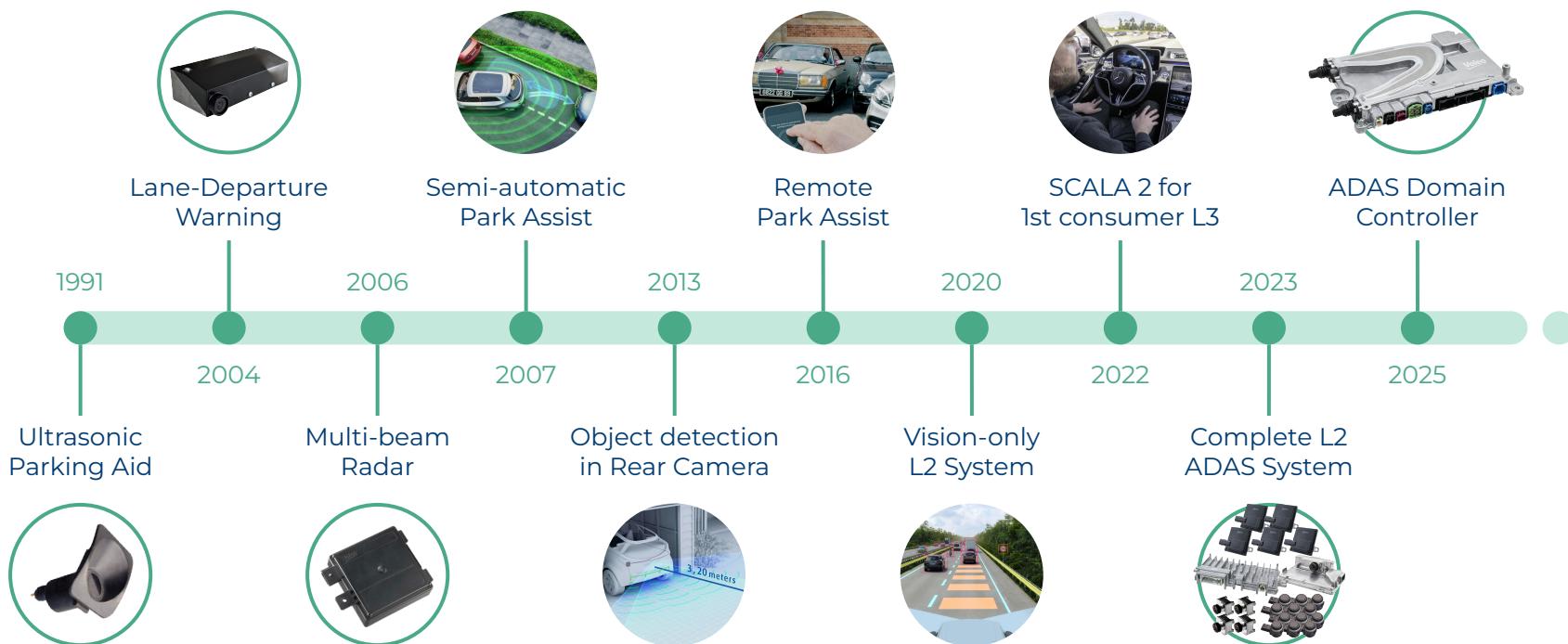




# Foundation models in the automotive industry

Andrei BURSUC  
[valeo.ai](http://valeo.ai)  
17 June 2024

# Valeo's history in ADAS



# Covering all ADAS segments

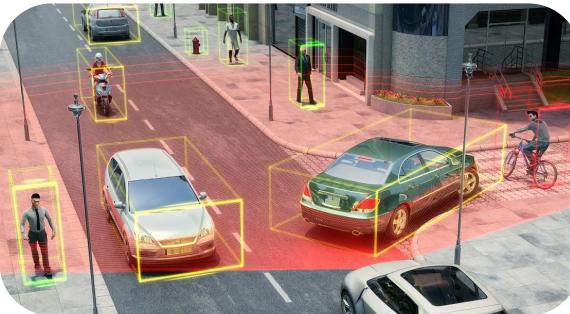
## One-stop solutions for **PARKING ASSISTANCE**



AI perception algorithms

Next generation sensors, cameras  
and processing units

## One-stop solutions for **SAFETY & ASSISTED DRIVING**



Scalable solutions from smart front  
camera to centralized architecture

Valeo algorithms for applications

## Perception & Functions for **AUTOMATED DRIVING**



Modules & back-up functions

High performance sensors, cameras  
and processing units

---

Expertise for solutions ranging from individual components to complete turnkey systems

# 1.5+ Billion sensors shipped in 30 years

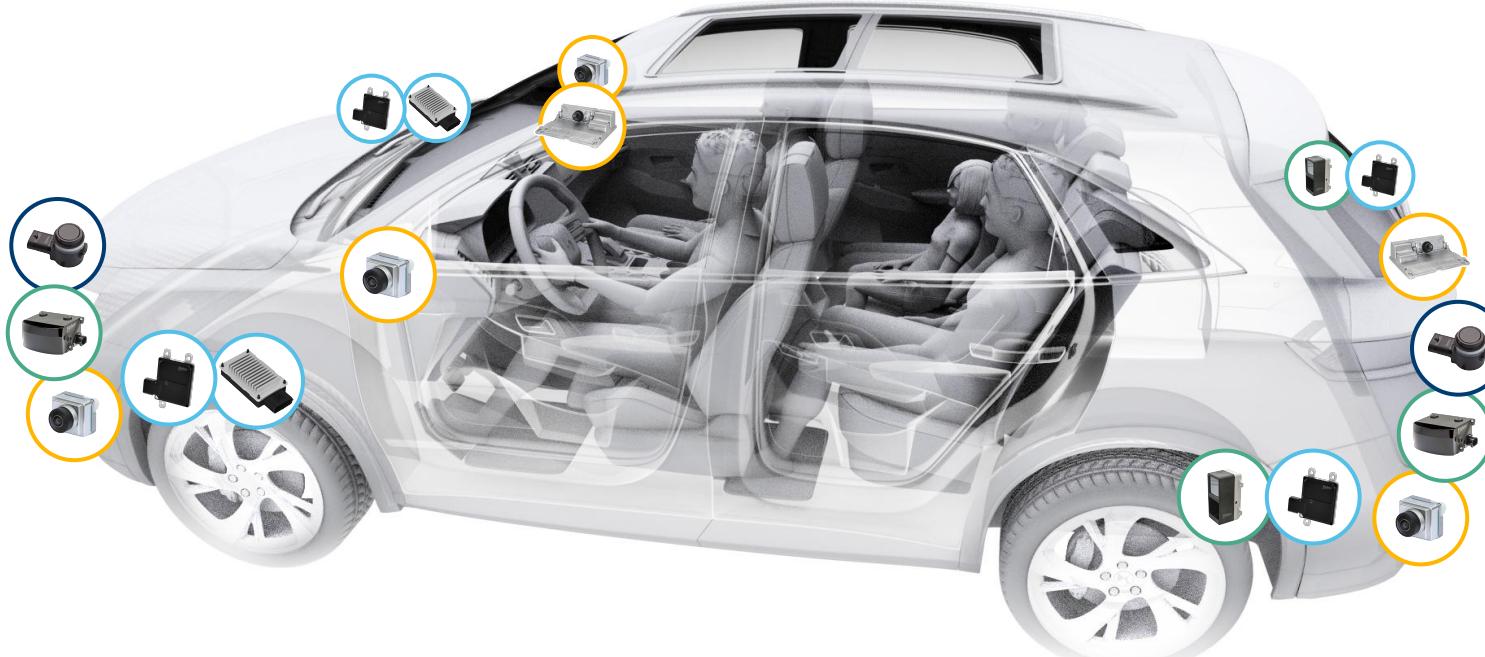


---

Another **1.5+ billion sensors** to be shipped in the next 5 years

---

# Valeo sensor suite



Ultrasonic  
sensors



ULS

Near field  
radars



RADARS

Mid range  
radars



Surround  
view cameras



CAMERAS

Long range  
cameras



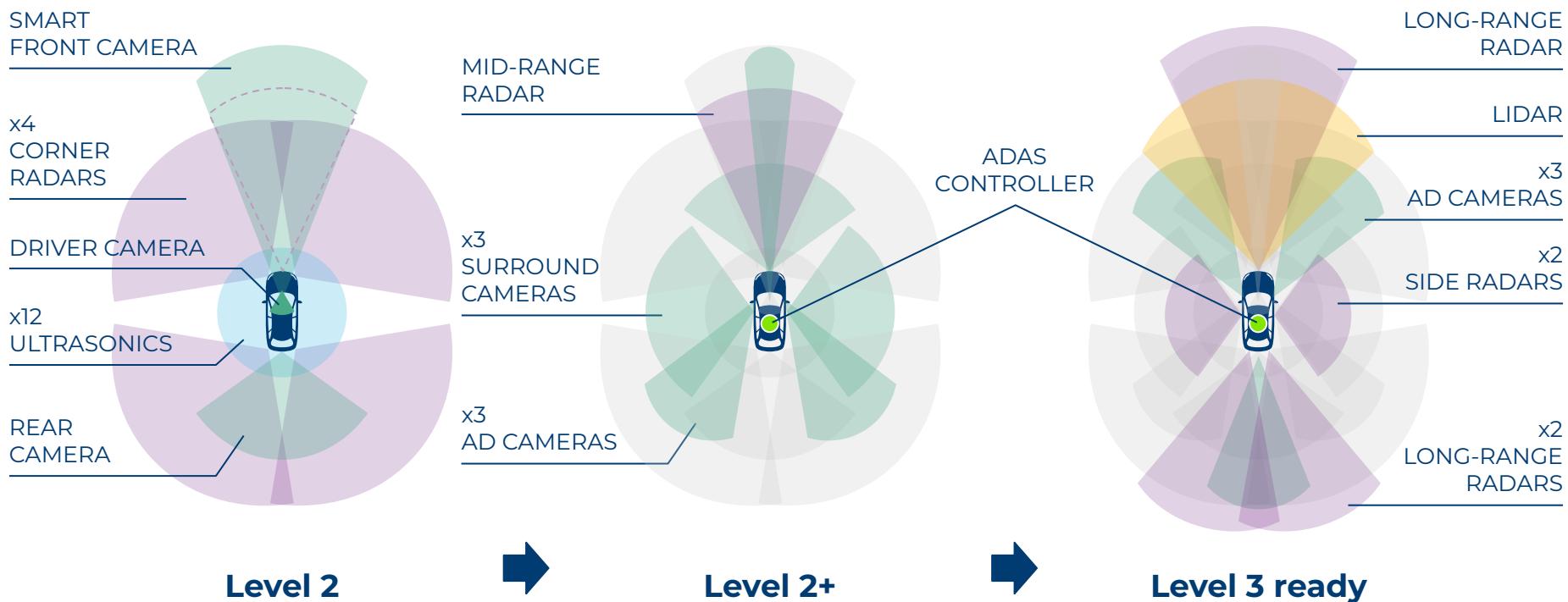
Near field  
lidars



LIDARS



# Scalable system architecture



# valeo.ai



- ~25 researchers & PhDs
- Dedicated to open research
- 10s of academic collabs across France and Europe
- Offices: Paris, Prague
- Topics: perception, data efficiency, forecasting, reliability, explainability



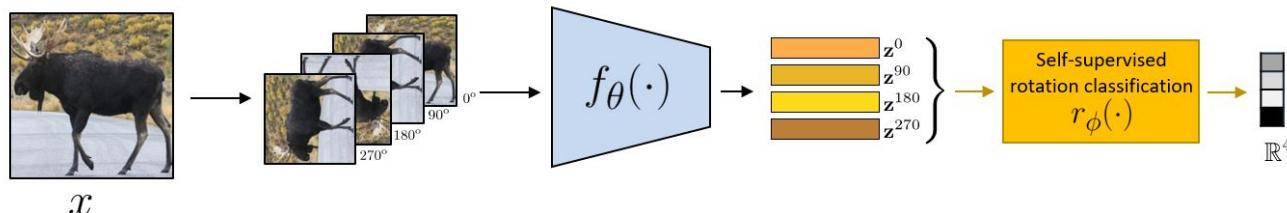
<https://valeoai.github.io>

**Foundation models:  
train once, use many times**

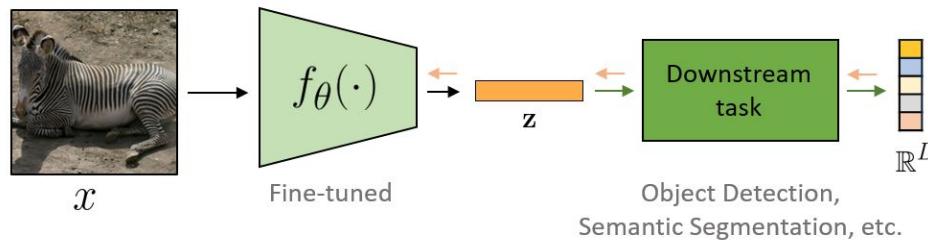
*“A foundation model is any model that is trained on broad data (generally using self-supervision at scale) that can be adapted (e.g., fine-tuned) to a wide range of downstream tasks.”* (*Bommasani et al., 2021*)

# Recap - self-supervised learning pipeline

Stage 1: Pretrain network on pretext task (without human labels)



Stage 2: Fine-tune network for new task with fewer labels



**Why would we want to do self-supervised learning?**

**For driving scenes:**

**Data is enormous,**

**Scalability of recent models (ChatGPT et al.),**

**When possible, manual annotation is difficult,**

**Annotation is impossible for some important tasks**

**Training datasets are evolving continuously**

# Difficult to keep the pace with an ever changing world

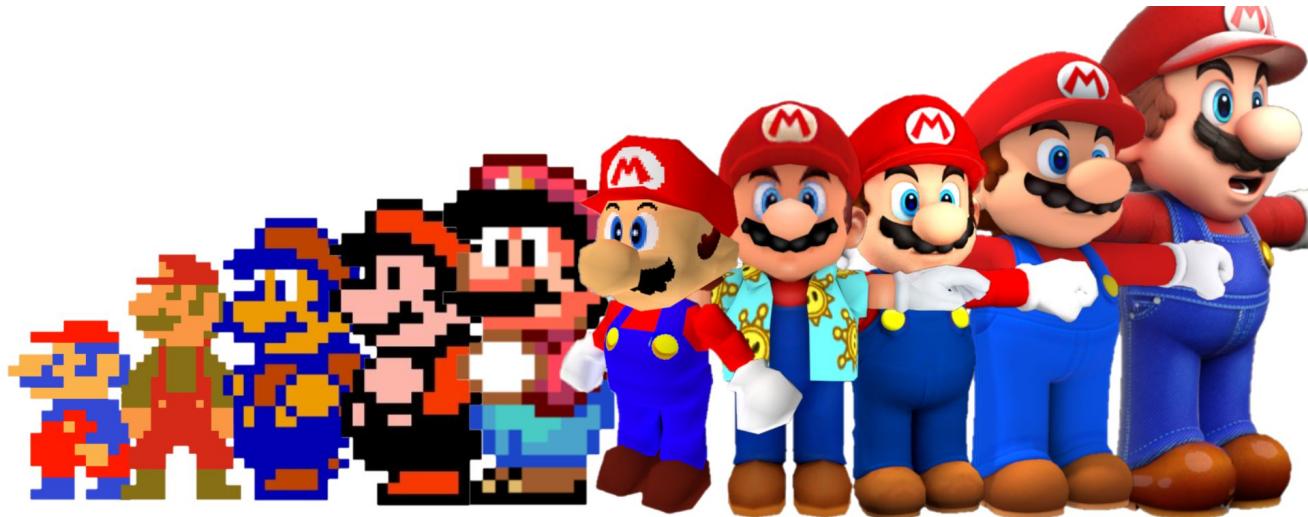
1980-1989



*Men's fashion trends 1980-1989*

- Data distributions shift all the time, e.g., fashion trends, vehicle types
- Infeasible to launch large annotation campaigns each time

