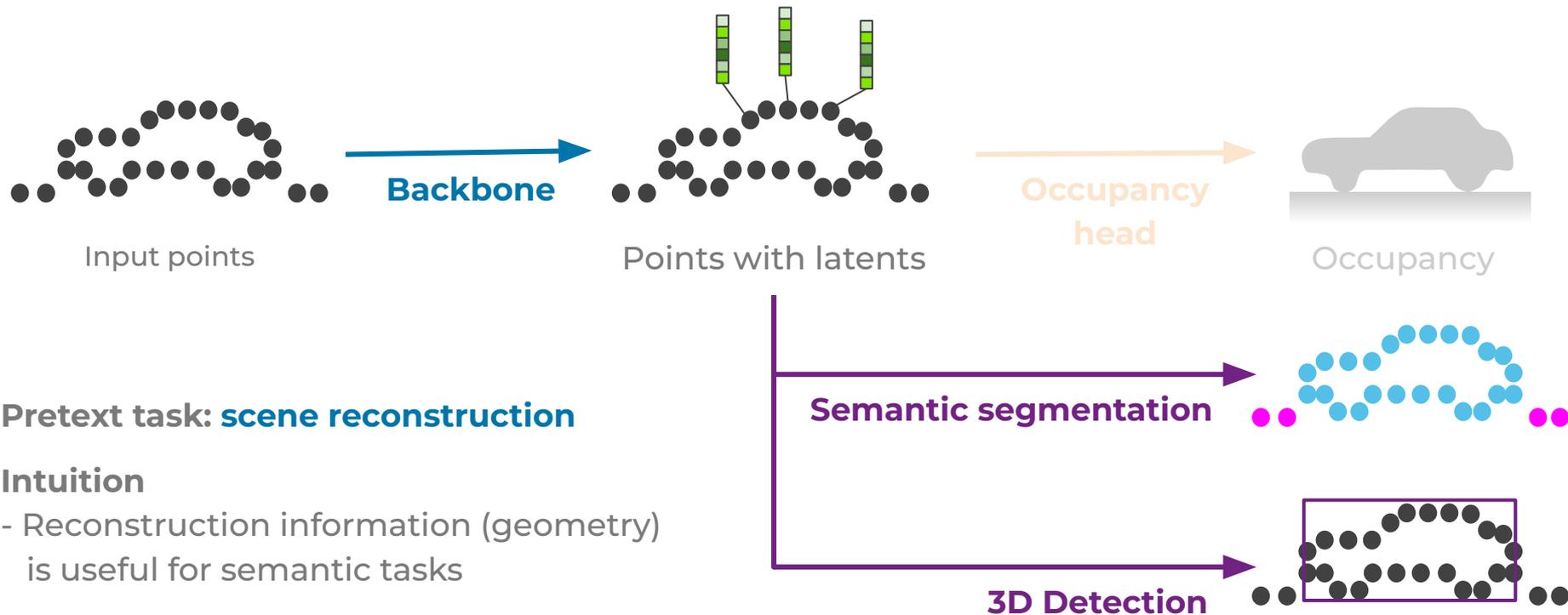
# Overview

- I. Tasks
- II. Self-supervised learning
  - A. Geometric reconstruction
  - B. Contrastive learning
  - C. Distillation
- III. Domain adaptation
- IV. OpenWorld



# ALSO: Automotive Lidar Self-Supervision by Occupancy Estimation

ALSO



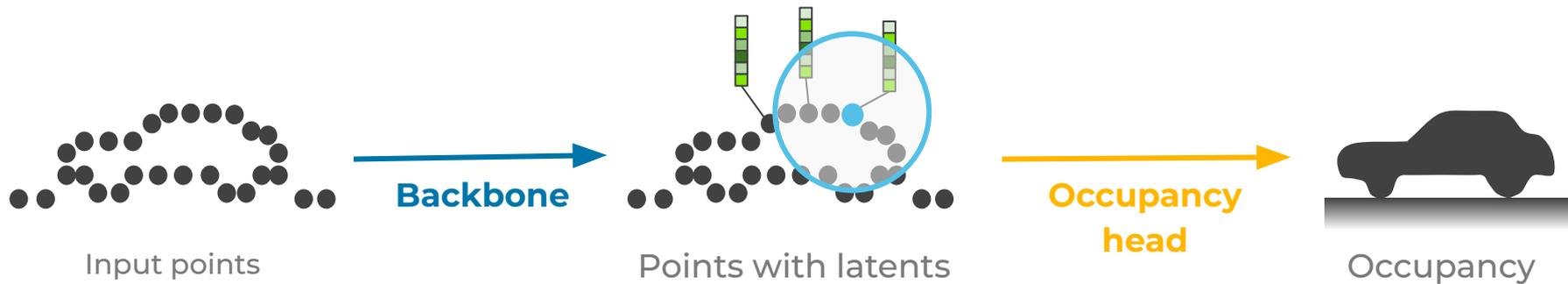
Pretext task: **scene reconstruction**

**Intuition**

- Reconstruction information (geometry) is useful for semantic tasks

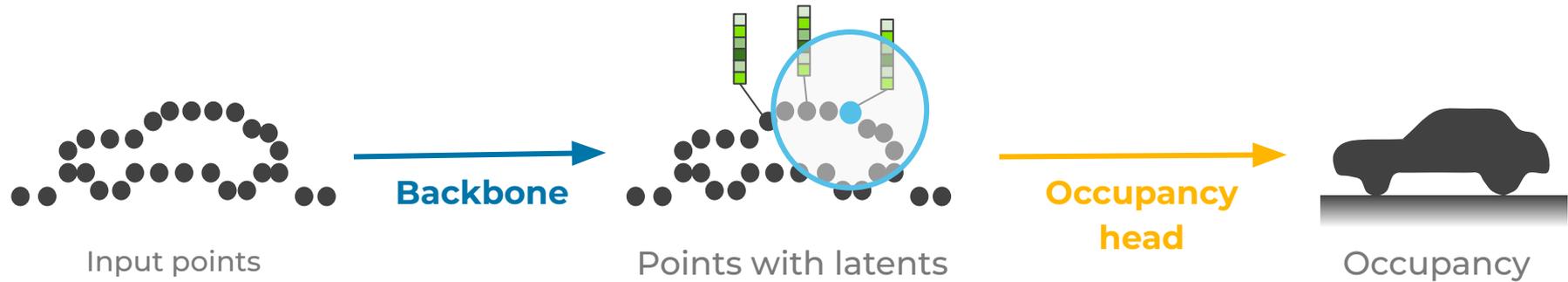
# Context reconstruction vs local reconstruction

ALSO



# Context reconstruction vs local reconstruction

ALSO



## Local reconstruction [POCO head]

- everywhere, from features of neighboring points
- ⇒ (too) detailed geometry

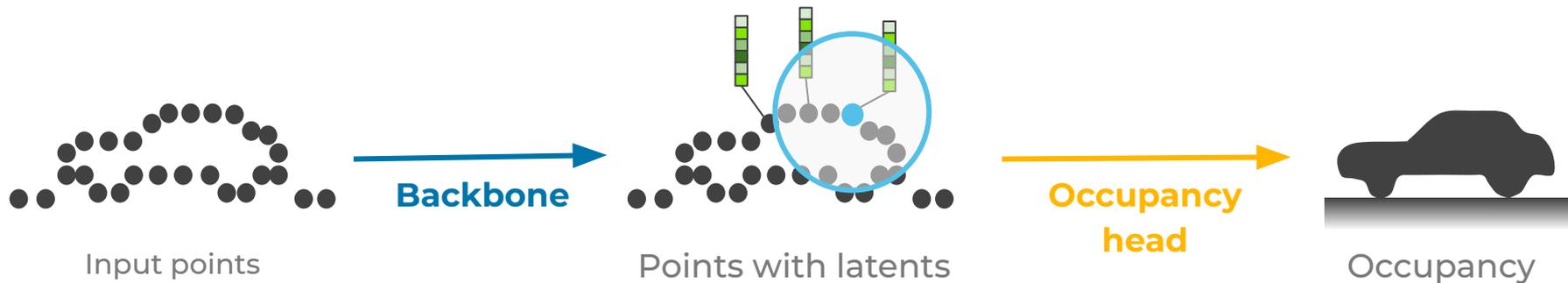
## Context reconstruction [ALSO head]

- of a 1 meter ball, from each single feature point
- ⇒ rough geometry, more suited for object recognition

*POCO: Point Convolution for Surface Reconstruction, A. Boulch, R. Marlet, CVPR 2022*

# Context reconstruction vs local reconstruction

ALSO

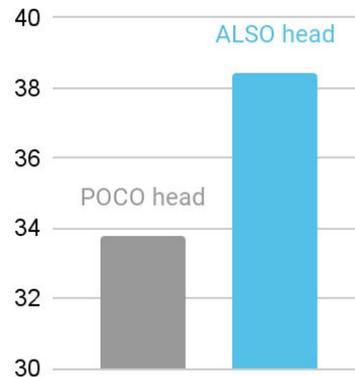


## Local reconstruction [POCO head]

- everywhere, from features of neighboring points
- ⇒ (too) detailed geometry

## Context reconstruction [ALSO head]

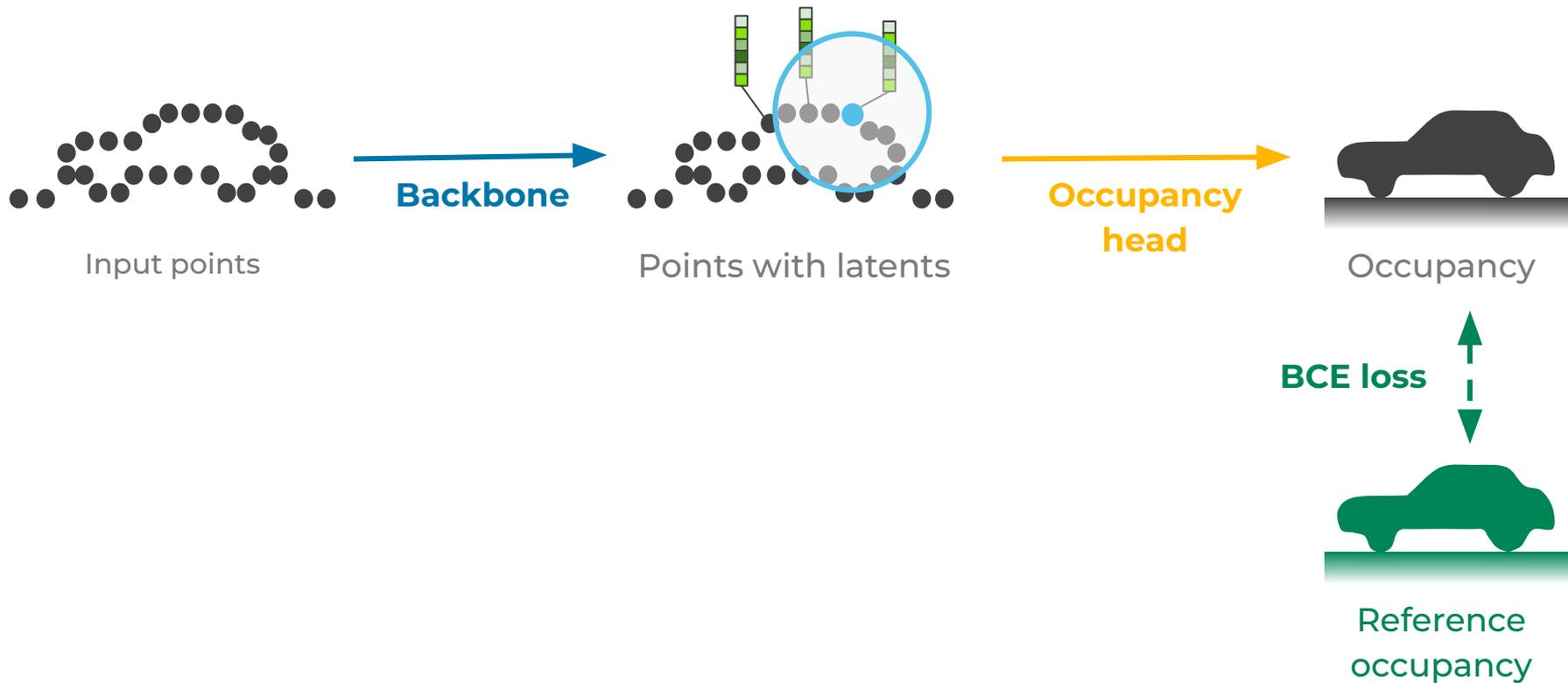
- of a 1 meter ball, from each single feature point
- ⇒ rough geometry, more suited for object recognition



POCO: Point Convolution for Surface Reconstruction, A. Boulch, R. Marlet, CVPR 2022

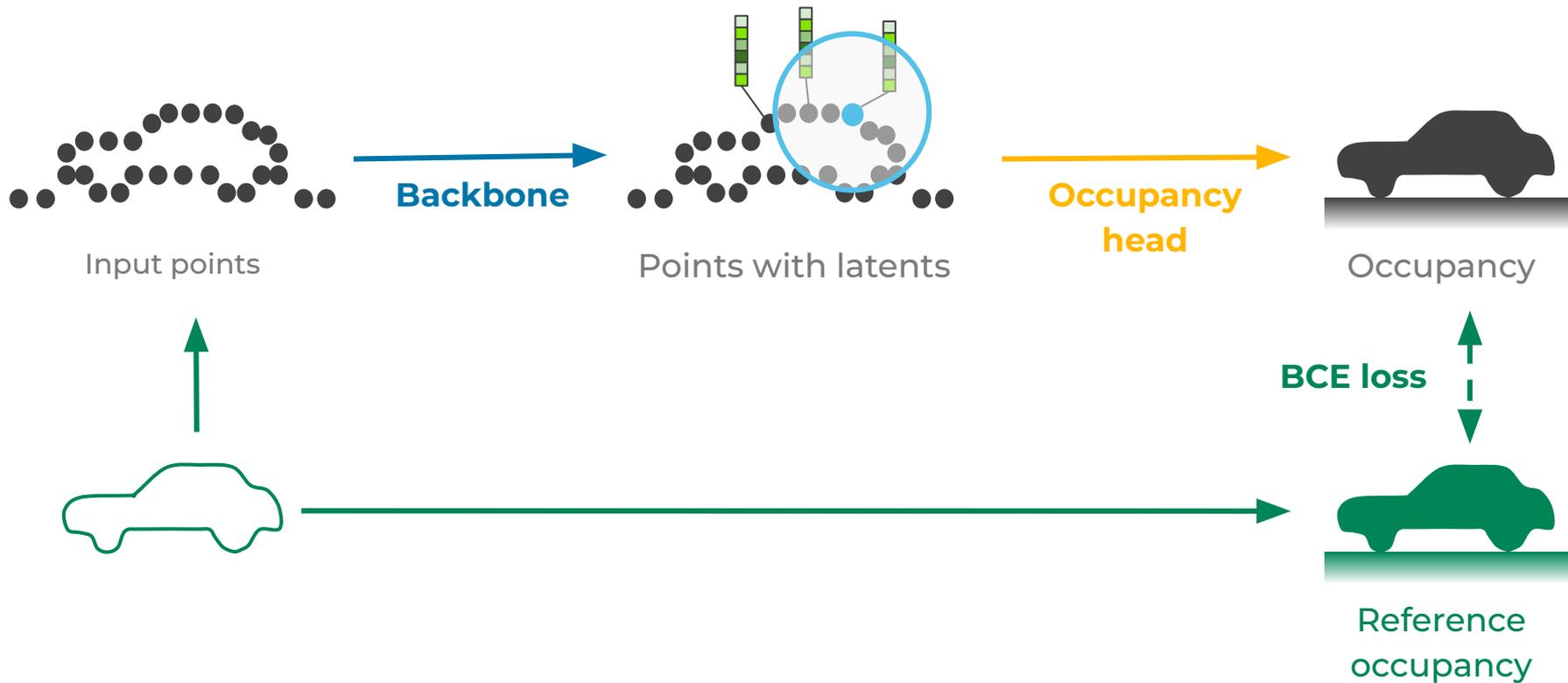
# Supervision

ALSO



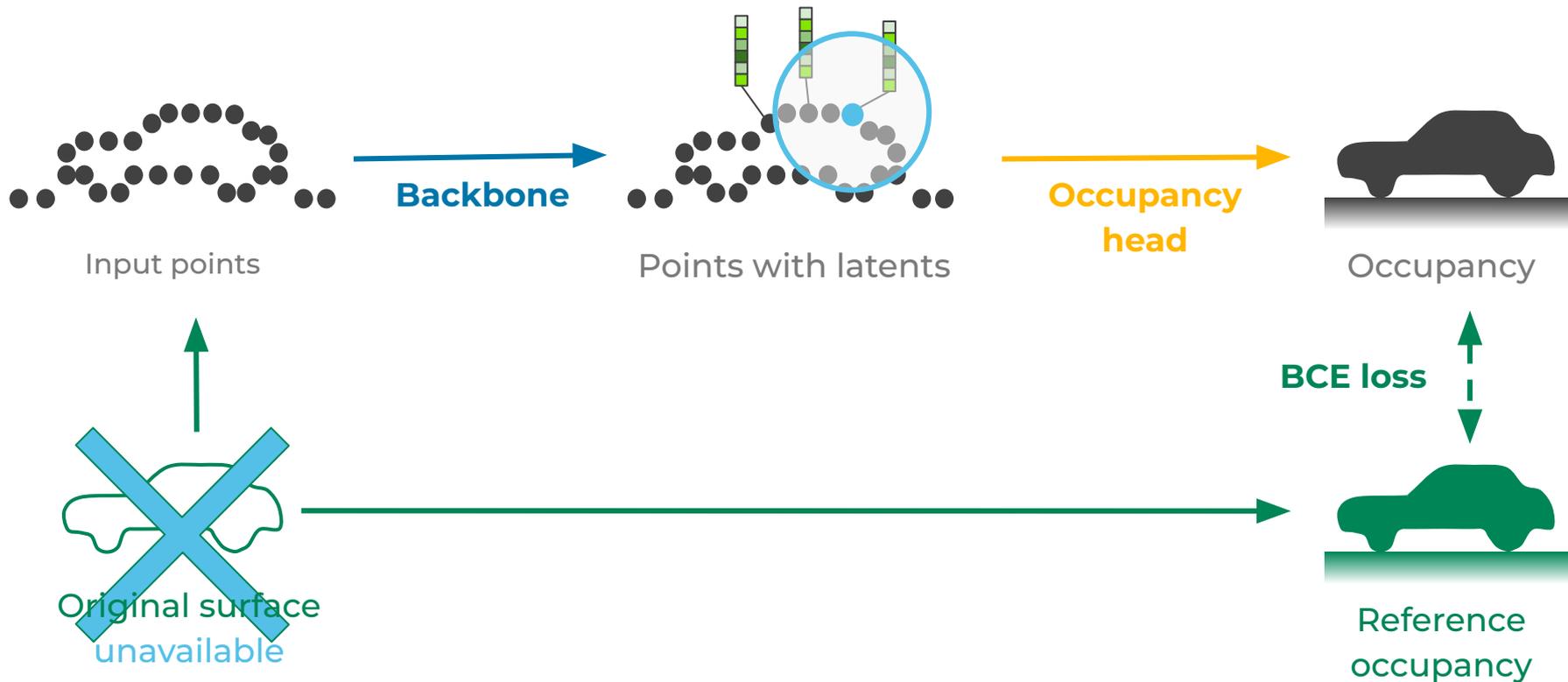
# Supervision

ALSO



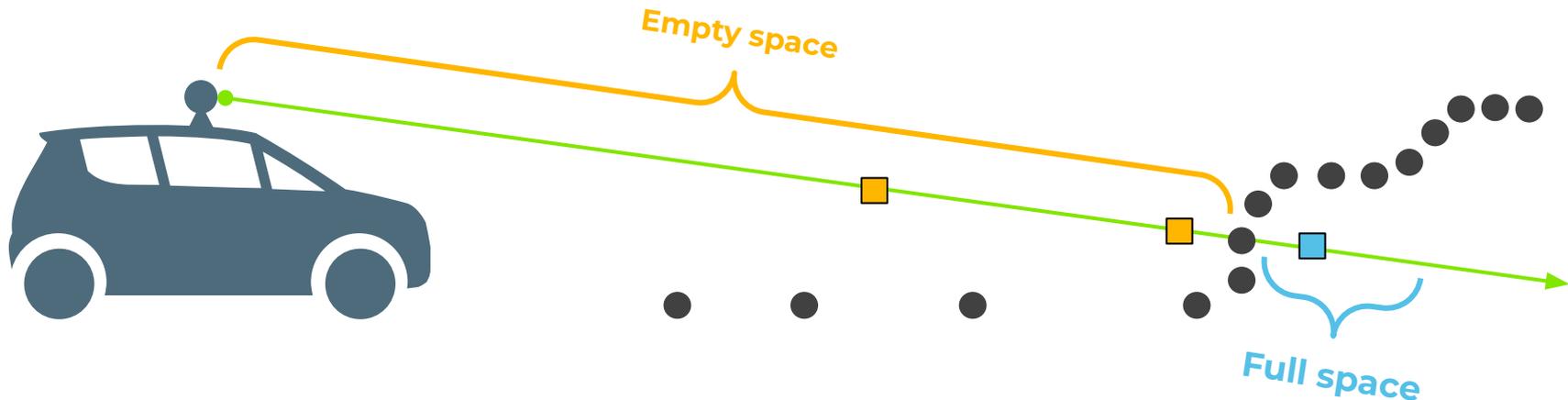
# Supervision

ALSO



# Self-supervised occupancy - Query point generation

ALSO

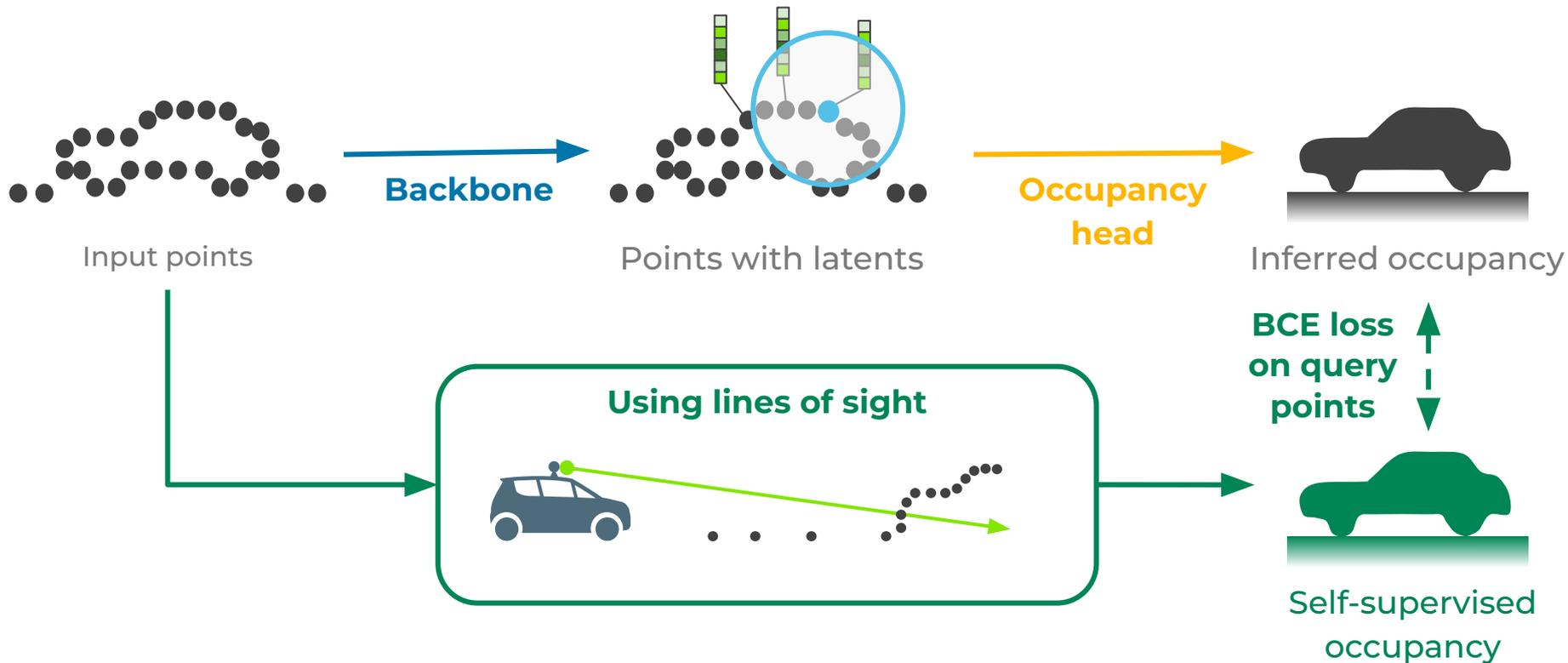


Along lidar lines of sight

- **Empty queries:** from sensor to observed point
- **Full queries:** just behind the point (max distance  $\delta = 0.1$  m)