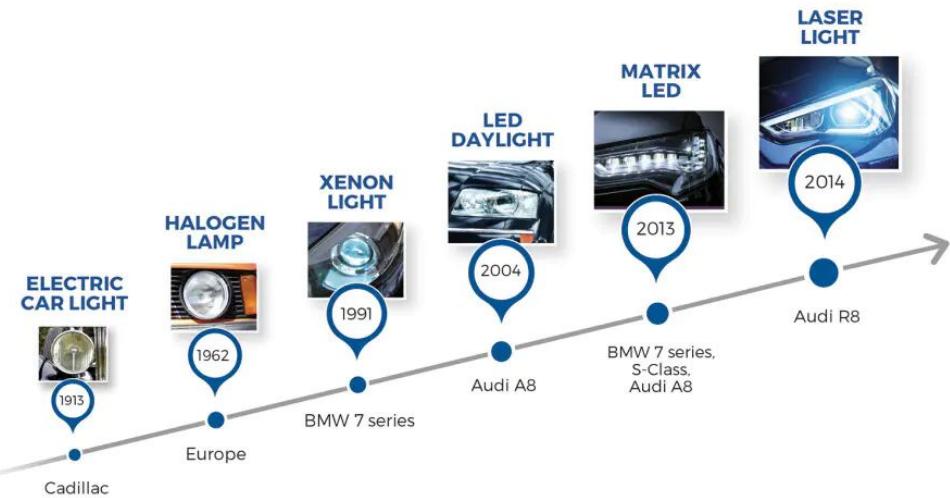
## Difficult to keep the pace with an ever changing world



*Super Mario from 1981 to 2017*

- Sensors specs are frequently upgraded
- Infeasible to launch large annotation campaigns each time

# Difficult to keep the pace with an ever changing world

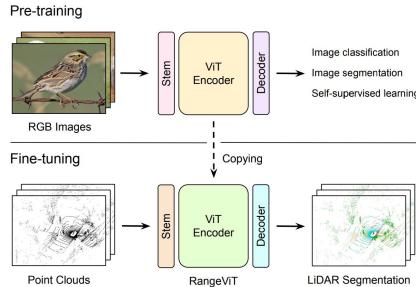


Halogen vs. LED

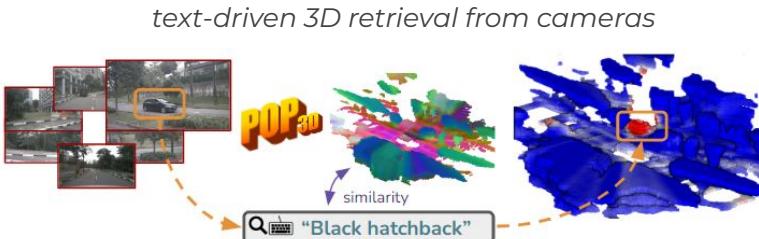
- Sensors specs are frequently upgraded
- Headlamps change the appearance of the scenes and of the vehicles to detect

# New use-cases possible with recent foundation models

Stage 2: Can also be distillation, auto-labeling, data mining, active learning, model initialization, etc.



RangeViT [CVPR'23]

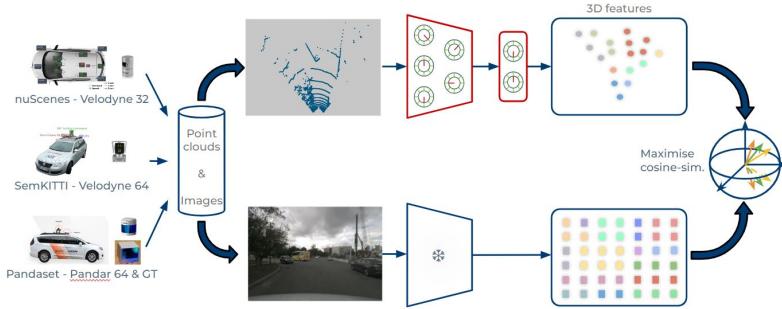


POP-3D [NeurIPS'23]

unsupervised semantic segmentation



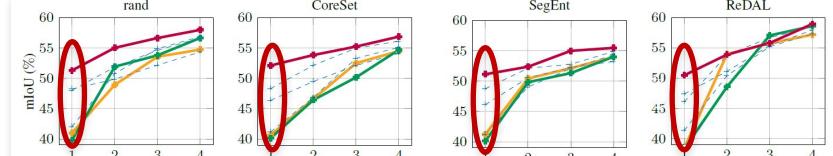
Drive&Segment [ECCV'22]



ScaLR [CVPR'24]

kickstart active learning

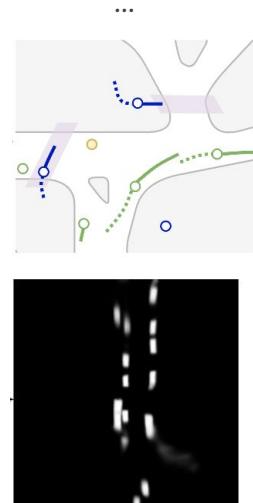
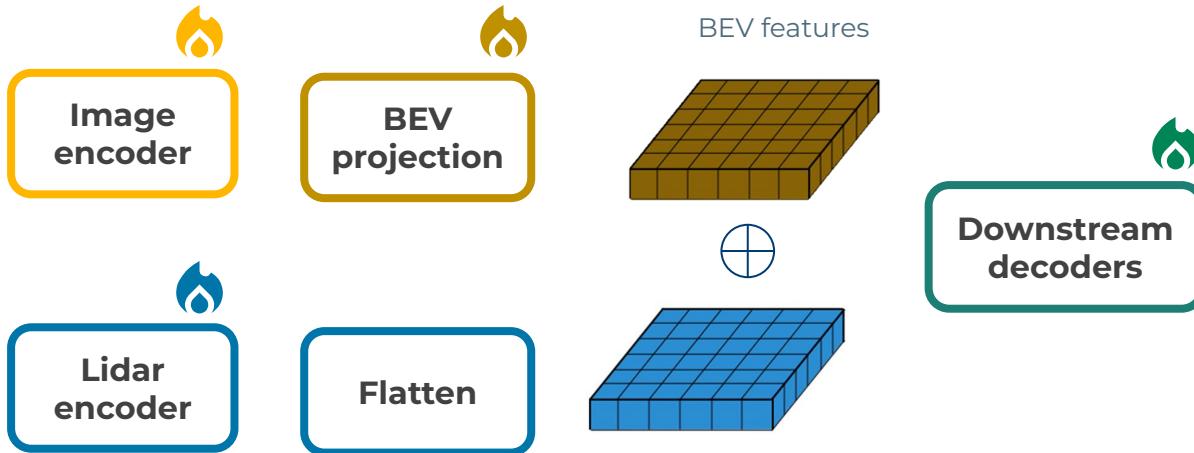
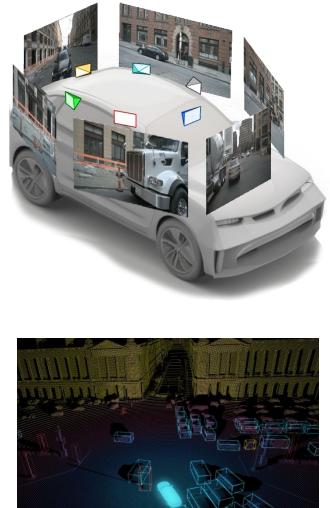
—+— random    —○— KMcentroid    —■— KMfurthest    —— SeedAL



SeedAL [ICCV'23]

# Modular pretraining

Perception +  
Forecasting

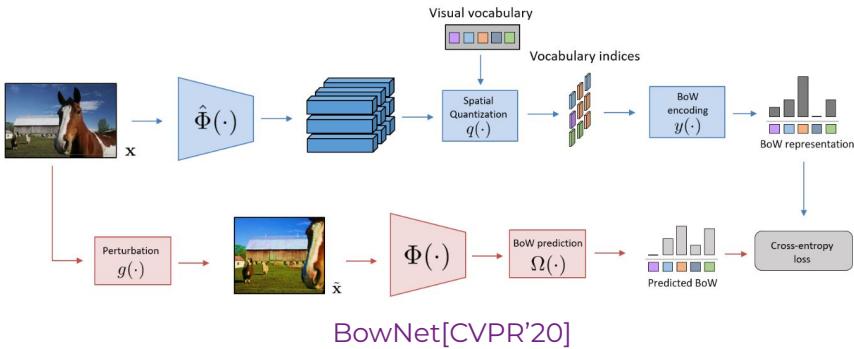


- Pretrain each branch individually to leverage knowledge from different distributions and data sources, sensors, sensor rigs
- Plug pretrained modules to downstream ones and continue training

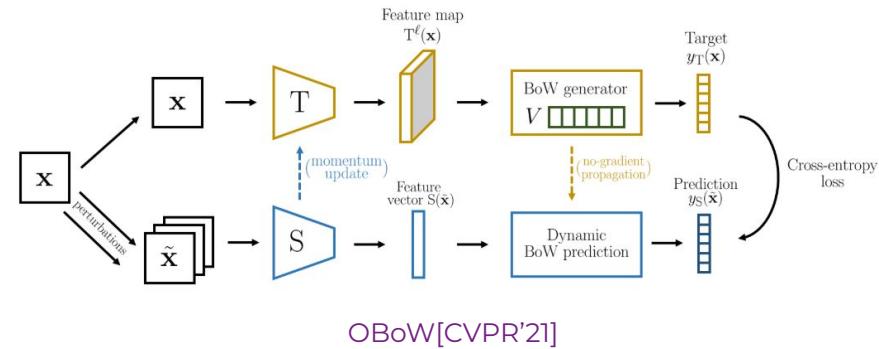
**One model for all**

# Pretraining image encoders

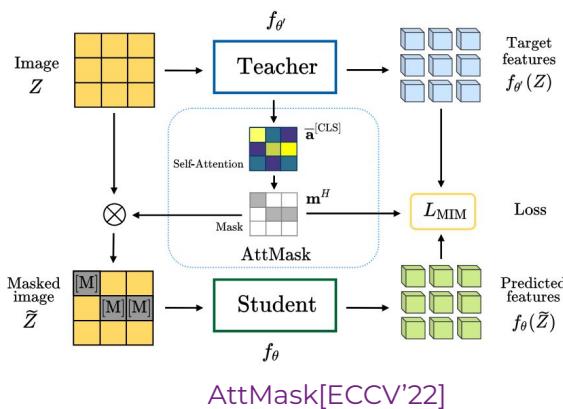
# Pretraining image encoders



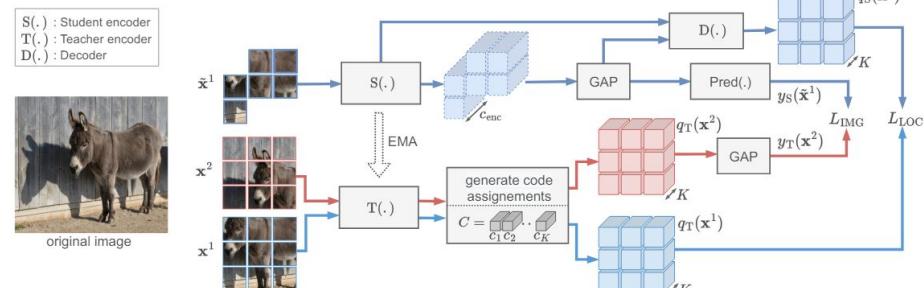
BowNet[CVPR'20]



OBoW[CVPR'21]



AttMask[ECCV'22]



MOCA[TMLR'24]

- Long-term exploration of (dense) SSL strategies with focus on 2D downstream performance and few-shot adaptation

## **Challenges in using self-supervised learning with autonomous driving data**



$\sim$  ImageNet



$\sim$  BDD100K



$\sim$  ImageNet



$\sim$  BDD100K



$\sim$  ImageNet



$\sim$  BDD100K



$\sim$  ImageNet



$\sim$  BDD100K



$\sim$  ImageNet



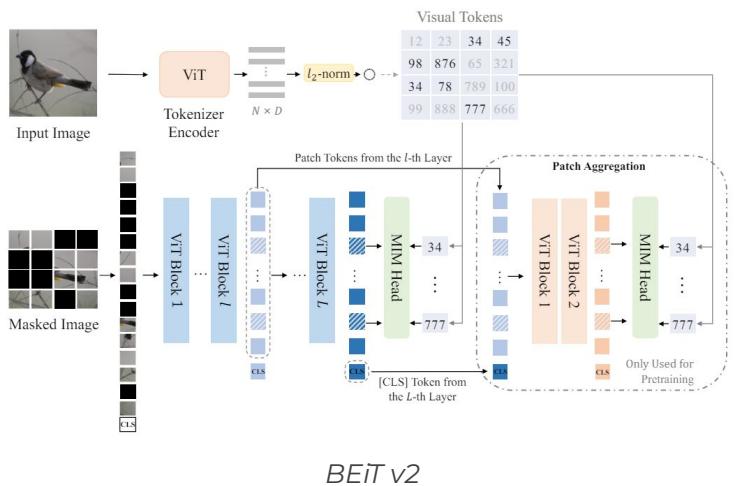
$\sim$  BDD100K

**From the ImageNet perspective, AD data is “boring”.  
Most SSL methods seem devised for well curated datasets.**



# Pretraining image encoders

- AD pretraining can struggle to match pretraining on large image collections (e.g., ImageNet1K/21K, Instagram, etc.) (Chen et al., 2021)
- Recent SSL methods (BEiTv2) distill knowledge from pretrained models (DINO, CLIP)
- Alternatively, finetuning large pretrained models is effective (Wei et al., 2024)