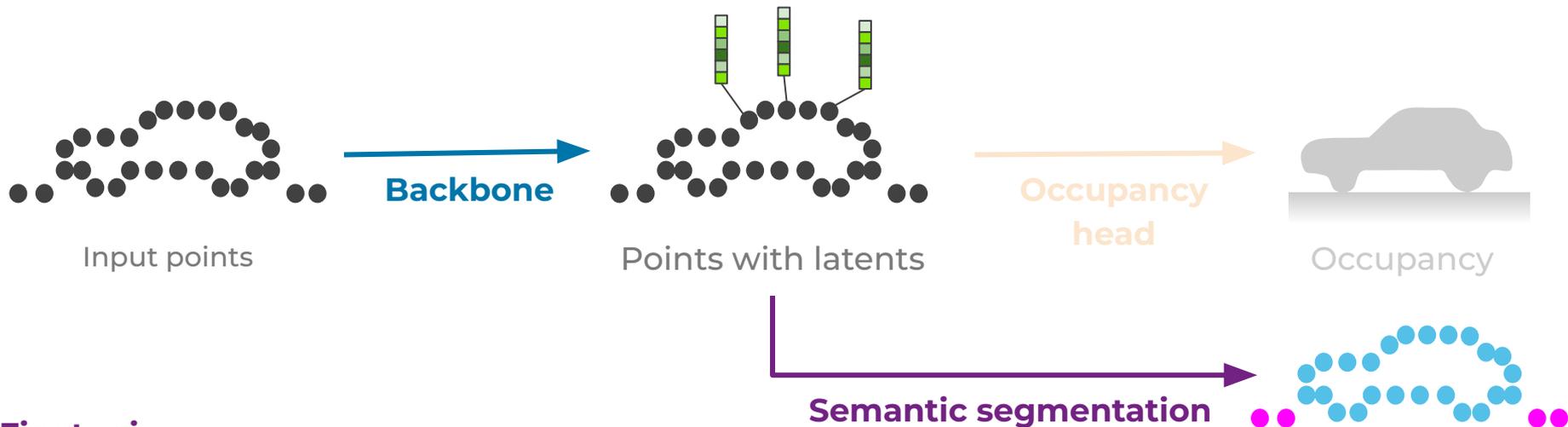
# Self-supervision

ALSO



## Downstream tasks

ALSO

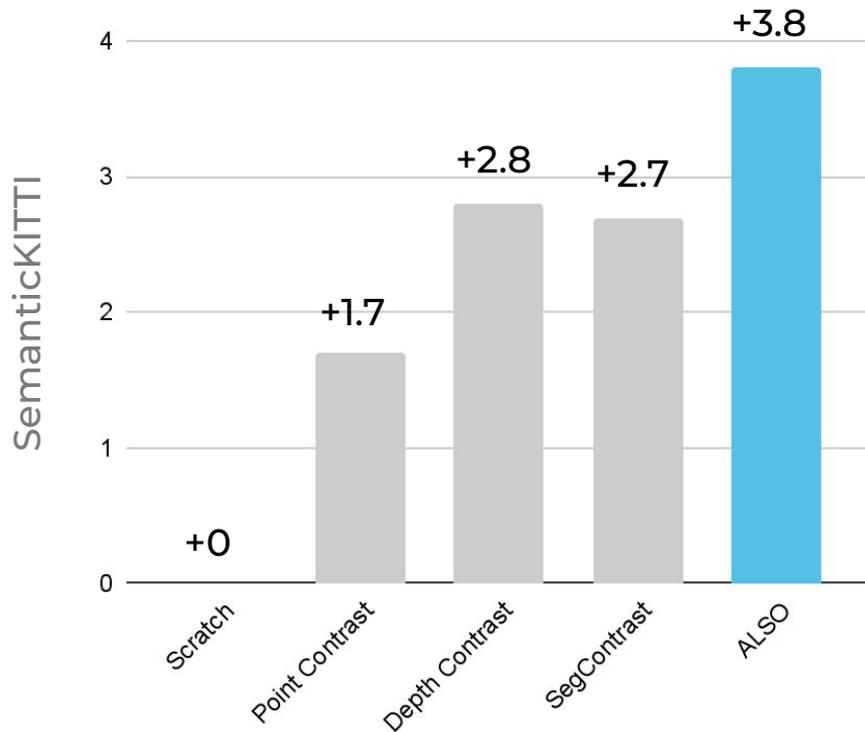
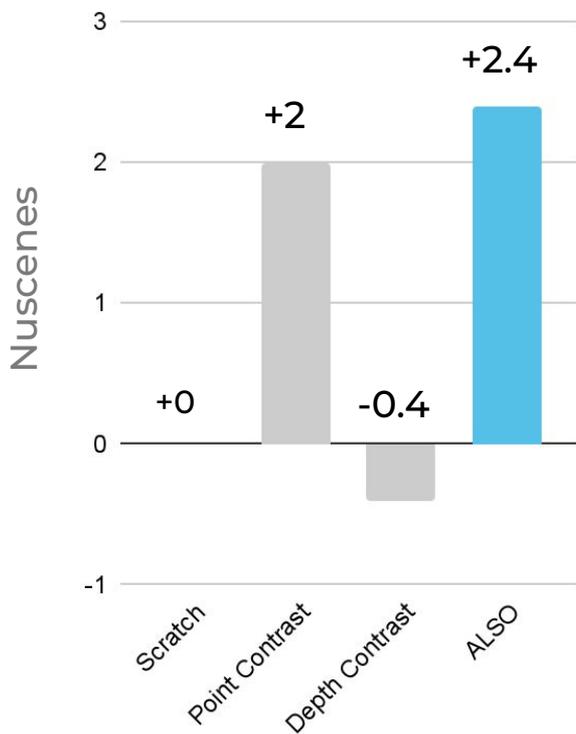


### Finetuning

- remove occupancy head
- add a single linear layer
- finetune the whole network

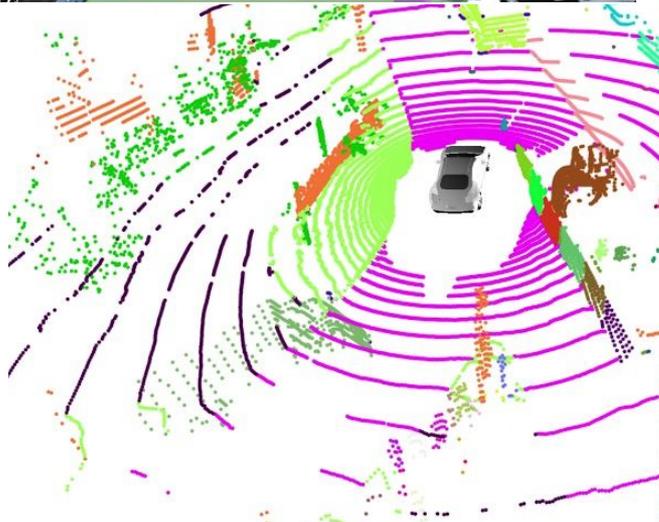
# Semantic segmentation - 1% annotated data

ALSO



# Overview

- I. Tasks
- II. Self-supervised learning
  - A. Geometric reconstruction
  - B. Contrastive learning
  - C. Distillation
- III. Domain adaptation
- IV. OpenWorld