BEiT v2

Backbone	Fine-tune Method	Trainable Params*	mIoU			
			Citys	BDD	Map	Avg.
CLIP [63] (ViT-Large)	Full	304.15M	51.3	47.6	54.3	51.1
	Freeze	0.00M	53.7	48.7	55.0	52.4
	Rein	2.99M	<b>57.1</b>	<b>54.7</b>	<b>60.5</b>	<b>57.4</b>
MAE [27] (Large)	Full	330.94M	53.7	<b>50.8</b>	58.1	54.2
	Freeze	0.00M	43.3	37.8	48.0	43.0
	Rein	2.99M	<b>55.0</b>	49.3	<b>58.6</b>	<b>54.3</b>
SAM [42] (Huge)	Full	632.18M	57.6	51.7	61.5	56.9
	Freeze	0.00M	57.0	47.1	58.4	54.2
	Rein	4.51M	<b>59.6</b>	<b>52.0</b>	<b>62.1</b>	<b>57.9</b>
EVA02 [18, 19] (Large)	Full	304.24M	62.1	56.2	64.6	60.9
	Freeze	0.00M	56.5	53.6	58.6	56.2
	Rein	2.99M	<b>65.3</b>	<b>60.5</b>	<b>64.9</b>	<b>63.6</b>
DINOv2 [58] (Large)	Full	304.20M	63.7	57.4	64.2	61.7
	Freeze	0.00M	63.3	56.1	63.9	61.1
	Rein	2.99M	<b>66.4</b>	<b>60.4</b>	<b>66.1</b>	<b>64.3</b>

Generalization for semantic segmentation

K. Chen et al., MultiSiam: Self-supervised Multi-instance Siamese Representation Learning for Autonomous Driving, ICCV 2021

Z. Peng et al., BEiT v2: Masked Image Modeling with Vector-Quantized Visual Tokenizers, arXiv 2022

Z. Wei et al., Stronger, Fewer, & Superior: Harnessing Vision Foundation Models for Domain Generalized Sem. Seg., CVPR 2024

**One model for all**

# **Pretraining Lidar encoders**

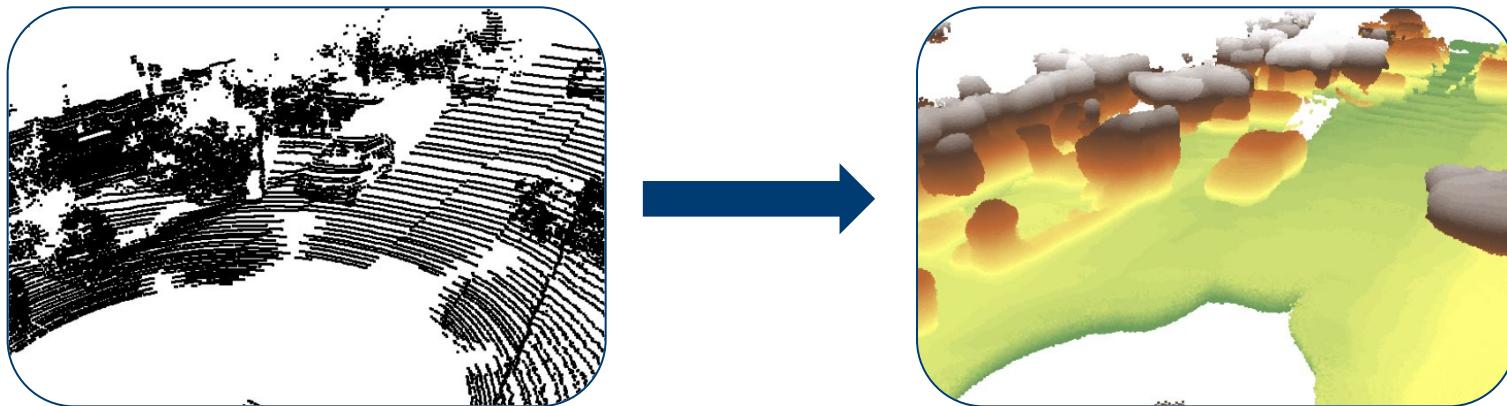
**The self-supervised 2-stage paradigm (pretrain + finetune)  
can be extend to other sensor modalities with some  
adequate pretext tasks**



# Automotive Lidar Self-Supervision by Occupancy Estimation

Intuition

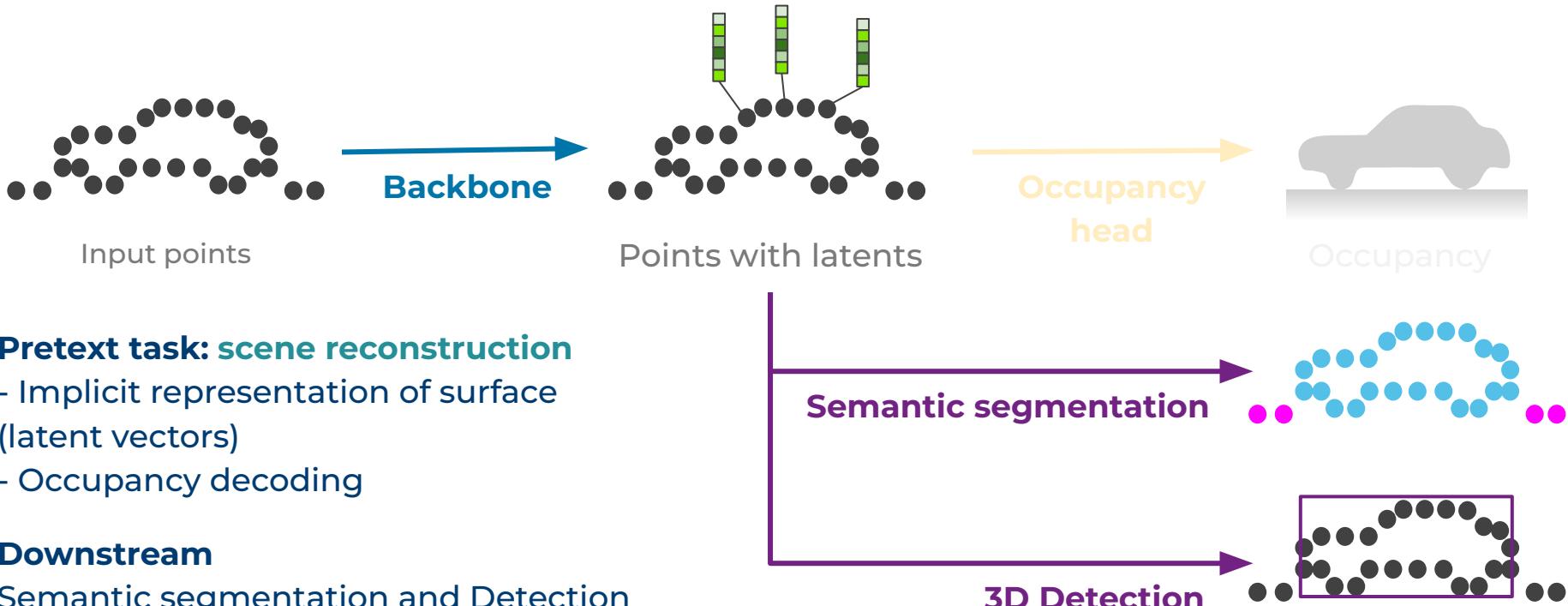
- Reconstruction information (geometry) requires semantic knowledge
- Point cloud ~ masked surface





# Automotive Lidar Self-Supervision by Occupancy Estimation

Overview



## Pretext task: scene reconstruction

- Implicit representation of surface (latent vectors)
- Occupancy decoding

## Downstream

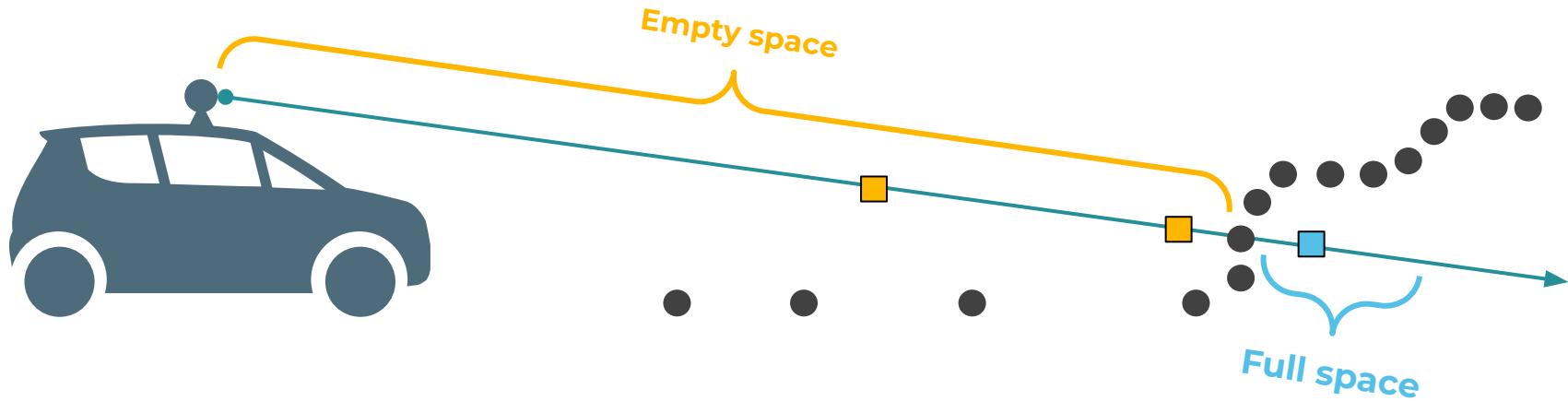
Semantic segmentation and Detection

---



# Automotive Lidar Self-Supervision by Occupancy Estimation

Pretext task



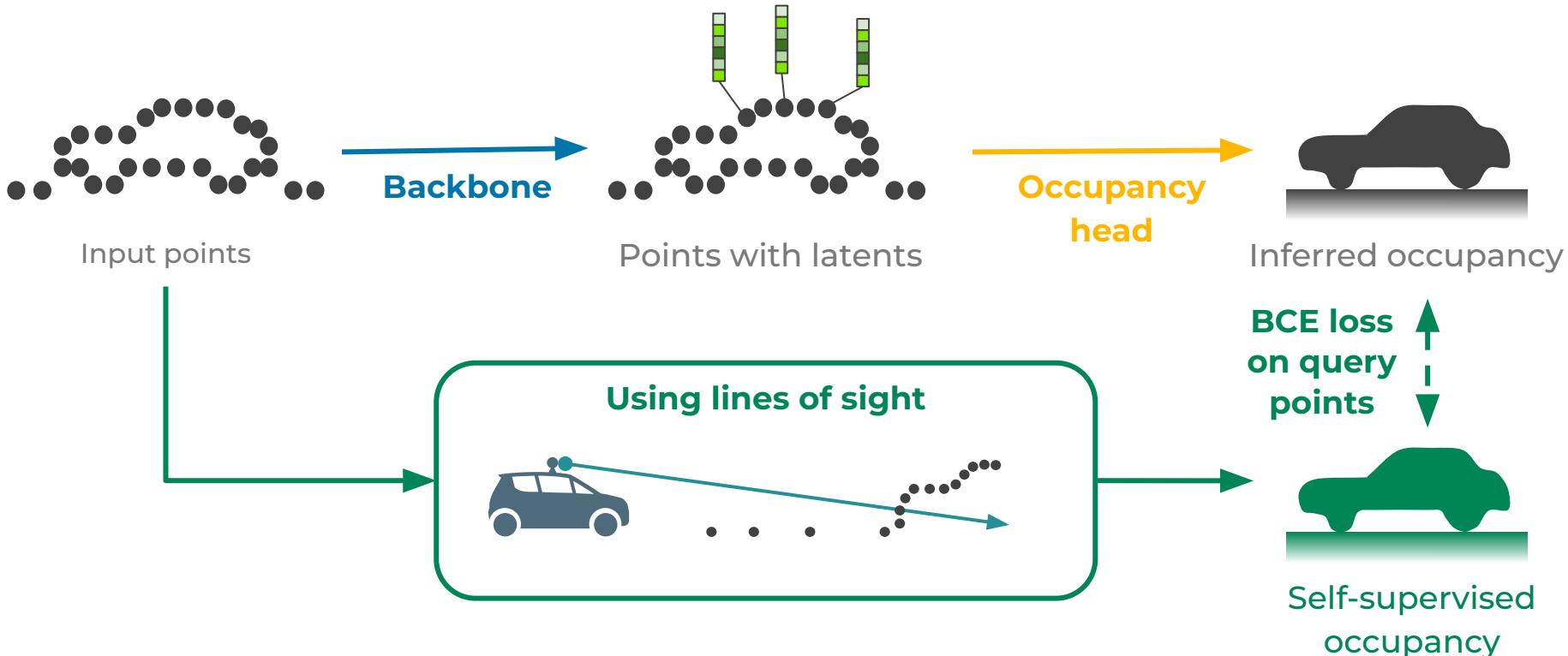
## Along lidar lines of sight

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



# Automotive Lidar Self-Supervision by Occupancy Estimation

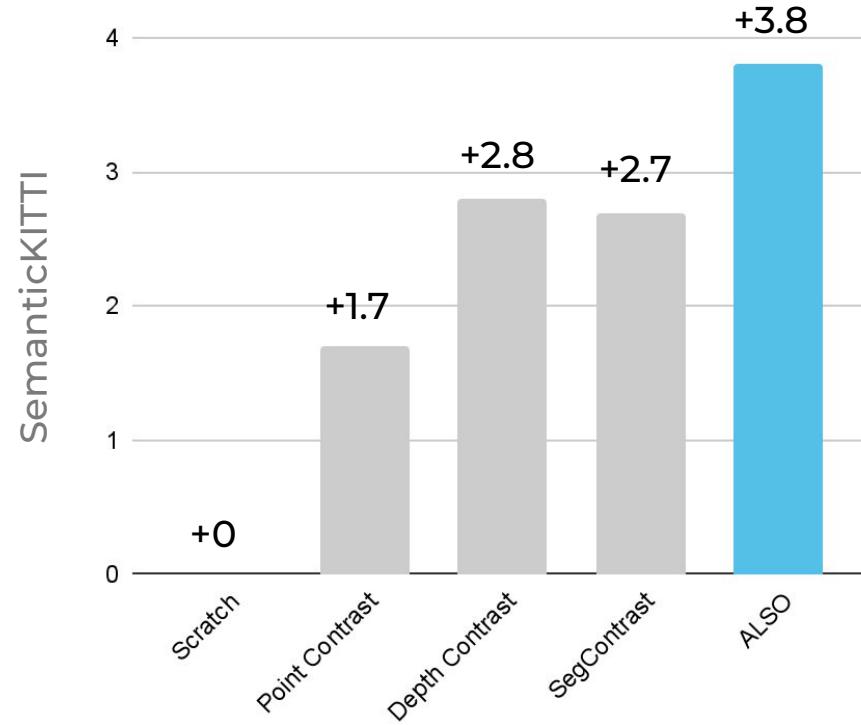
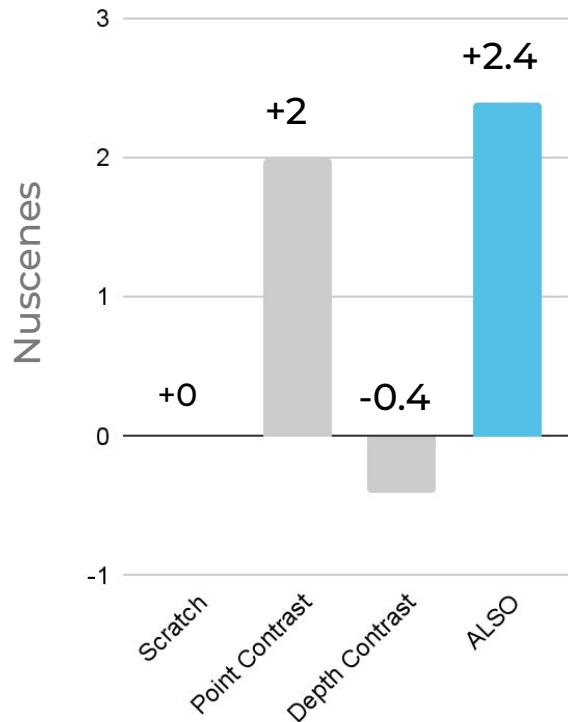
Pretext task