# Contrastive learning

Contrastive methods

Car

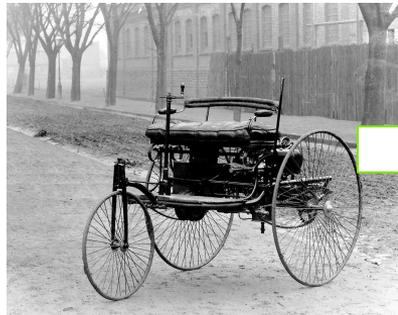


$Z_1$  Dissimilar  $Z_3$



Bus

Similar

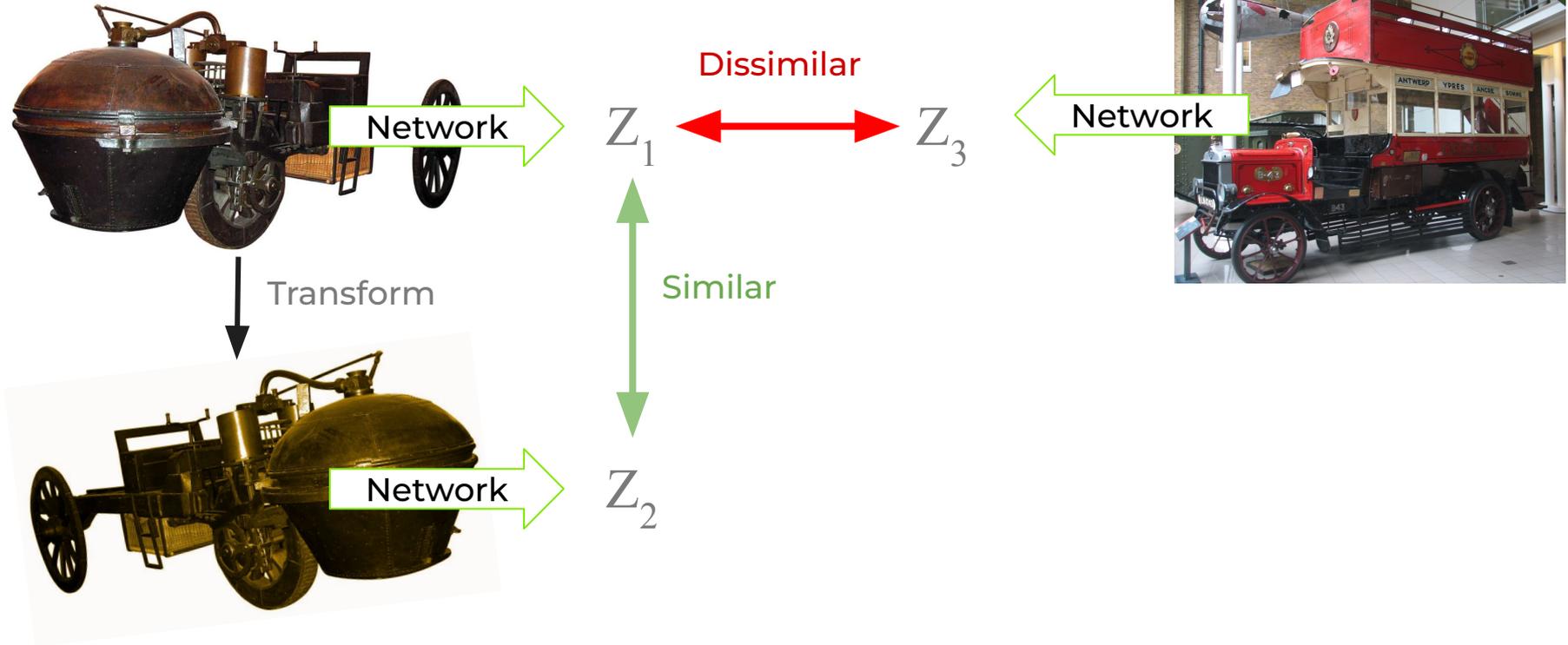


$Z_2$

Car

# Contrastive learning

## Contrastive methods



# Contrastive learning

## Contrastive methods

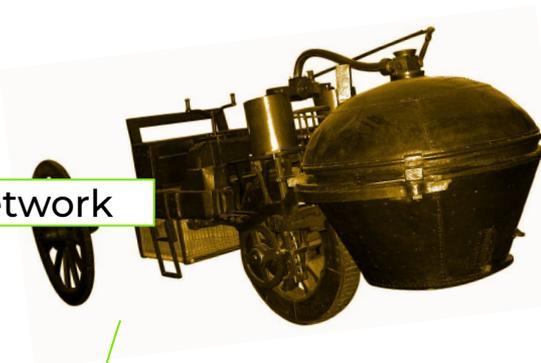


Network

$q$

$k_+$

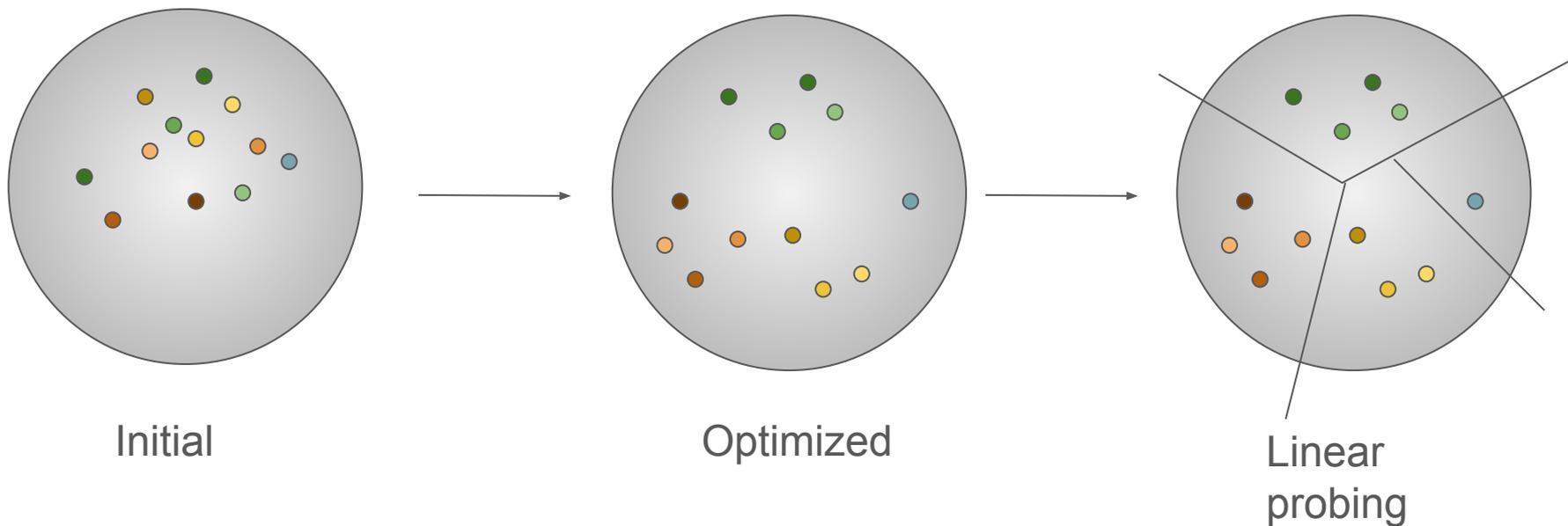
Network



$$\mathcal{L}_q = -\log \frac{\exp(q \cdot k_+ / \tau)}{\sum_{i=0}^K \exp(q \cdot k_i / \tau)}$$

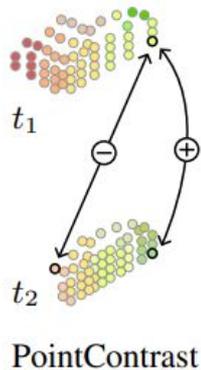


# Contrastive learning



# Contrastive self-supervised learning for point clouds

## Contrastive methods

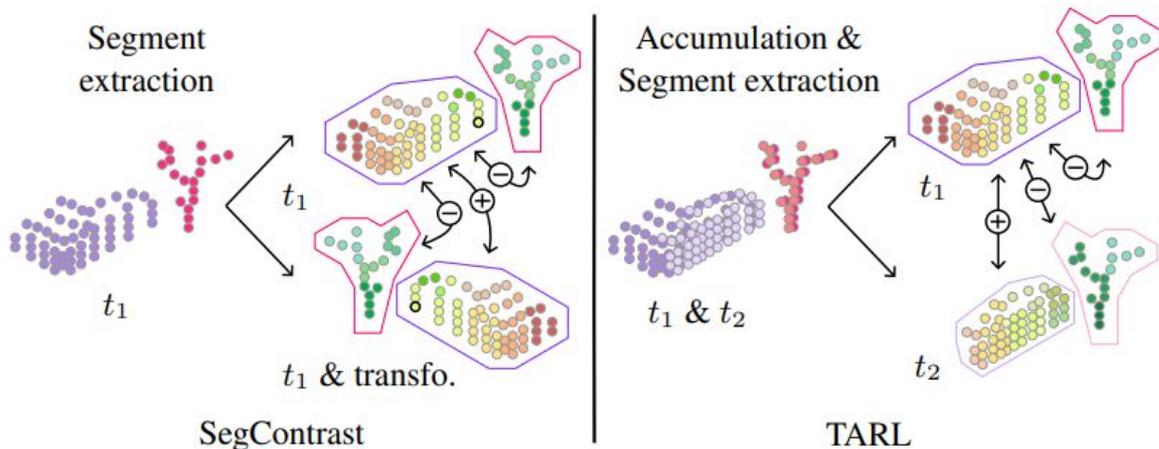


### PointContrast

- +++ Easy to implement
- - - Contrast in the same object

# Contrastive self-supervised learning for point clouds

## Contrastive methods

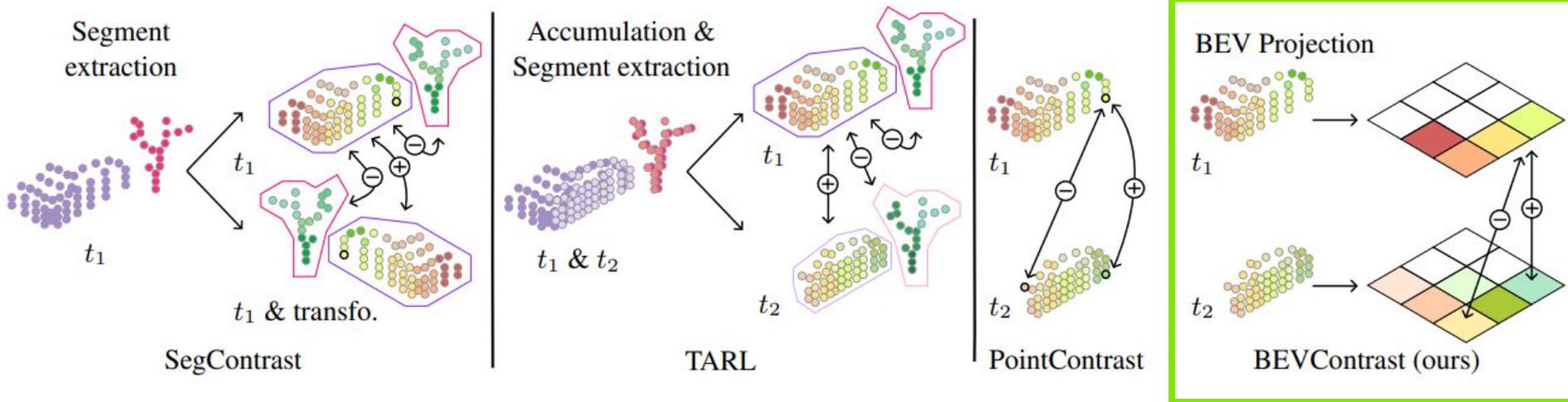


### SegContrast / TARL

- +++ Efficient thanks to object segmentation (temporal for TARL)
- - - Difficult to set up  $\rightarrow$  rely on HDBScan (hyperparameters)

# Contrastive self-supervised learning for point clouds

## Contrastive methods

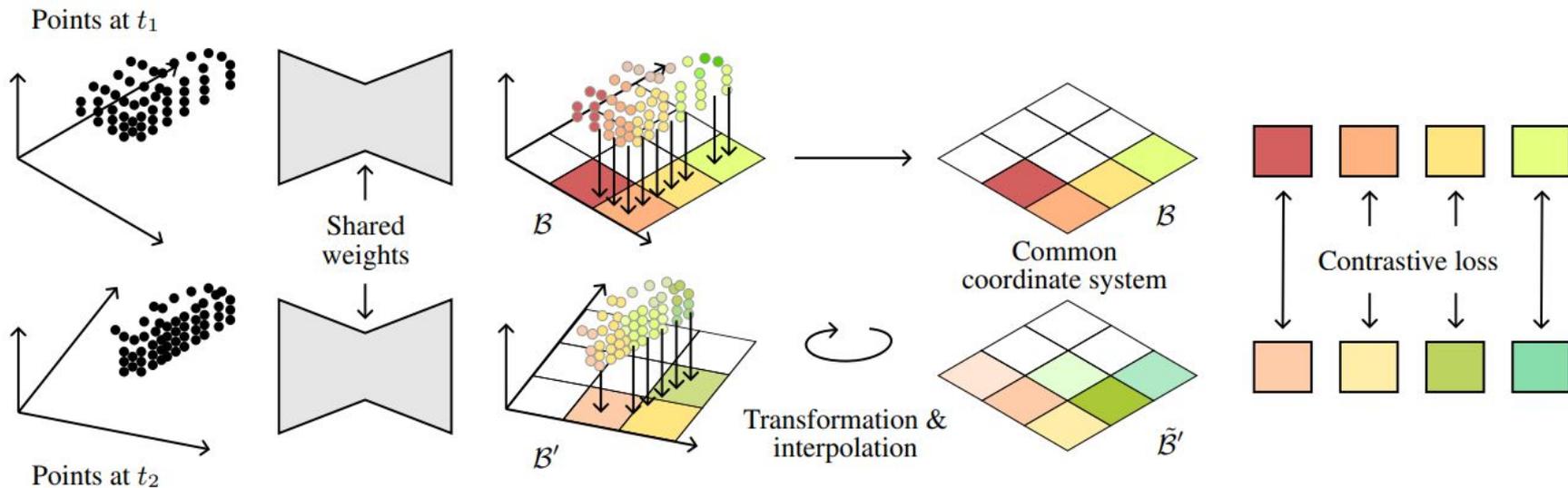


### BEVContrast

- +++ Simple  $\rightarrow$  easy projection in BEV
- +++ Object separation approximation with BEV cells

# BEVContrast

## Contrastive methods

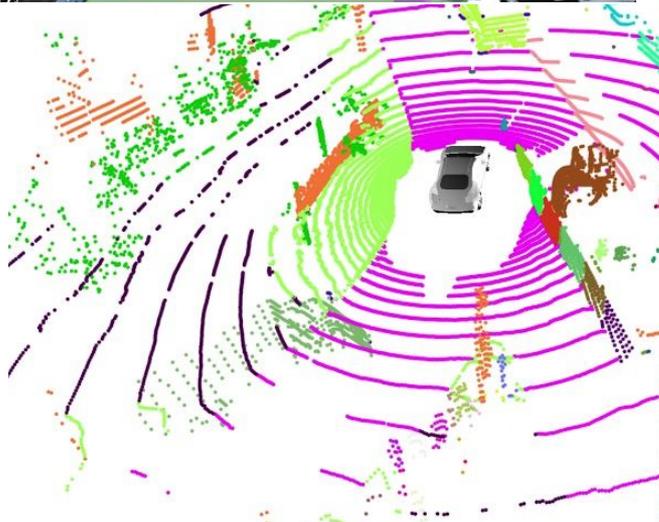


## Contrastive methods

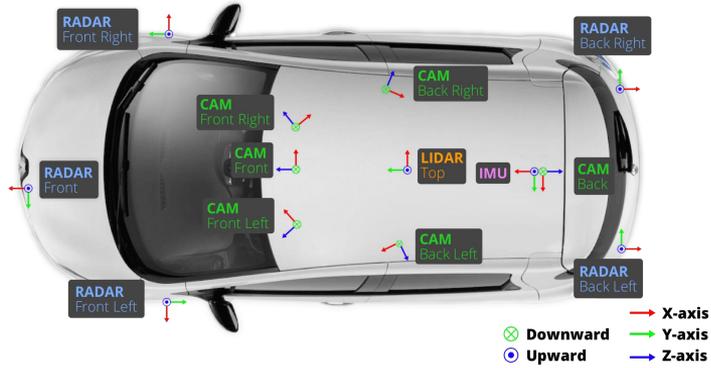
Dataset	Method	0.1%		1%		10%		50%		100%	
nuScenes	No pre-training	21.6	±0.5	35.0	±0.3	57.3	±0.4	69.0	±0.2	71.2	±0.2
	PointContrast <sup>†</sup> [40]	27.1	±0.5	37.0	±0.5	58.9	±0.2	69.4	±0.3	71.1	±0.2
	DepthContrast <sup>†</sup> [46]	21.7	±0.3	34.6	±0.5	57.4	±0.5	69.2	±0.3	71.2	±0.2
	ALSO [3]	26.2	±0.5	37.4	±0.3	59.0	±0.4	69.8	±0.2	71.8	±0.2
	<i>BEVContrast (ours)</i>	26.6	±0.5	37.9	±0.4	59.0	±0.6	70.5	±0.2	72.2	±0.1
SemanticKITTI	No pre-training	30.0	±0.2	46.2	±0.6	57.6	±0.9	61.8	±0.4	62.7	±0.3
	PointContrast <sup>‡</sup> [40]	32.4	±0.5	47.9	±0.5	59.7	±0.5	62.7	±0.3	63.4	±0.4
	SegContrast [29]	32.3	±0.3	48.9	±0.3	58.7	±0.5	62.1	±0.4	62.3	±0.4
	DepthContrast <sup>†</sup> [46]	32.5	±0.4	49.0	±0.4	60.3	±0.5	62.9	±0.5	63.9	±0.4
	STSSL [39]	32.0	±0.4	49.4	±1.1	60.0	±0.6	62.9	±0.7	63.3	±0.3
	ALSO [3]	35.0	±0.1	50.0	±0.4	60.5	±0.1	63.4	±0.5	63.6	±0.5
	TARL [30]	37.9	±0.4	52.5	±0.5	61.2	±0.3	63.4	±0.2	63.7	±0.3
	<i>BEVContrast (ours)</i>	39.7	±0.9	53.8	±1.0	61.4	±0.4	63.4	±0.6	64.1	±0.4