# Downstream performance

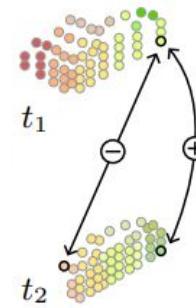
1% annotated training data





# Contrastive self-supervised learning for point clouds

BEVContrast



PointContrast

## PointContrast

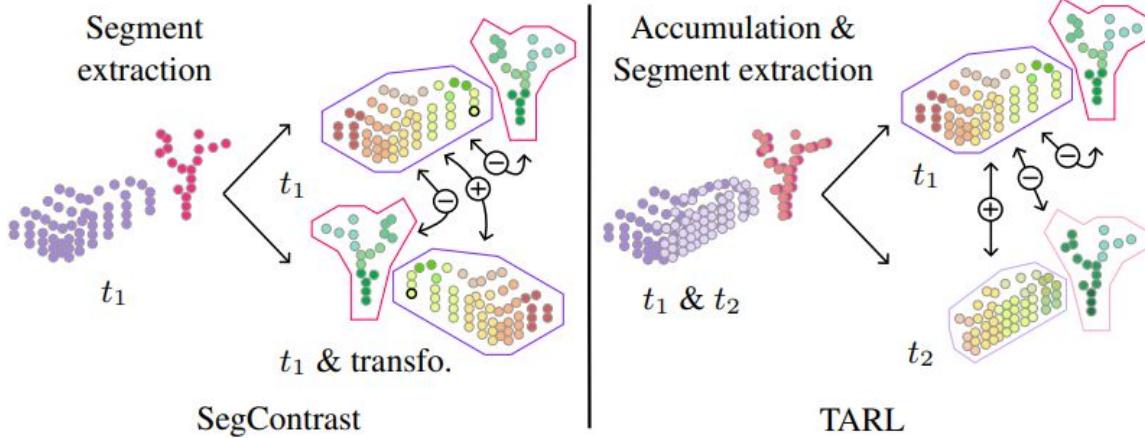
- +++ Easy to implement
- - - Contrast in the same object

---



# Contrastive self-supervised learning for point clouds

BEVContrast



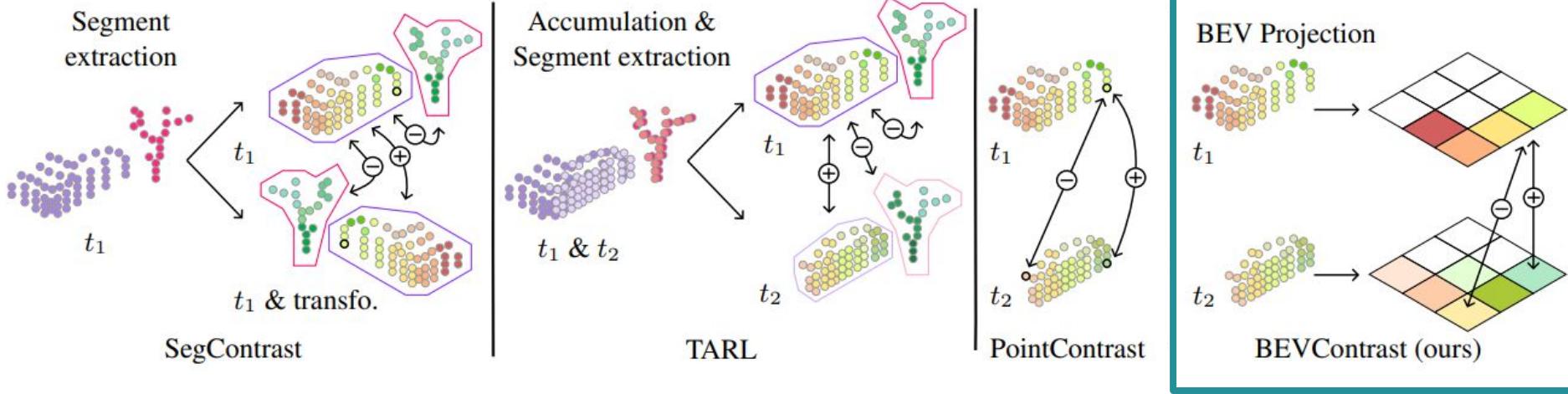
## SegContrast / TARL

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



# Contrastive self-supervised learning for point clouds

BEVContrast



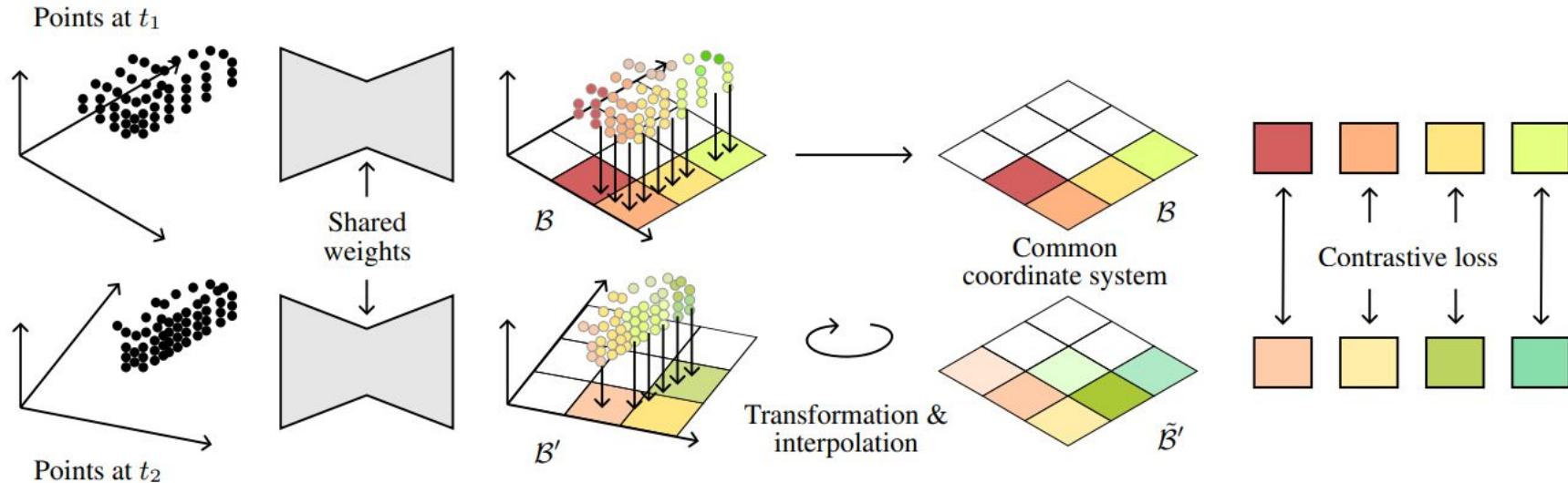
## BEVContrast

- +++ Simple → easy projection in BEV
- +++ Object separation approximation with BEV cells



# Contrastive self-supervised learning for point clouds

BEVContrast



+++ Simple → easy projection in BEV

+++ Object separation approximation with BEV cells

# BEVContrast



Semantic segmentation: qualitative results

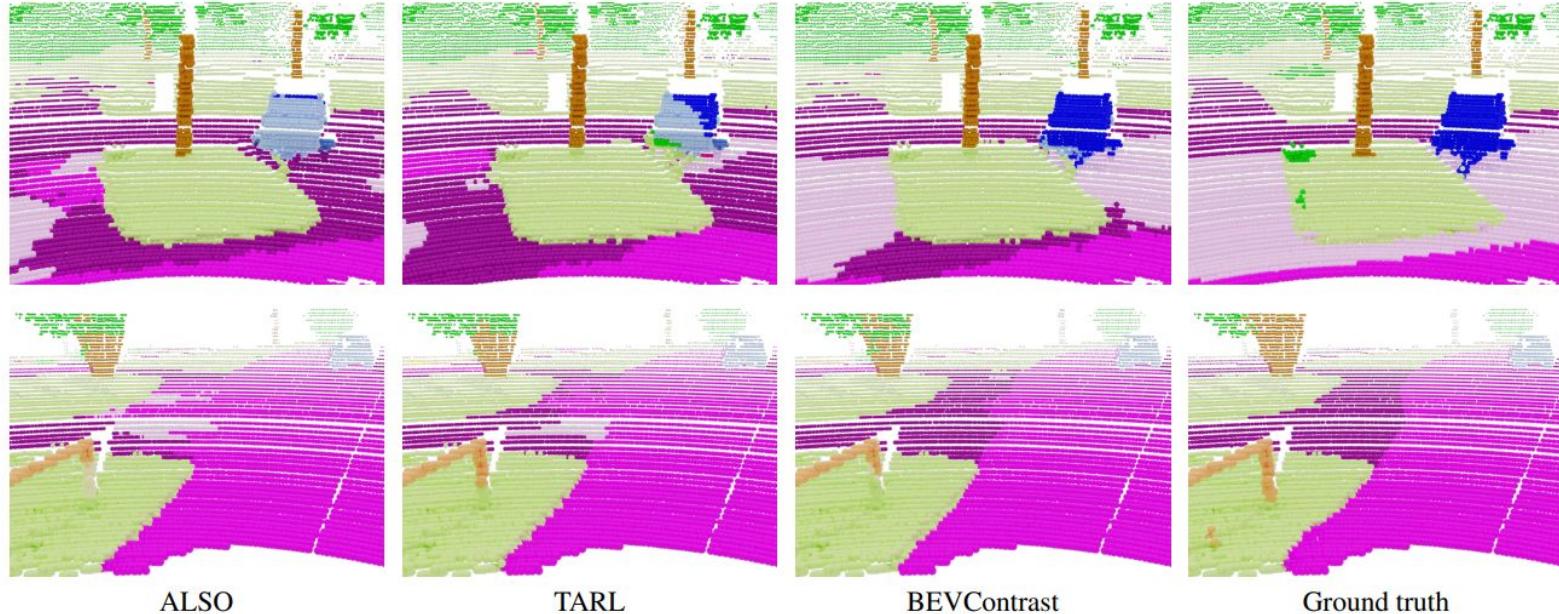


Figure 3. Semantic segmentation visualizations on SemanticKITTI after pre-training on the full training dataset and fine-tuning using 1% of the annotated data.

car	other veh.	sidewalk	drivable surf.	parking	terrain	vegetation	trunk	fence
-----	------------	----------	----------------	---------	---------	------------	-------	-------

# BEVContrast



Semantic segmentation: quantitative results

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

**One model from many for all**  
**Multi-modal training**

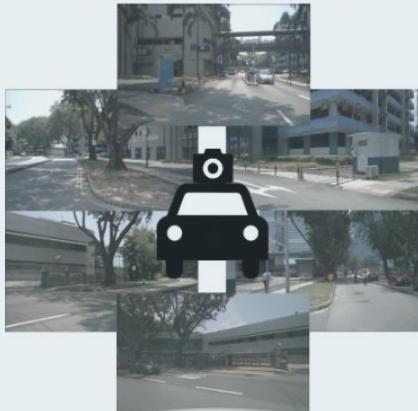
# Lidar-to-Image models

# Open-Vocabulary 3D Occupancy Prediction from Images

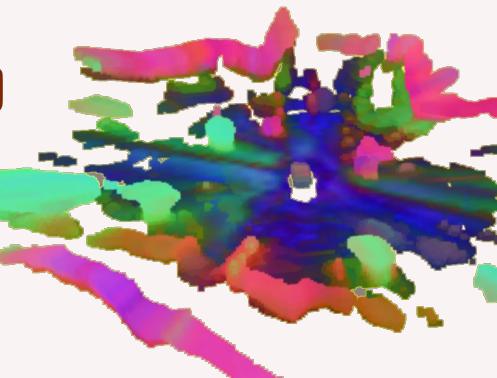


Goals

**INPUT:**  
surround  
-view  
images

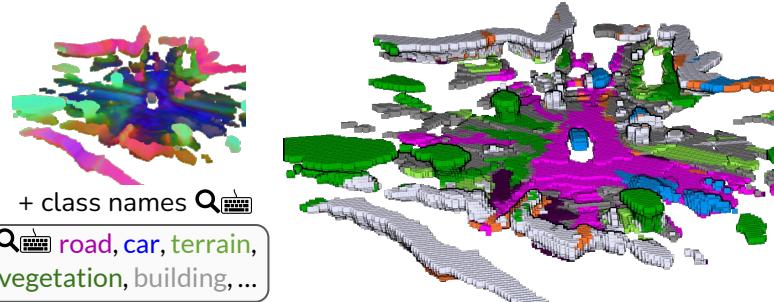


**POP<sub>3D</sub>**

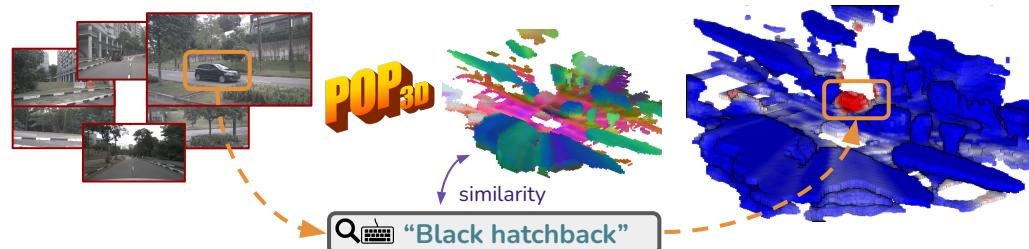


**OUTPUT:**  
**3D voxel field with:**  
- occupancy  
- open-vocabulary  
features

**TASK #1:** zero-shot semantic occupancy segmentation



**TASK #2:** text-driven 3D retrieval from cameras

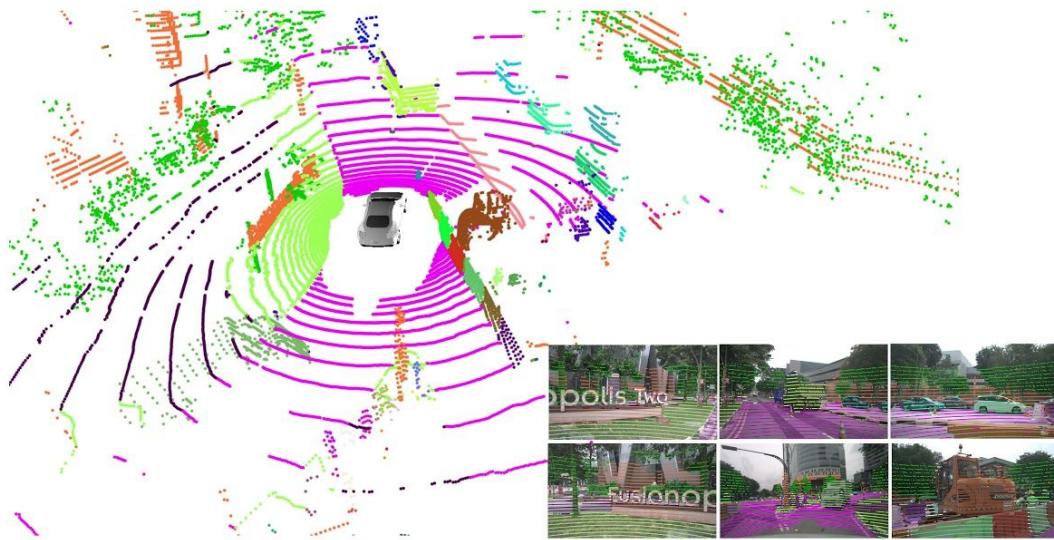


# Open-Vocabulary 3D Occupancy Prediction from Images



## Motivation

- Only sparse LiDAR annotations available (costly and difficult to scale).
- Annotated labels for a predefined closed-vocabulary.