# Overview

- I. Tasks
- II. Self-supervised learning
  - A. Geometric reconstruction
  - B. Contrastive learning
  - C. Distillation
- III. Domain adaptation
- IV. OpenWorld



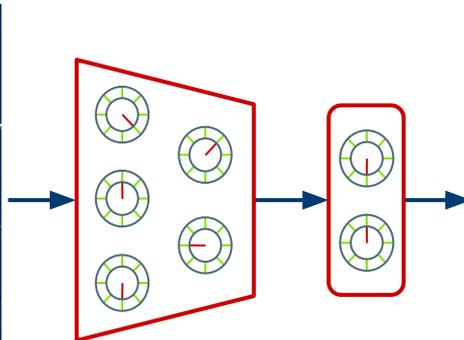
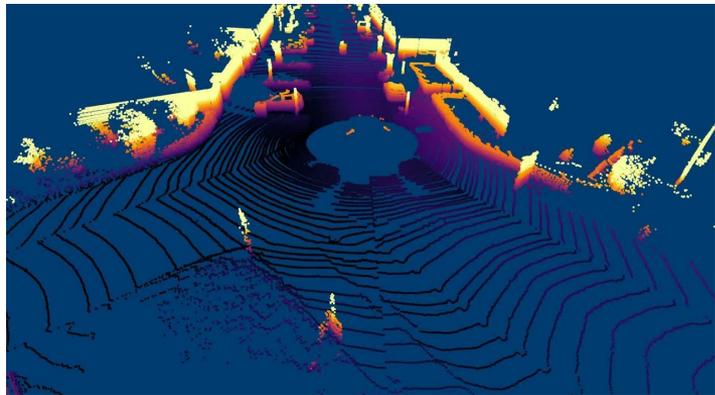
# Sensor setup



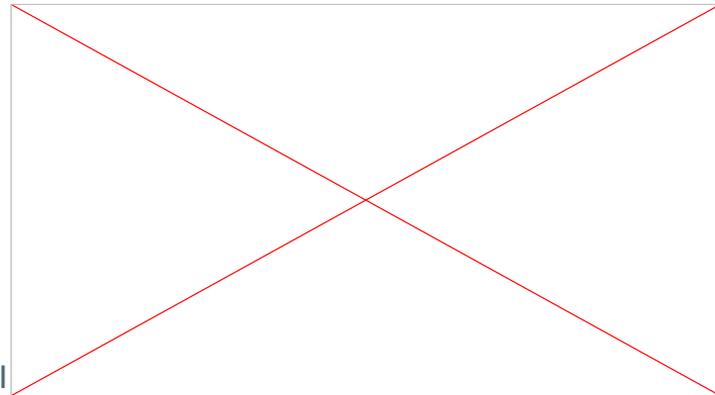
Example from nuScenes



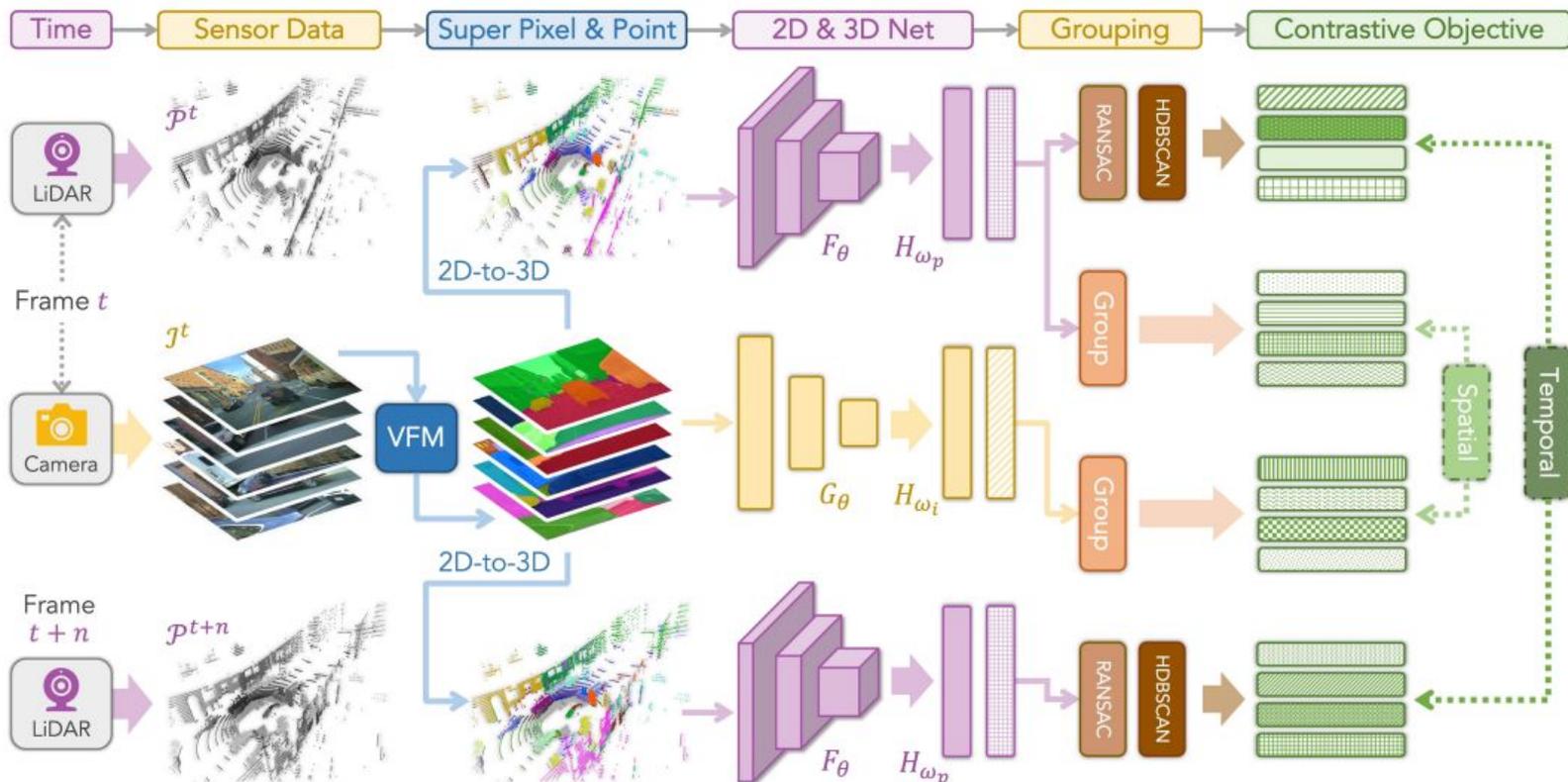
Task of interest: point cloud semantic segmentation