# Open-Vocabulary 3D Occupancy Prediction from Images

Overview

**open-vocabulary** 3D semantic occupancy prediction

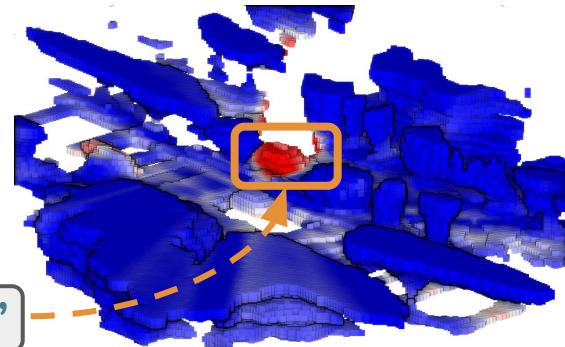
*Training:*

**unlabeled image-LiDAR data** and a **pre-trained image-language model**

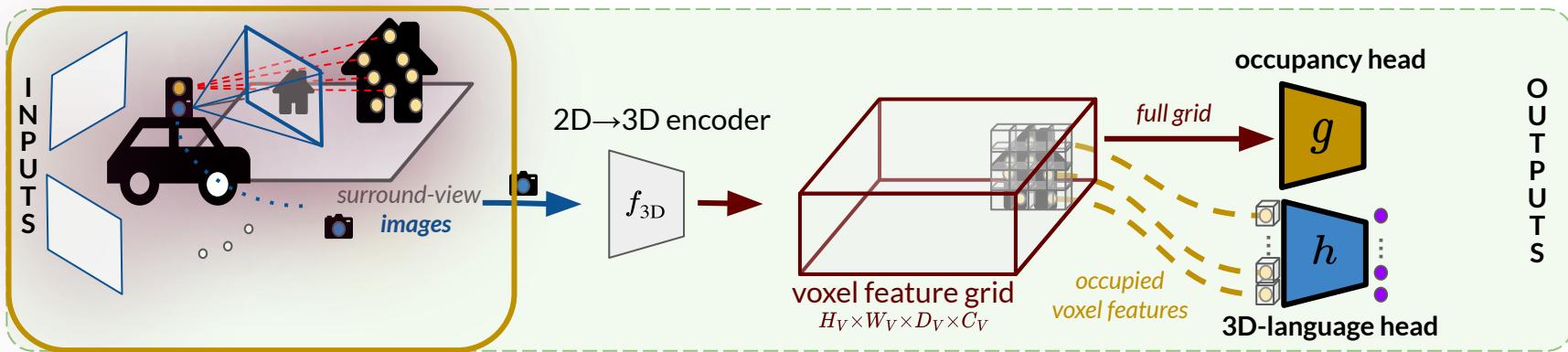
*Inference:*

**only images + text queries**

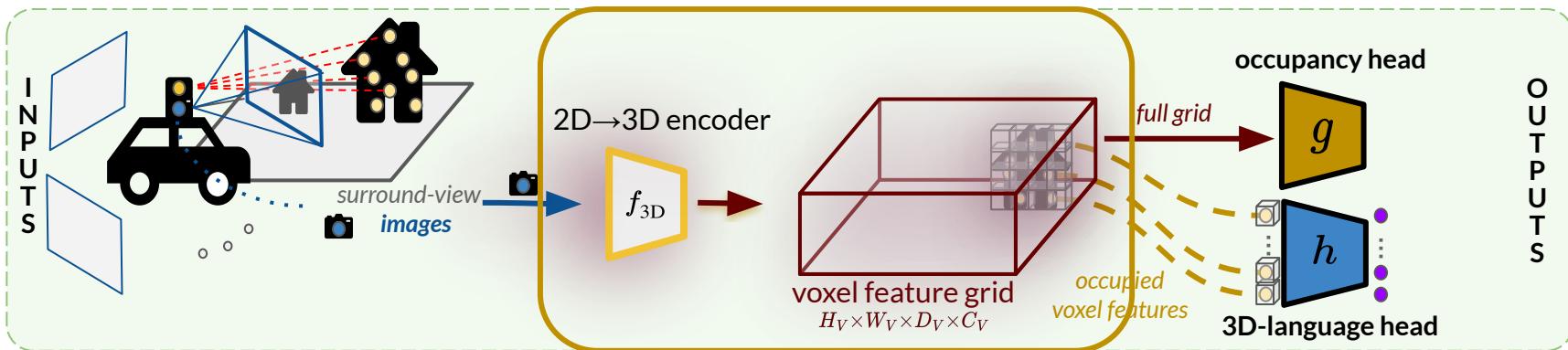
*Example of text-driven retrieval*



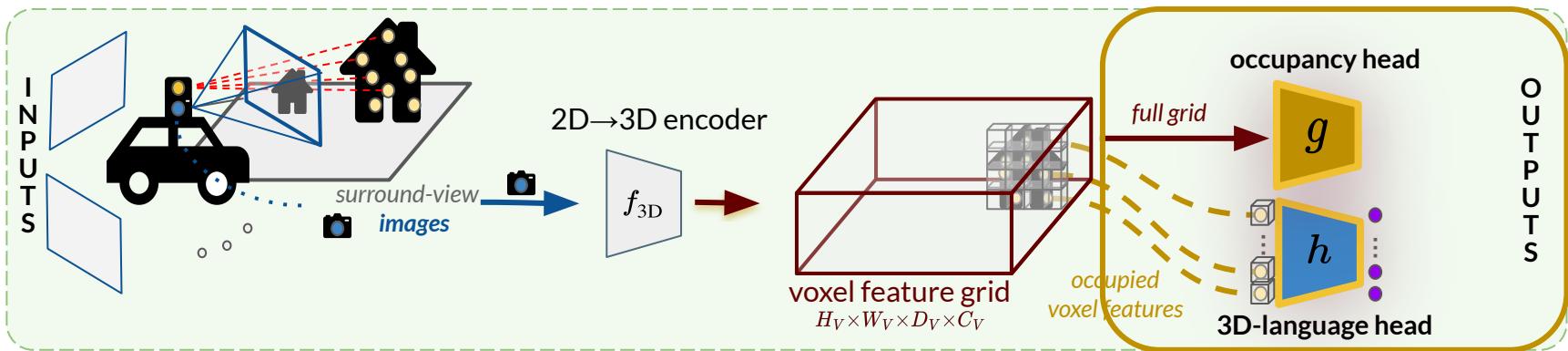
# POP-3D Architecture



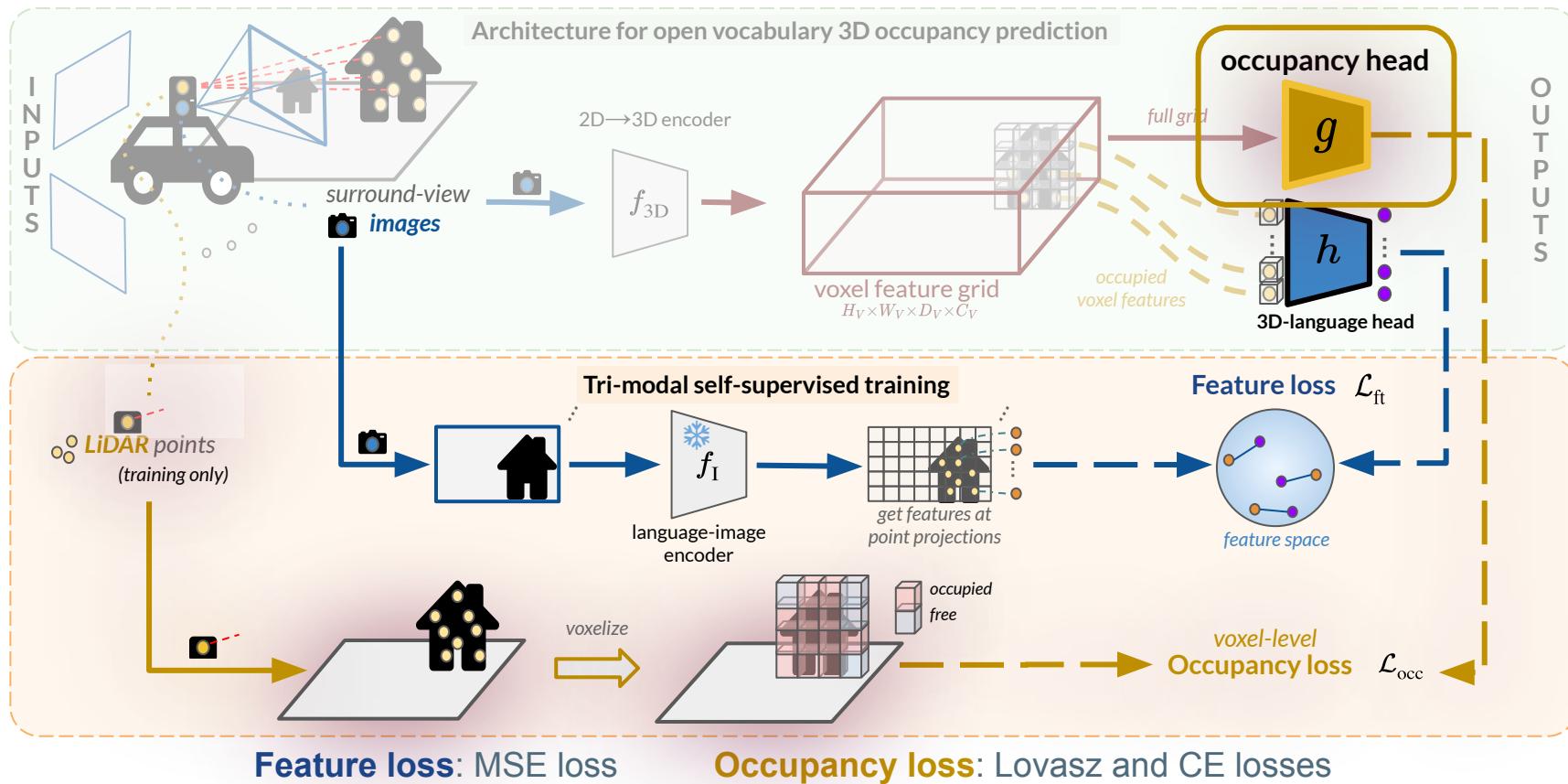
# POP-3D Architecture



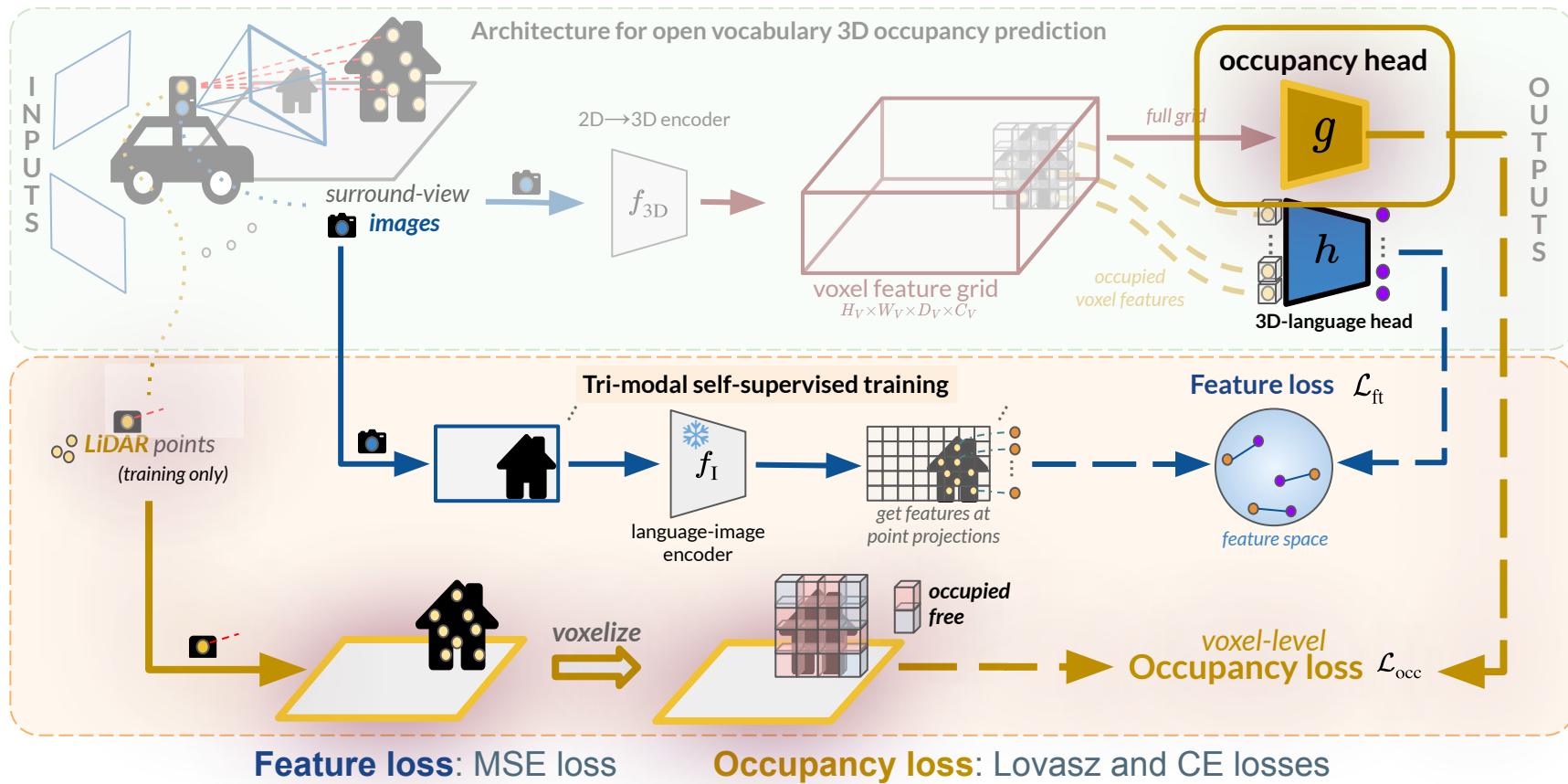
# POP-3D Architecture



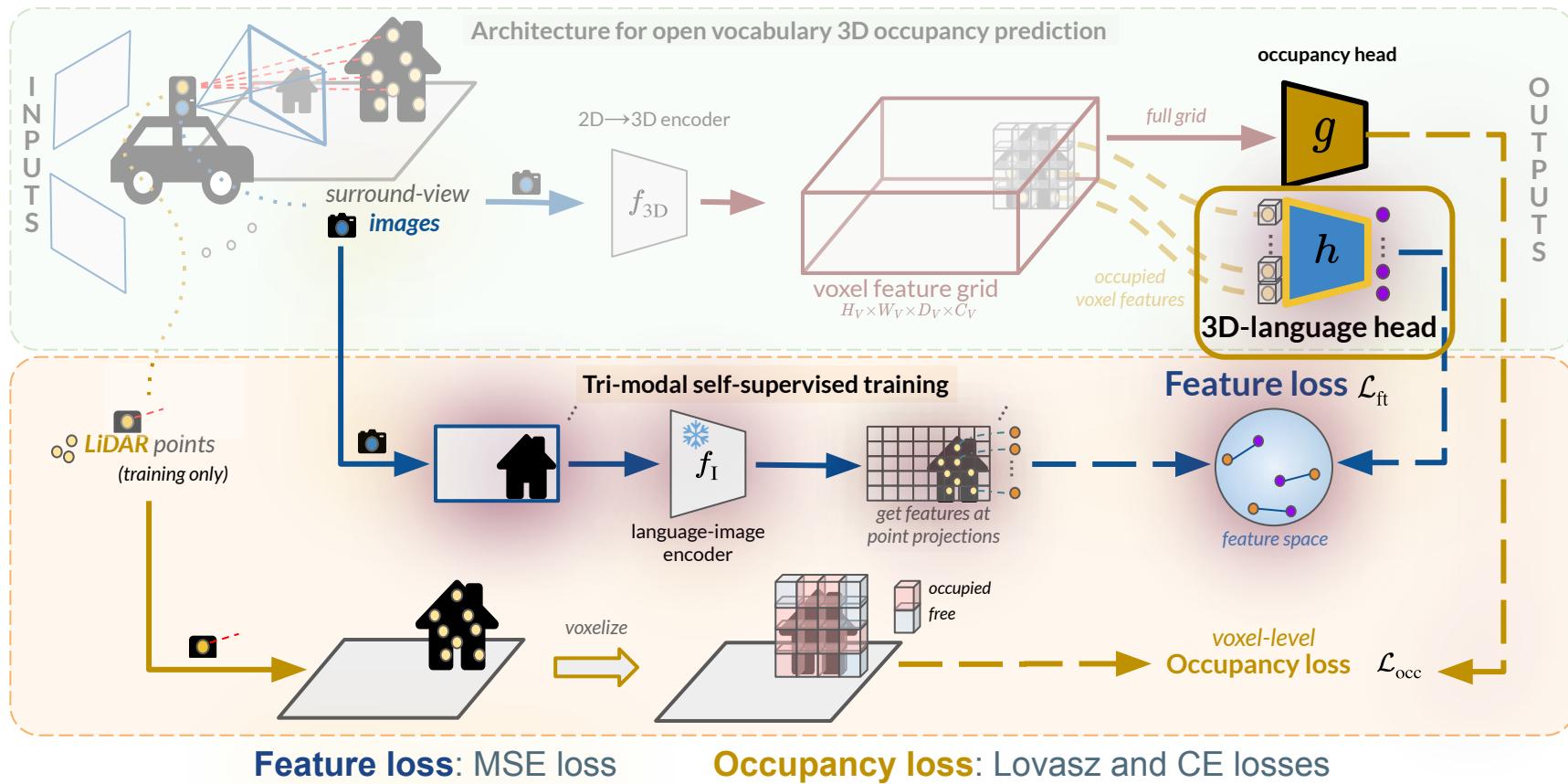
# POP-3D Losses



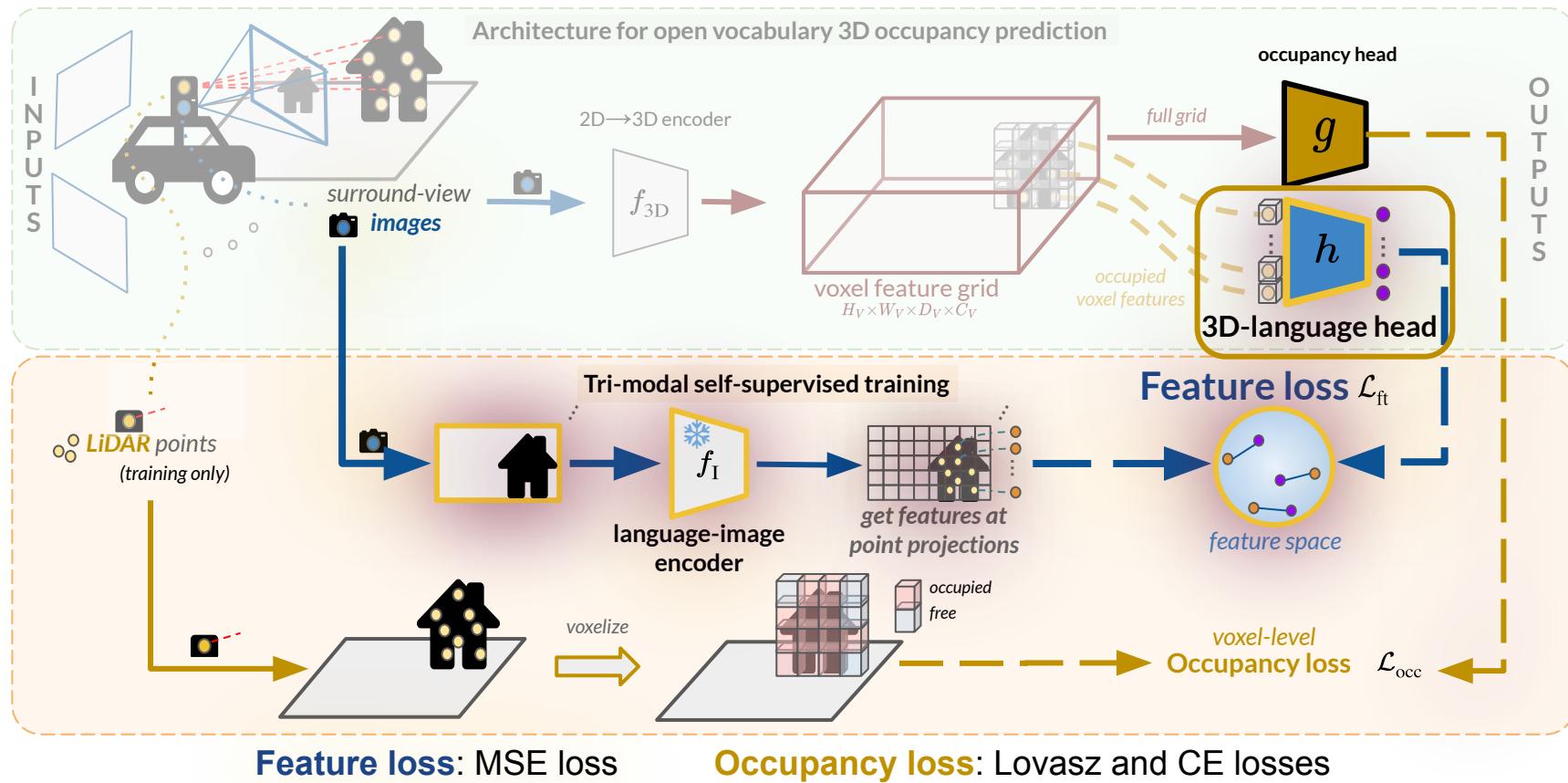
# POP-3D Losses



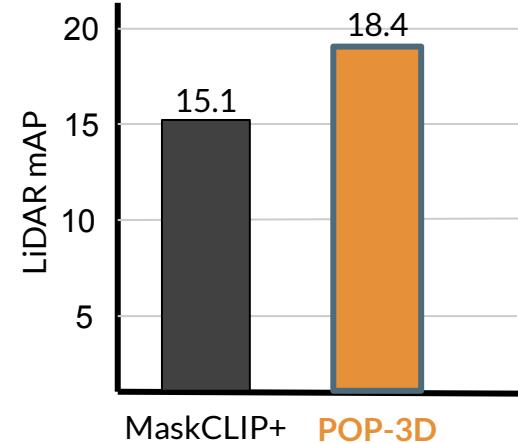
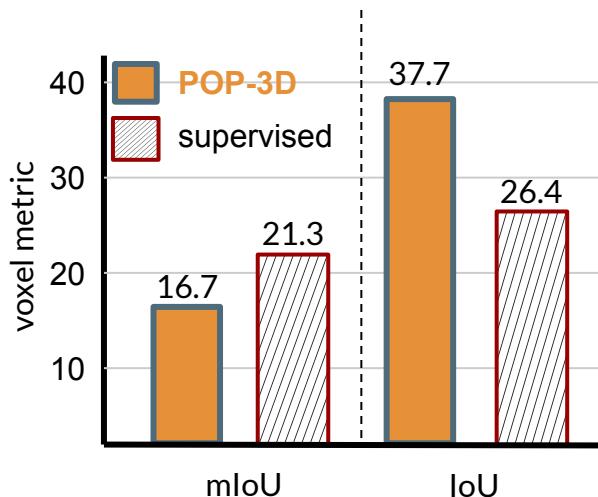
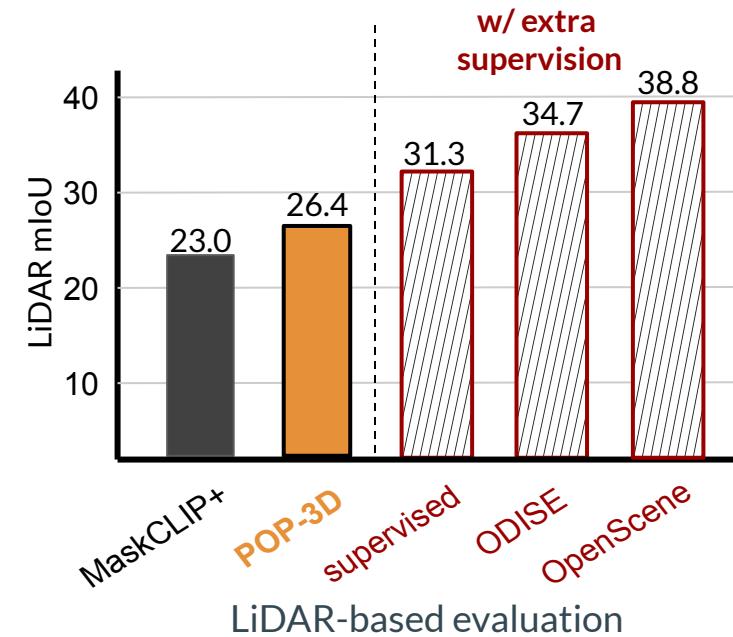
# POP-3D Losses



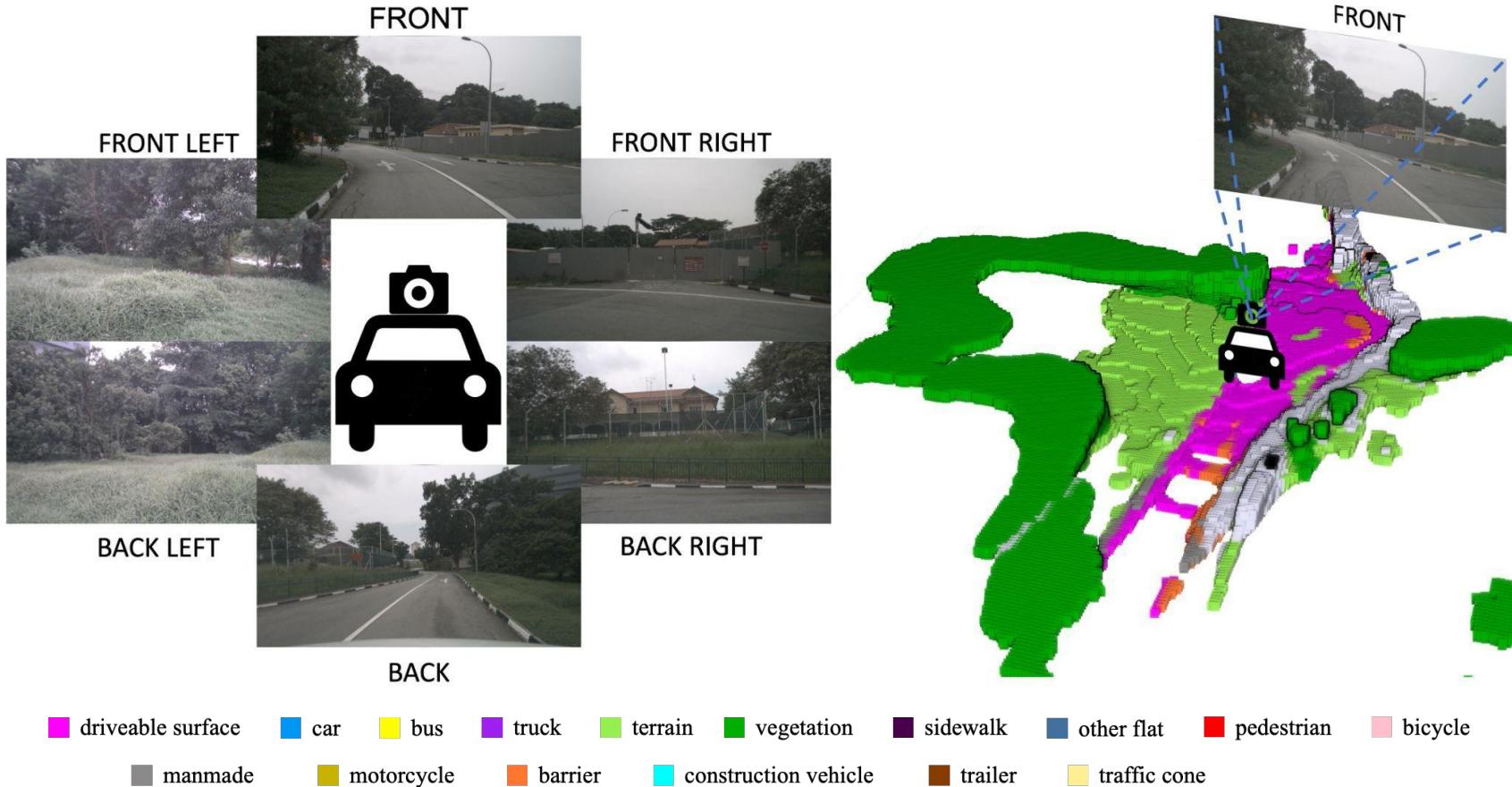
# POP-3D Losses



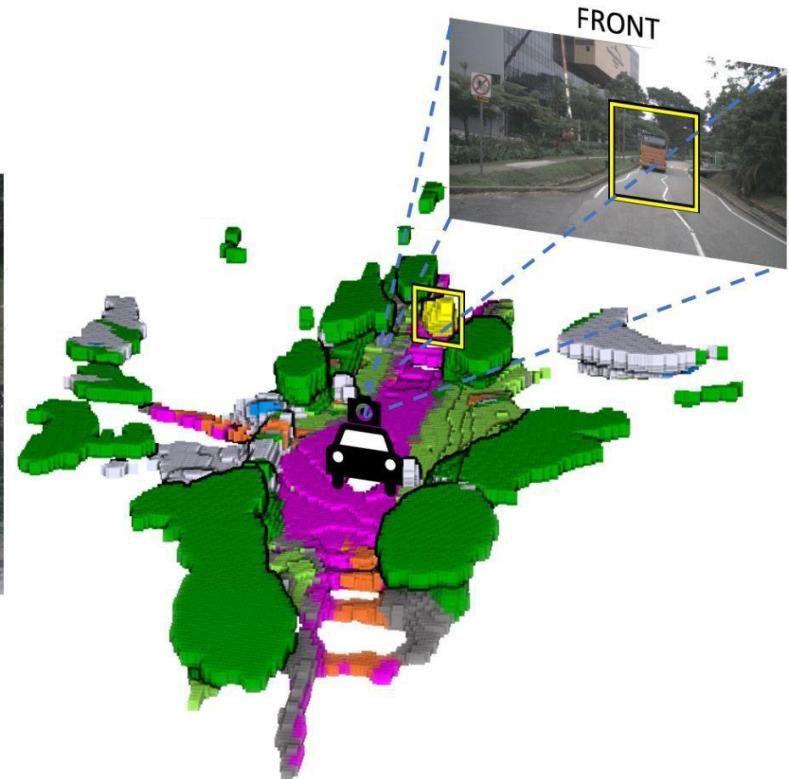
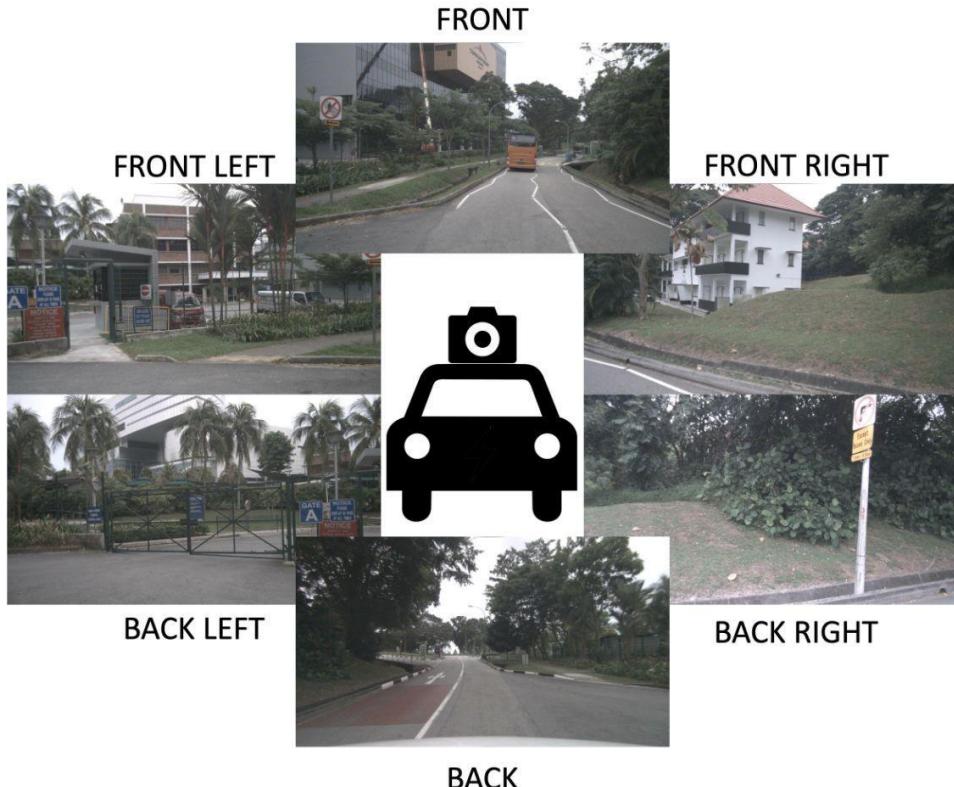
# POP-3D Quantitative results (nuscenes)



# POP-3D Qualitative results: zero-shot semantic segmentation

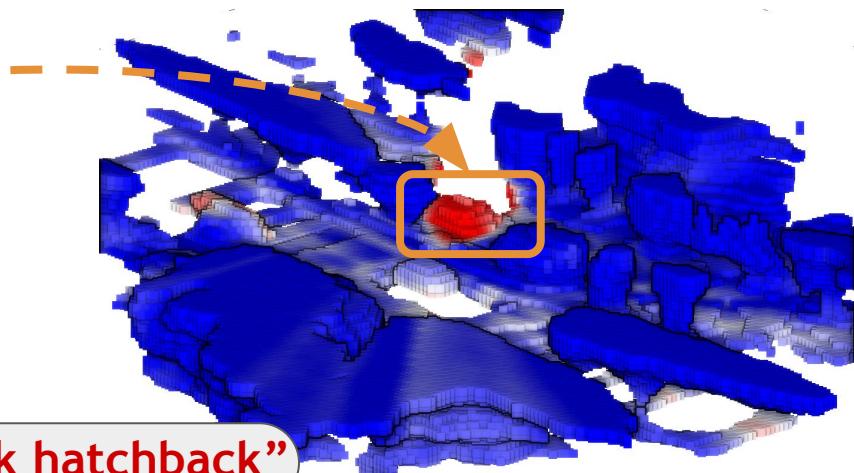


# POP-3D Qualitative results: zero-shot semantic segmentation



- driveable surface   ■ car   ■ bus   ■ truck   ■ terrain   ■ vegetation   ■ sidewalk   ■ other flat   ■ pedestrian   ■ bicycle
- manmade   ■ motorcycle   ■ barrier   ■ construction vehicle   ■ trailer   ■ traffic cone

# POP-3D Qualitative results: retrieval

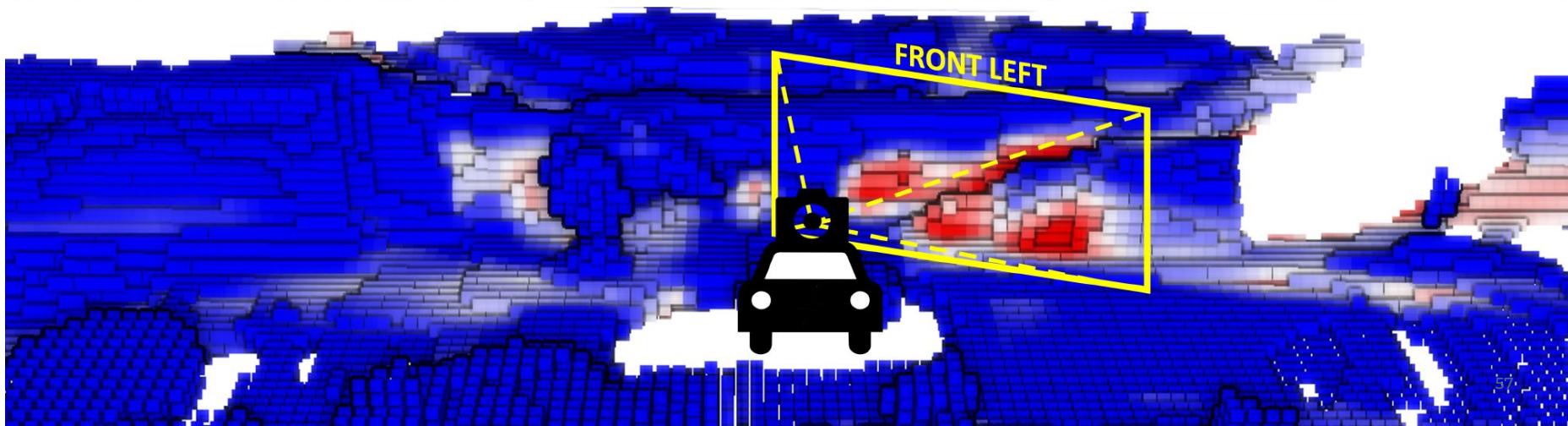


🔍⌨️: “Black hatchback”

# POP-3D Qualitative results: retrieval



🔍💻 “stairs”

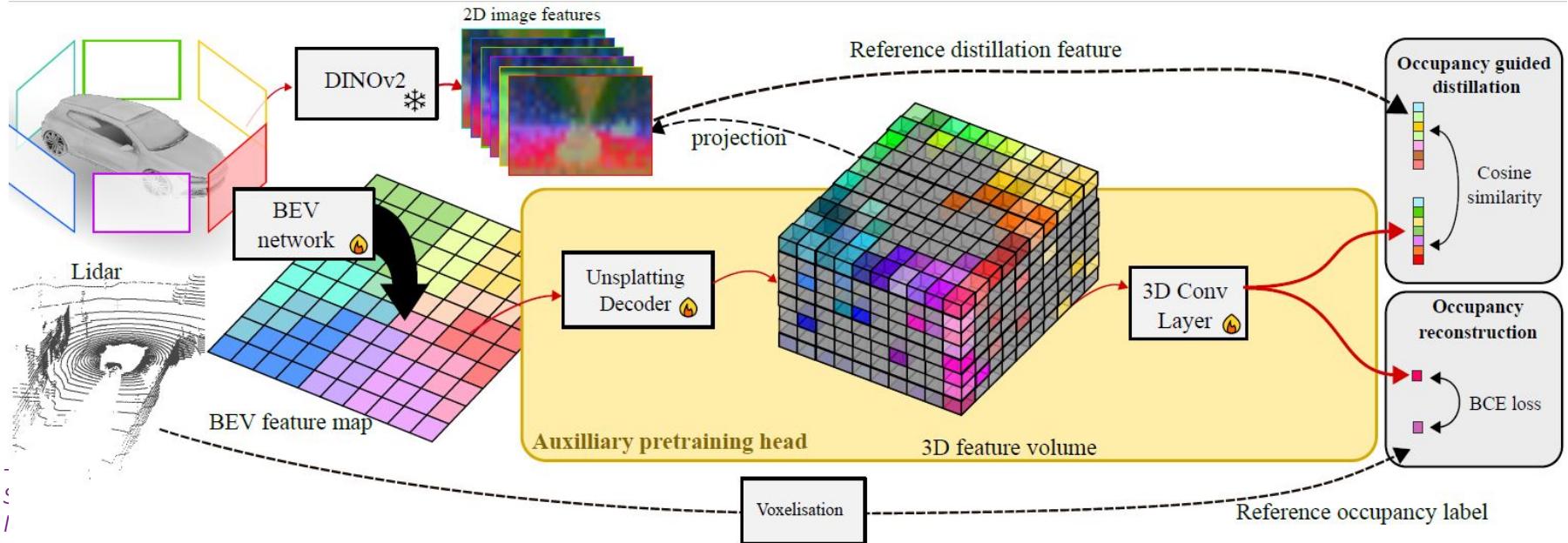


# Self-supervised Occupancy Feature Prediction for Pretraining BEV



OccFeat

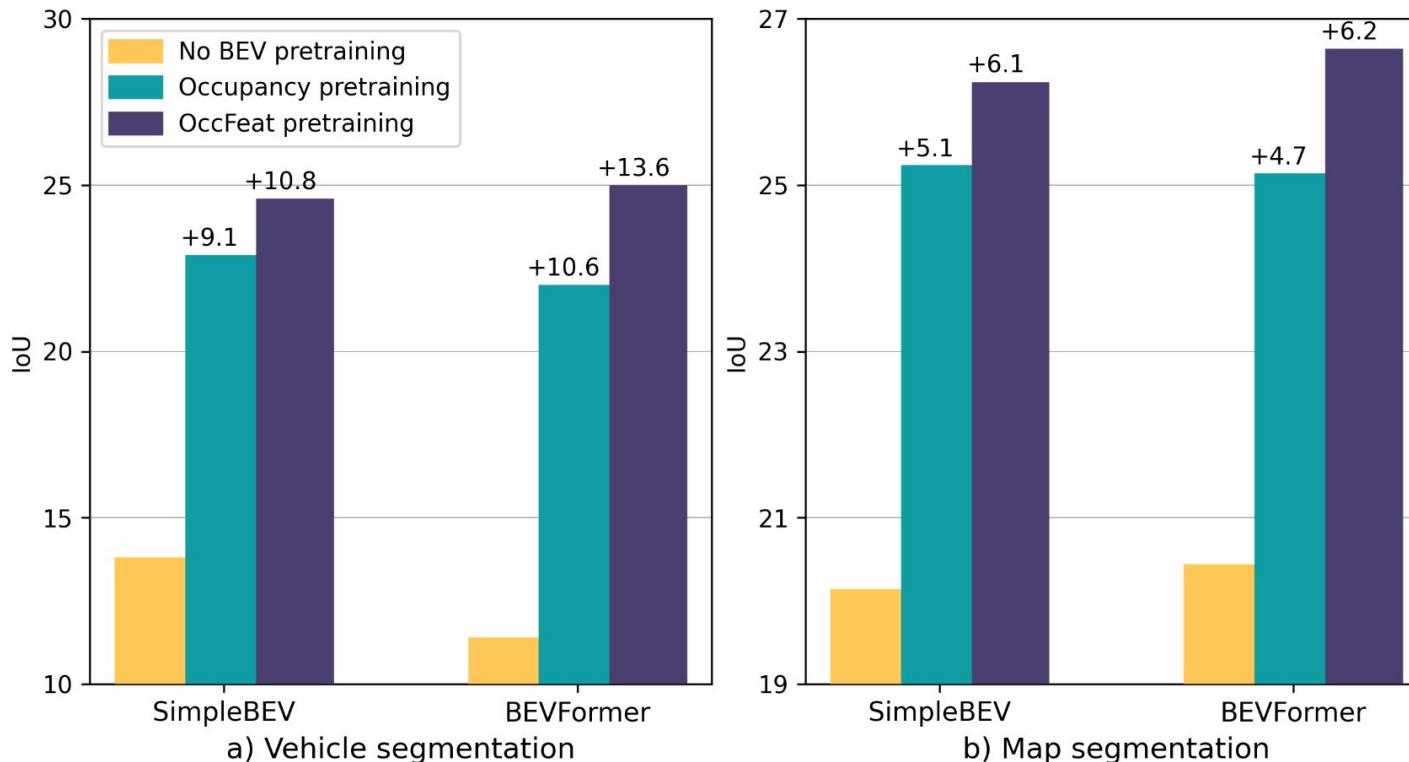
- Architecture-agnostic pretraining that teaches the model semantic-aware 3D geometry information
- Leveraging aligned Image & LiDAR data pretrained SSL image encoder (DINOv2)