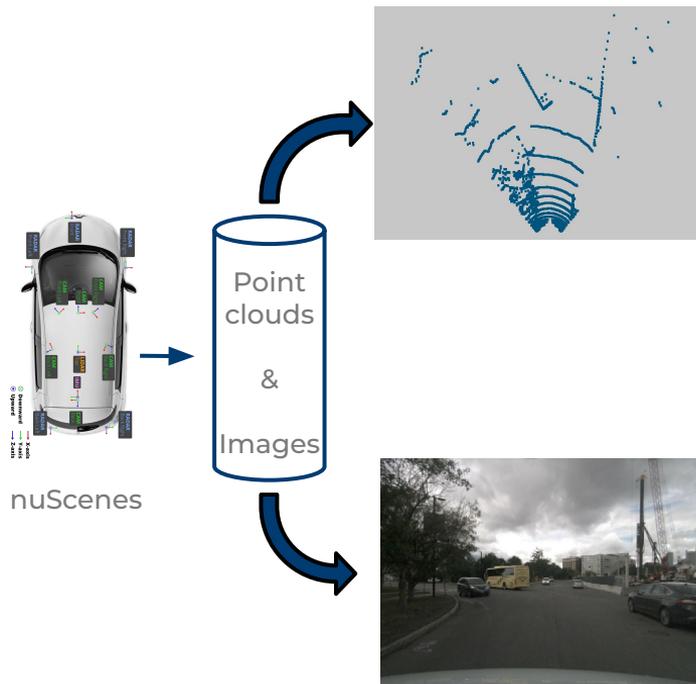
Example from SemanticKITTI



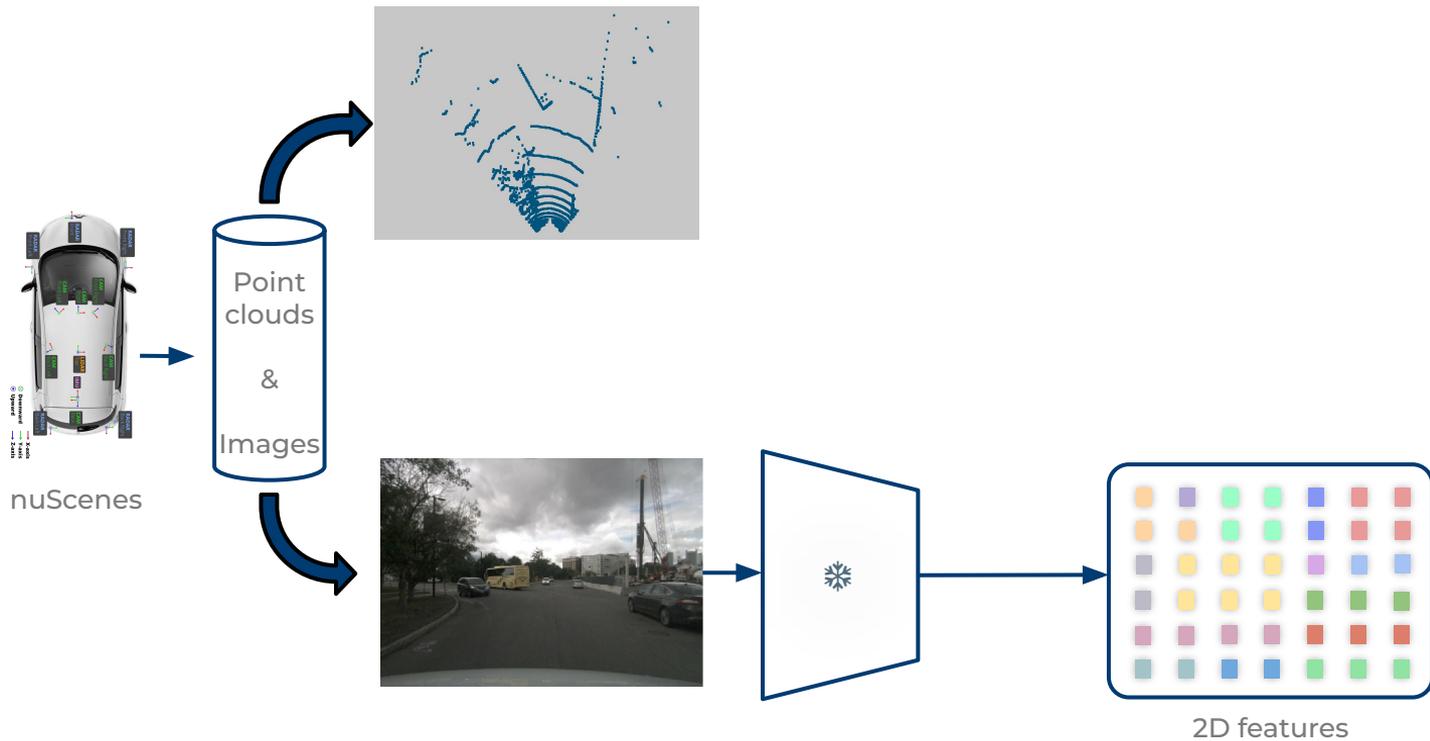
# SEAL - Segment any point cloud



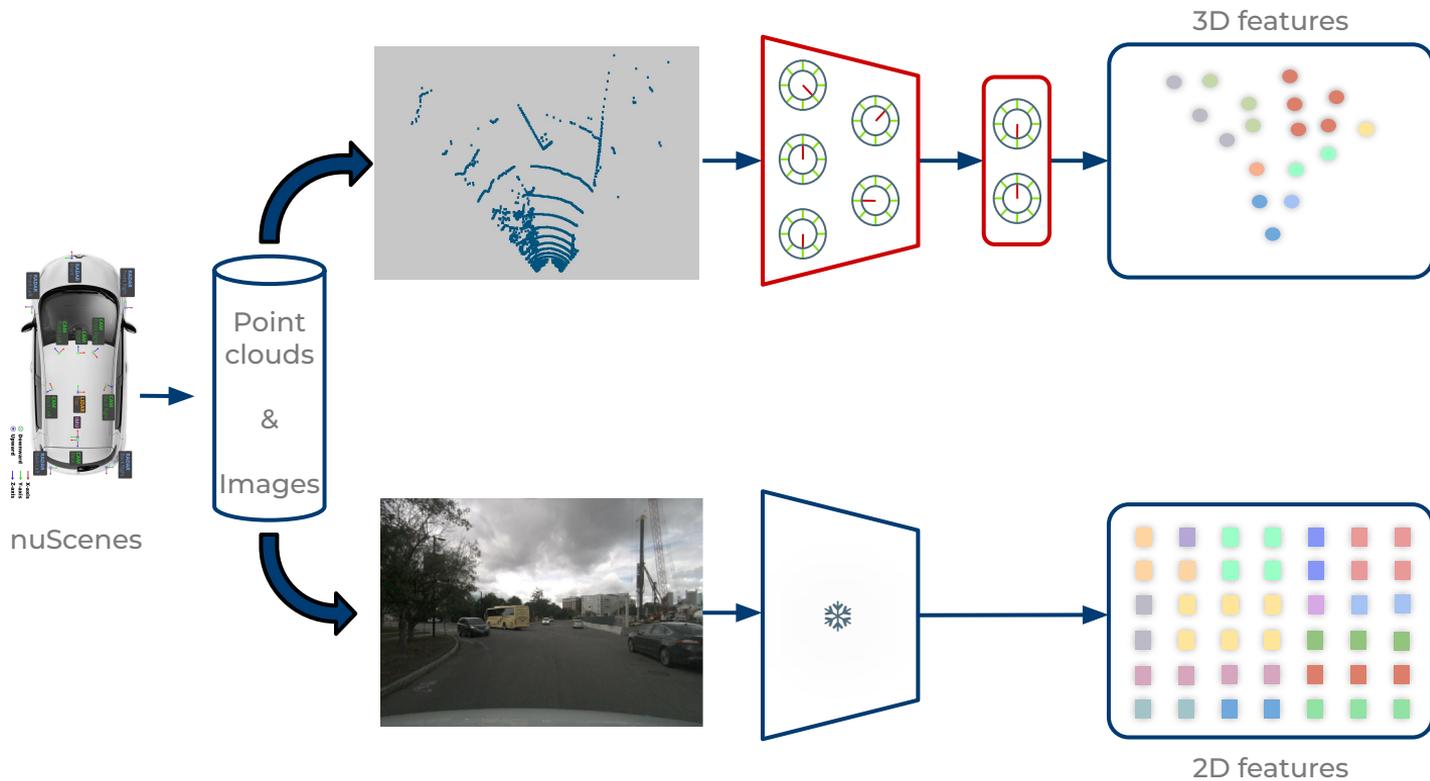
# ScaLR



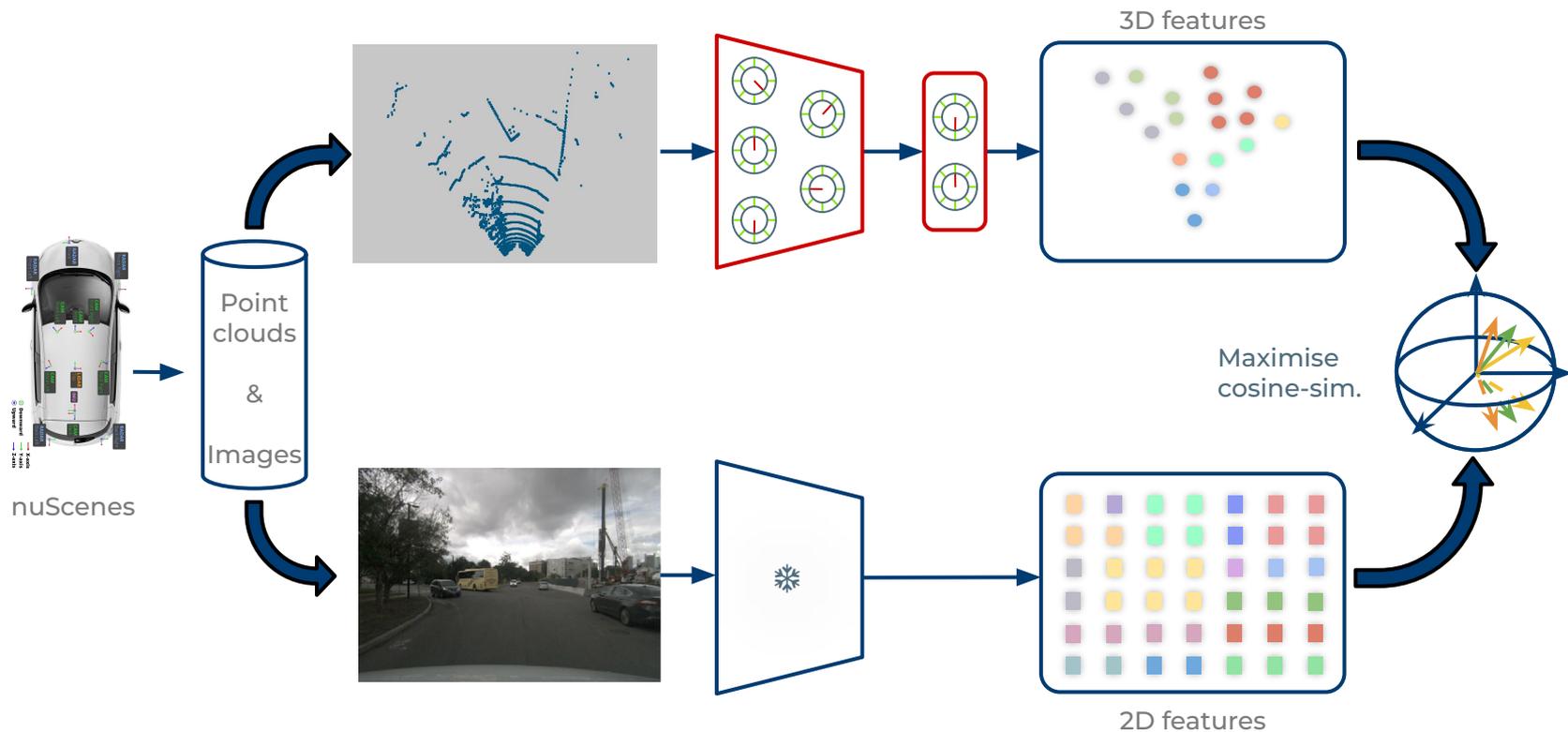
# ScaLR



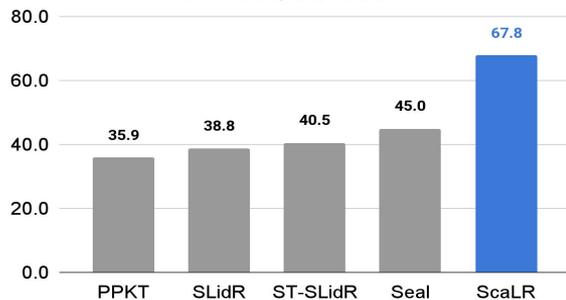
# ScaLR



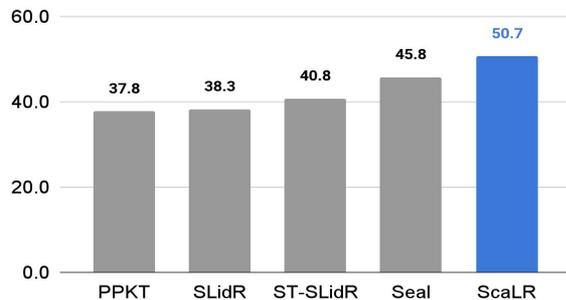
# ScaLR



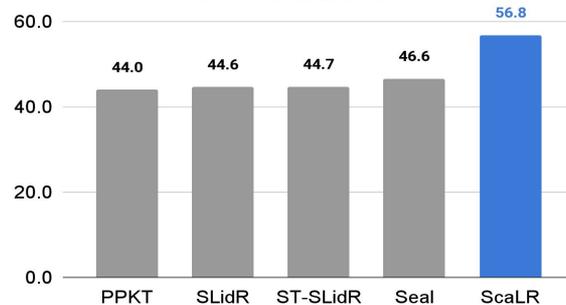
### Linear probing on nuScenes



### Finetuning with 1% labels on nuScenes

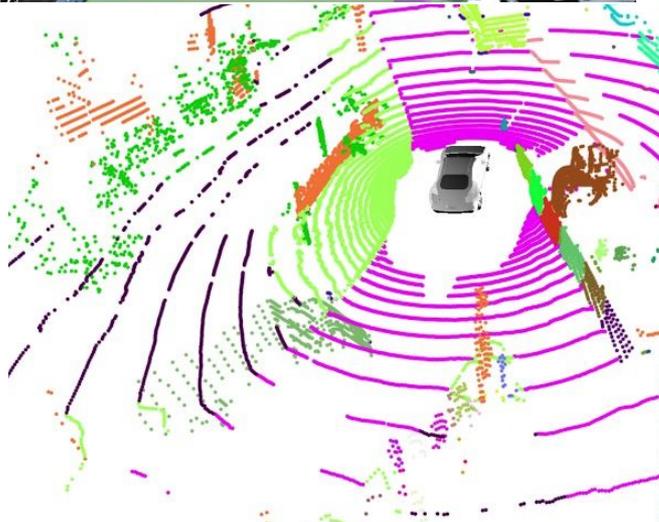


### Finetuning with 1% labels on SemKITTI

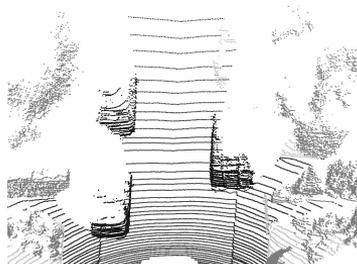
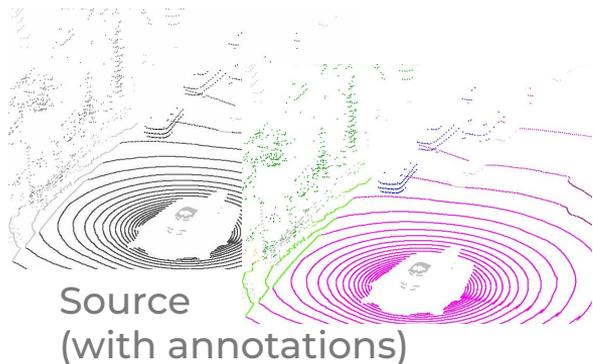


# Overview

- I. Tasks
- II. Self-supervised learning
  - A. Geometric reconstruction
  - B. Contrastive learning
  - C. Distillation
- III. Domain adaptation
- IV. OpenWorld



# Unsupervised domain adaptation



Network



## Domain gap

Different sensors

Different locations

Different objects

⇒ source model performs poorly

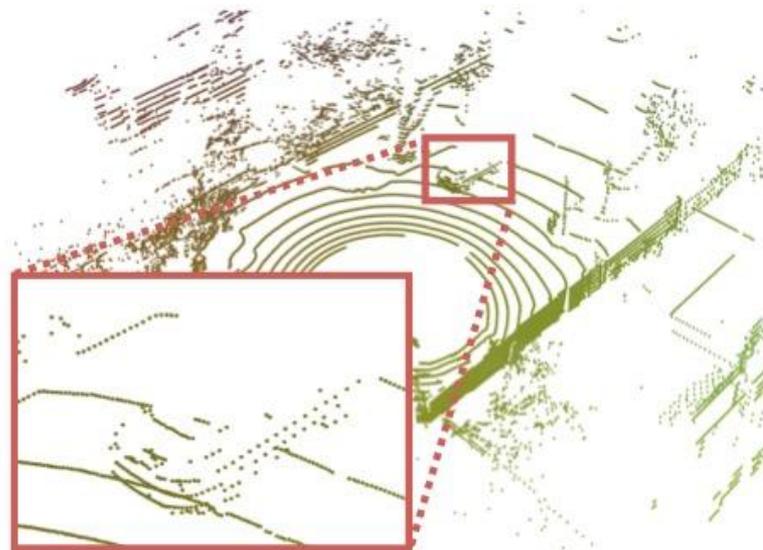
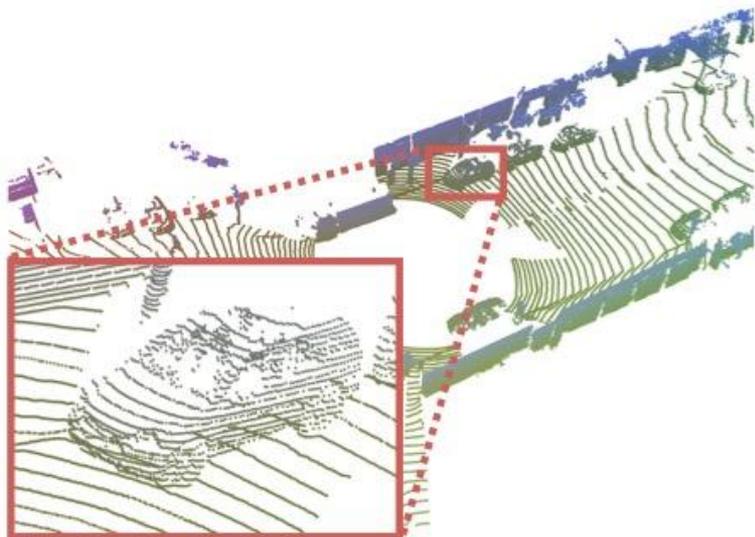
## Adaptation

Use target data to enhance performances

⇒ prevent collapse

# Example - sensor gap

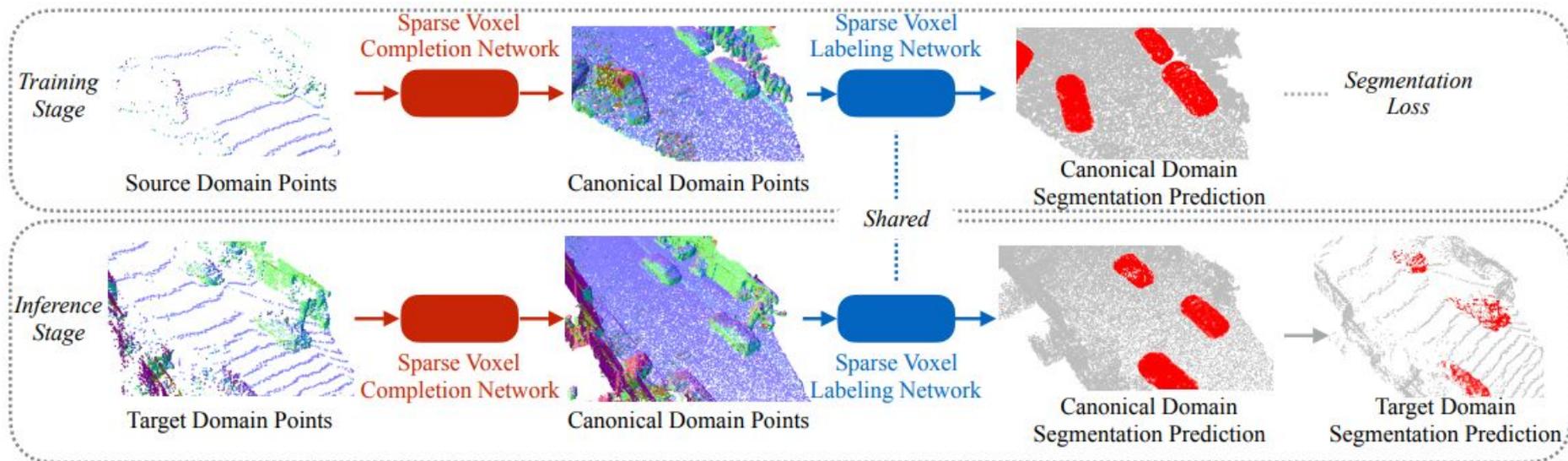
Domain adaptation



(a) captured by a 64-beam LiDAR (b) captured by a 32-beam LiDAR

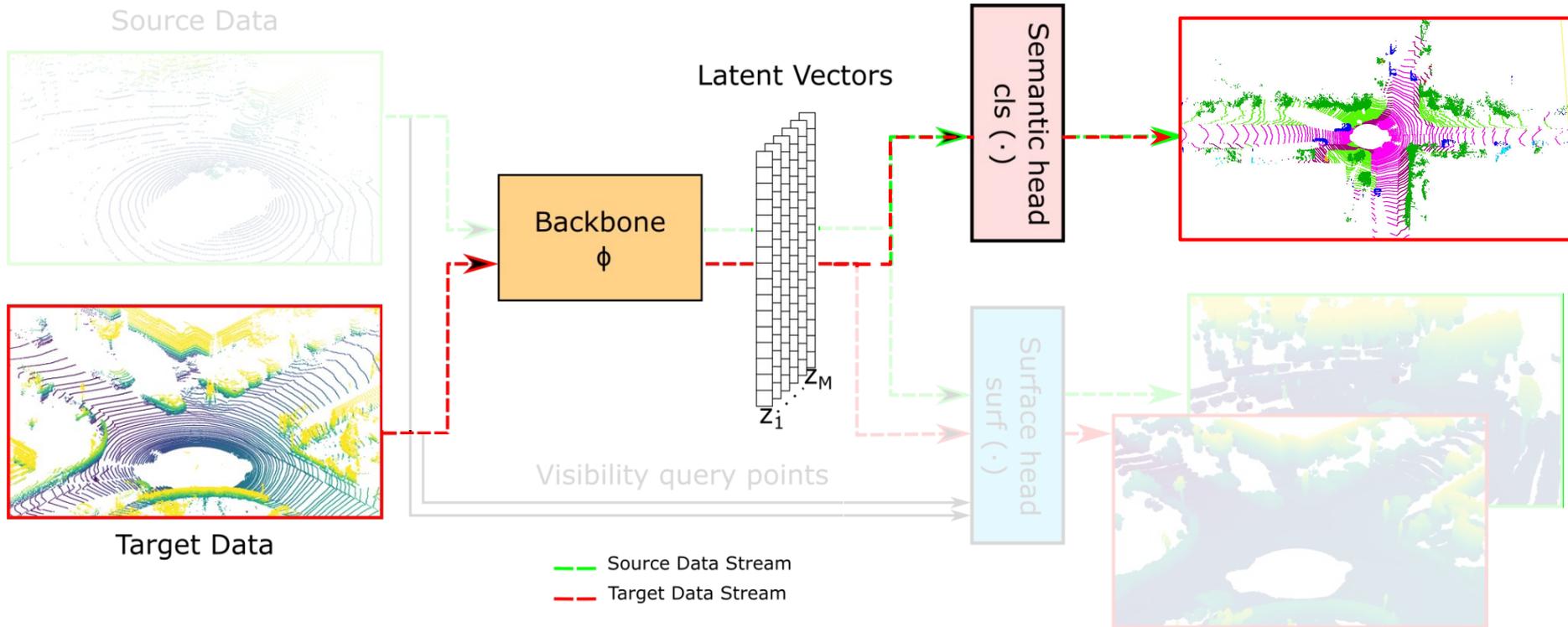
# Complete and label

Domain adaptation



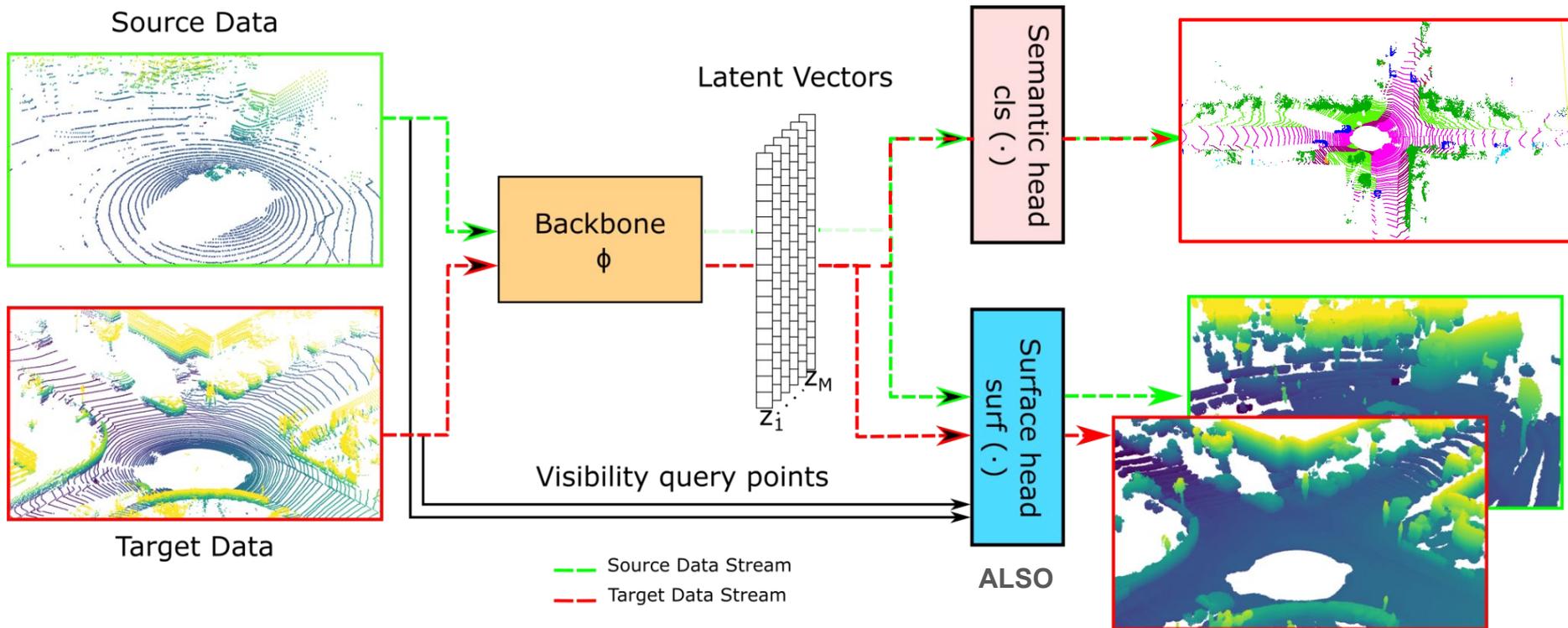
# SALUDA

Domain adaptation



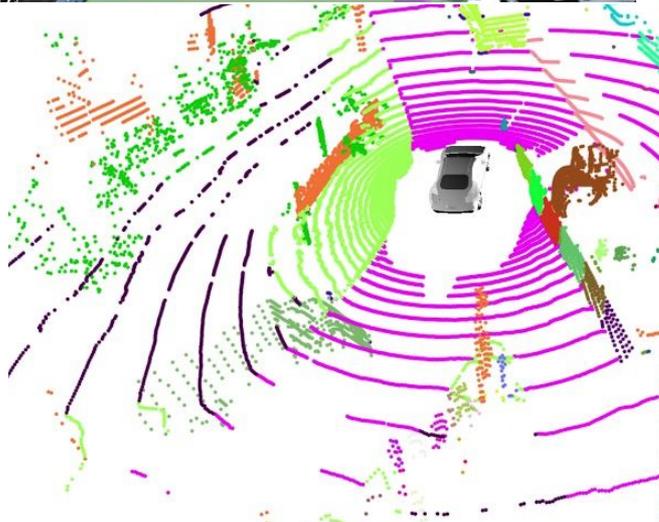
# SALUDA

Domain adaptation



# Overview

- I. Tasks
- II. Self-supervised learning
  - A. Geometric reconstruction
  - B. Contrastive learning
  - C. Distillation
- III. Domain adaptation
- IV. OpenWorld

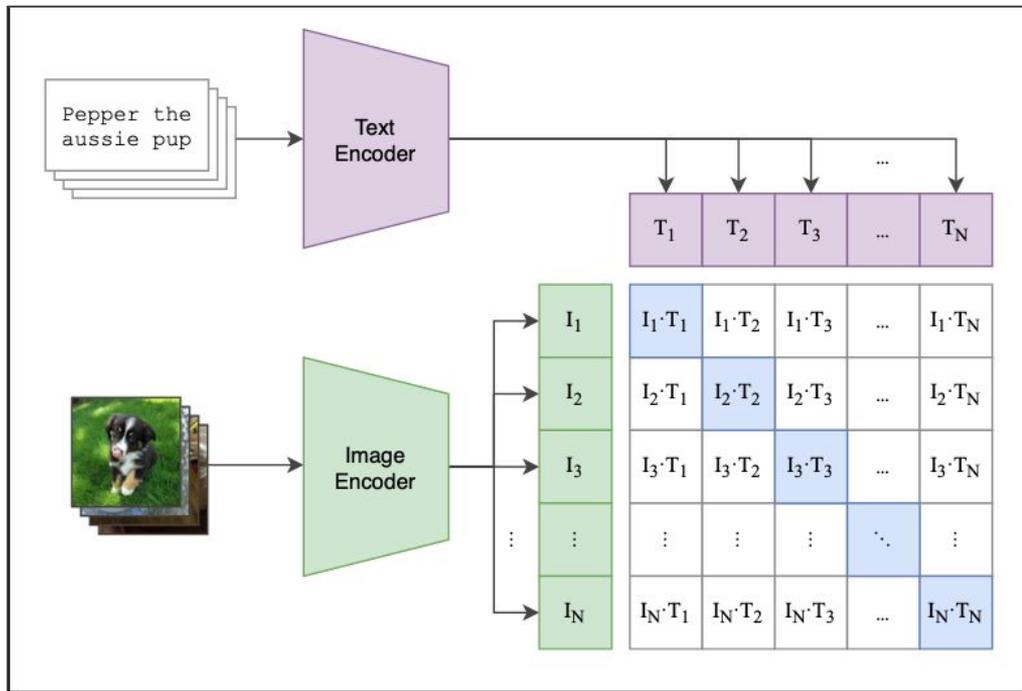




# Open vocabulary

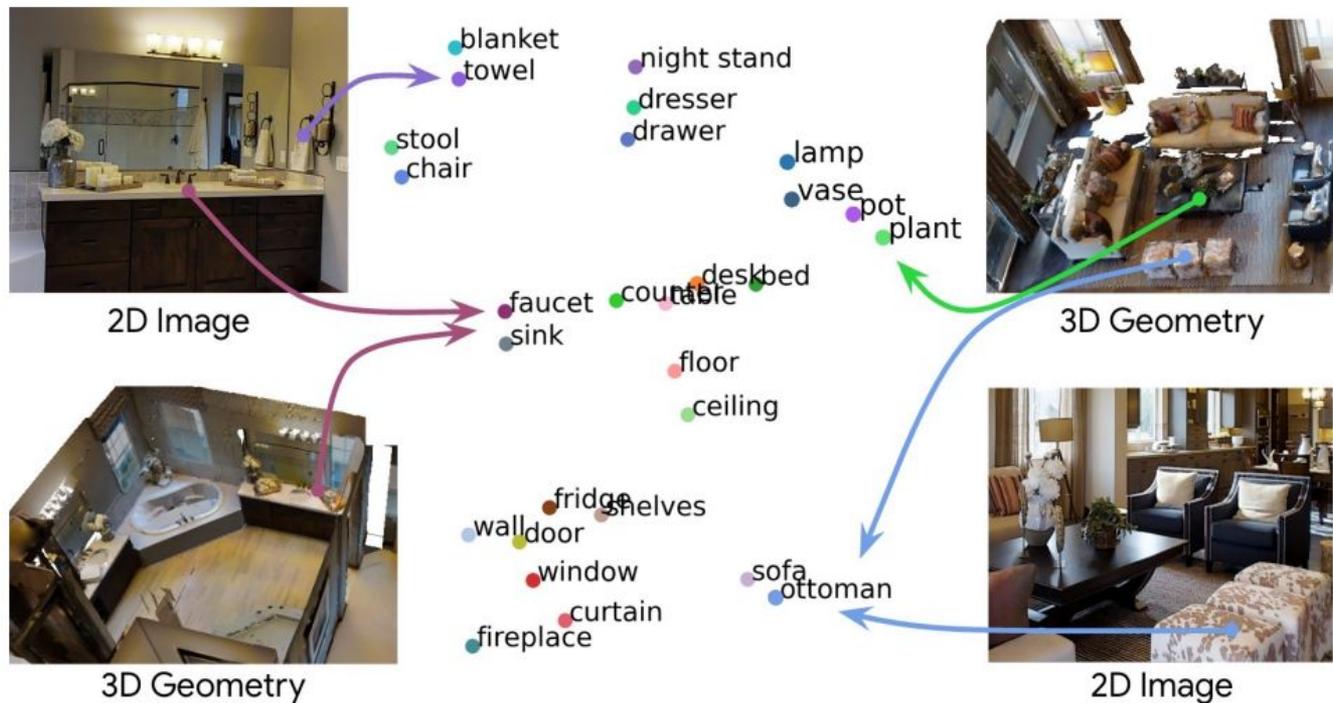
The world is complex

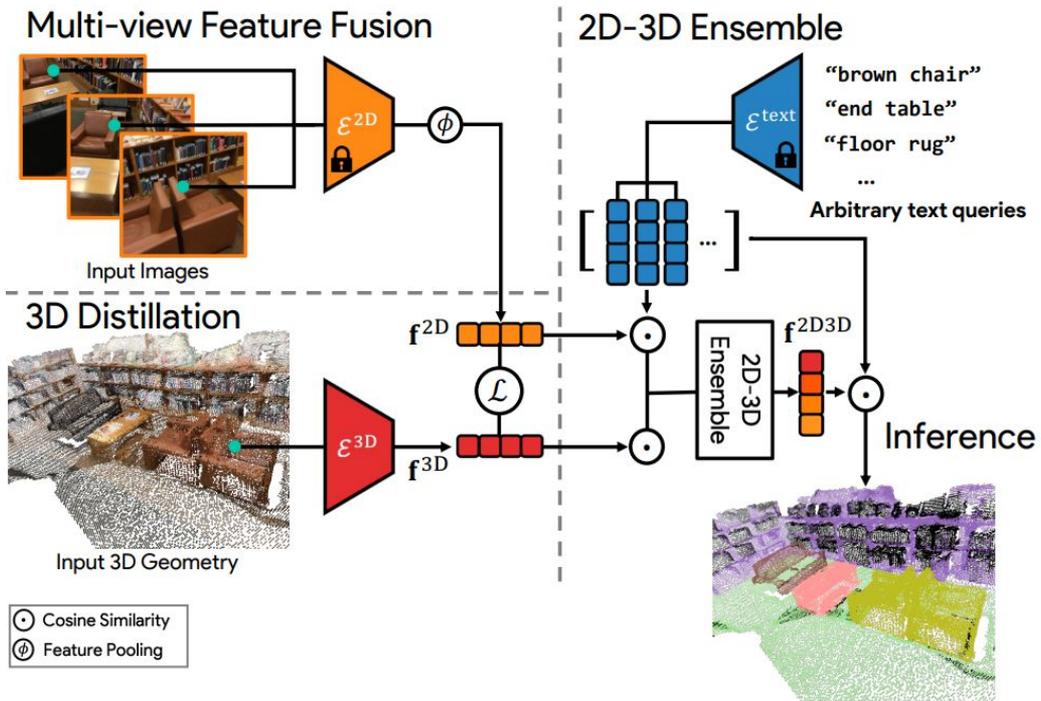




- Contrastive learning
  - Contrast positive/negative pairs
- Trained using 400 millions (**image / text**) pairs extracted from internet
  - Meta-data
  - Legends

# OpenScene

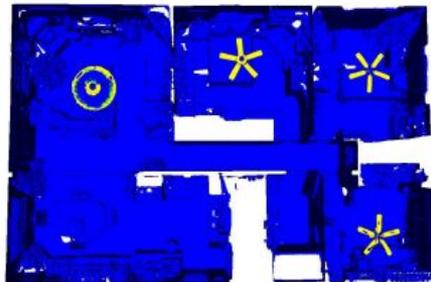




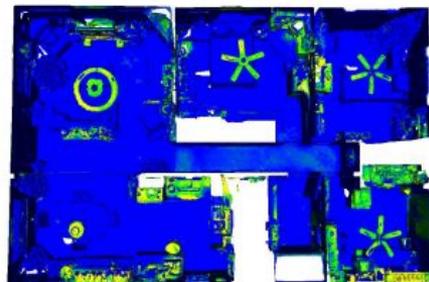
# OpenScene



Input 3D Point Cloud



"fan" - Object



"metal" - Material



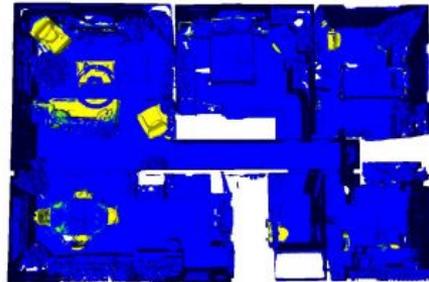
"kitchen" - Room Type



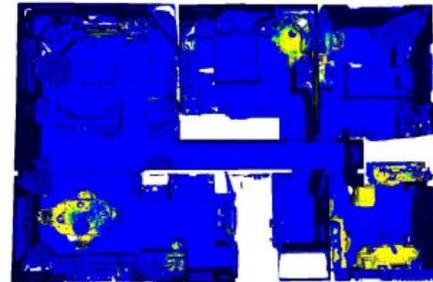
Zero-shot Semantic Segmentation



"soft" - Property



"sit" - Affordance



"work" - Activity

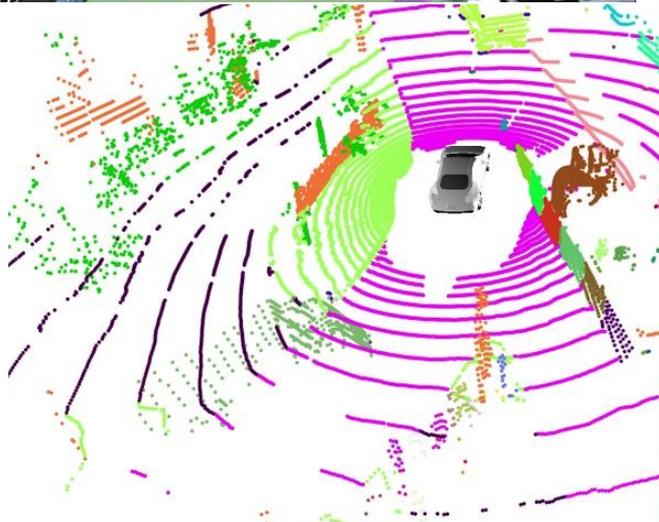
# OpenScene

Link to video

<https://pengsongyou.github.io/openscene>

# Overview

- I. Tasks
- II. Self-supervised learning
  - A. Geometric reconstruction
  - B. Contrastive learning
  - C. Distillation
- III. Domain adaptation
- IV. OpenWorld
- V. Conclusion



# Conclusion

A very short overview of some tasks

- There are many other
- Only few methods were presented, not state-of-the-art anymore

Practical session

Open Vocabulary on point cloud with MaskCLIP