# Downstream performance - 1% labelled samples (nuscenes)

OccFeat

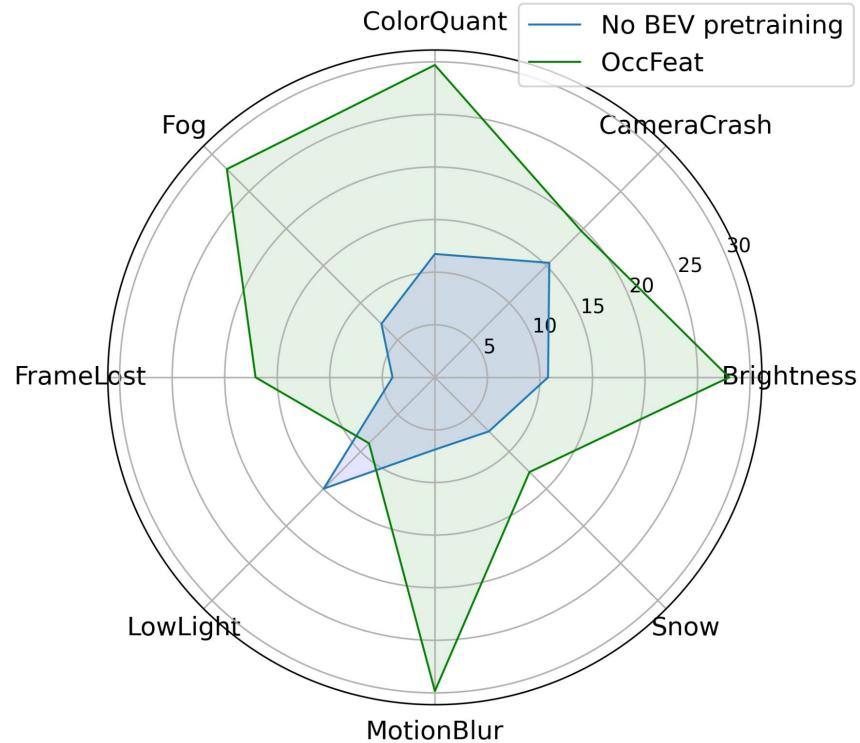




# Robustness to distribution shift

OccFeat

- Segmentation results on nuScenes-C dataset (8 corruptions on nuScenes validation set)
- OccFeat vs. no BEV pretraining



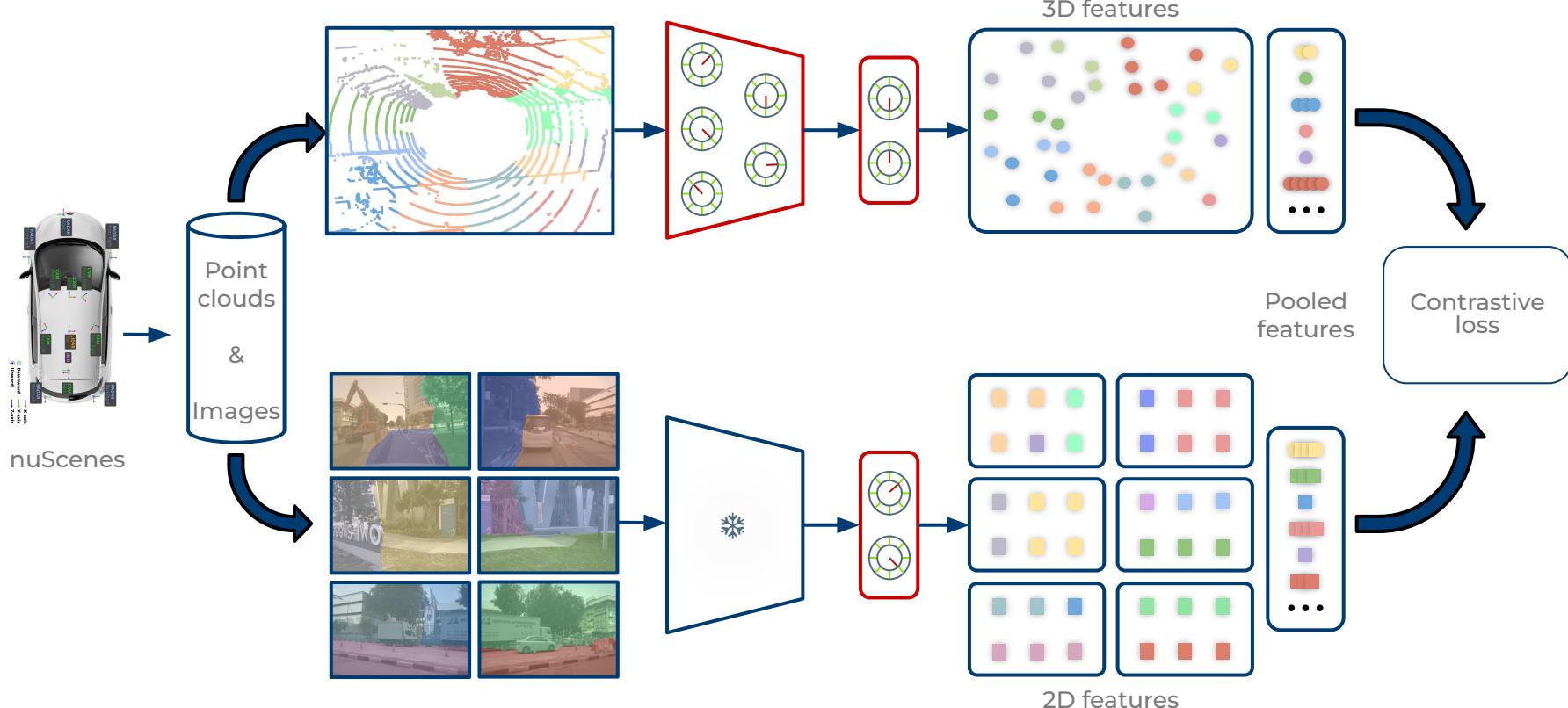
---

# Image-to-Lidar models



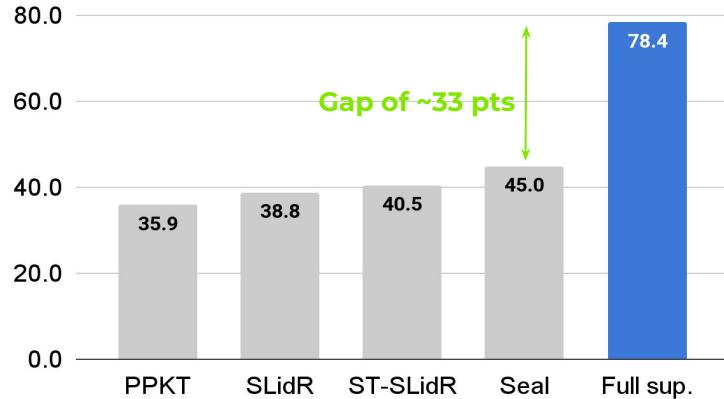
# Image-to-Lidar distillation

How SoTA methods work





# Image-to-Lidar distillation

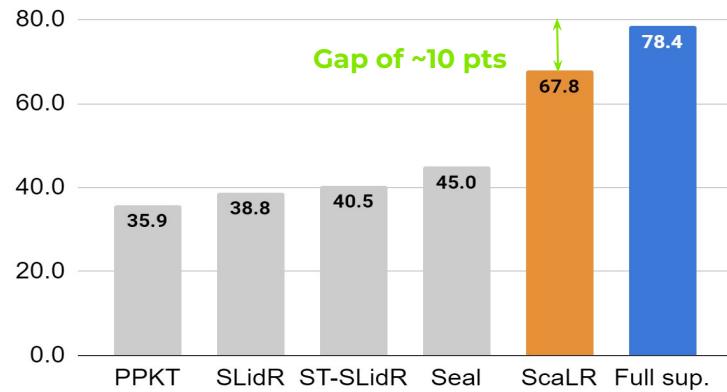
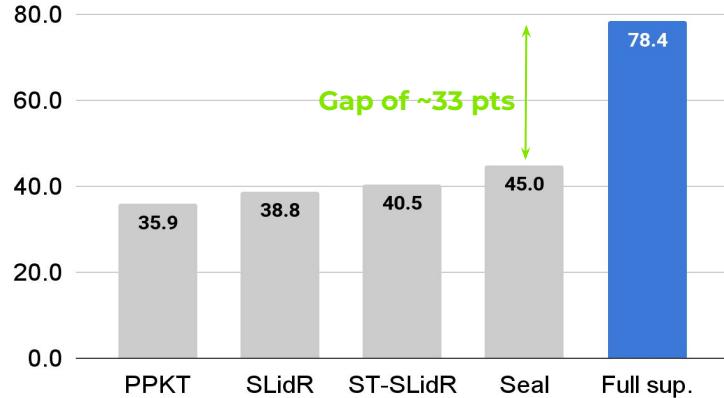


## Motivation:

Distilled 3D features << supervised 3D features



# Image-to-Lidar distillation



## Motivation:

Distilled 3D features << supervised 3D features

## Main result:

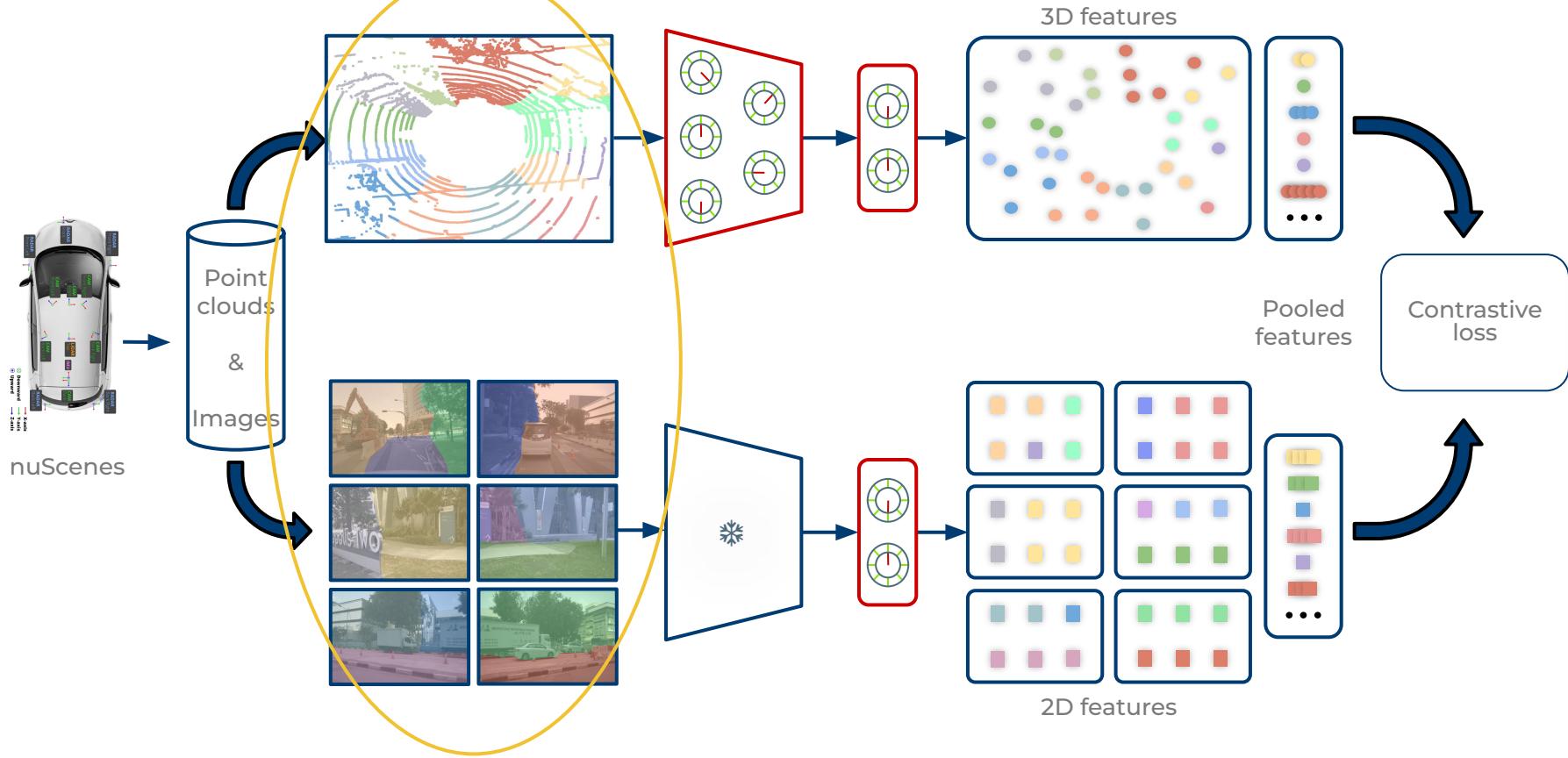
Large improvement of feature quality: +22.8 mIoU pts



**With ScaLR we revisit this pipeline and propose the following simplifications:**

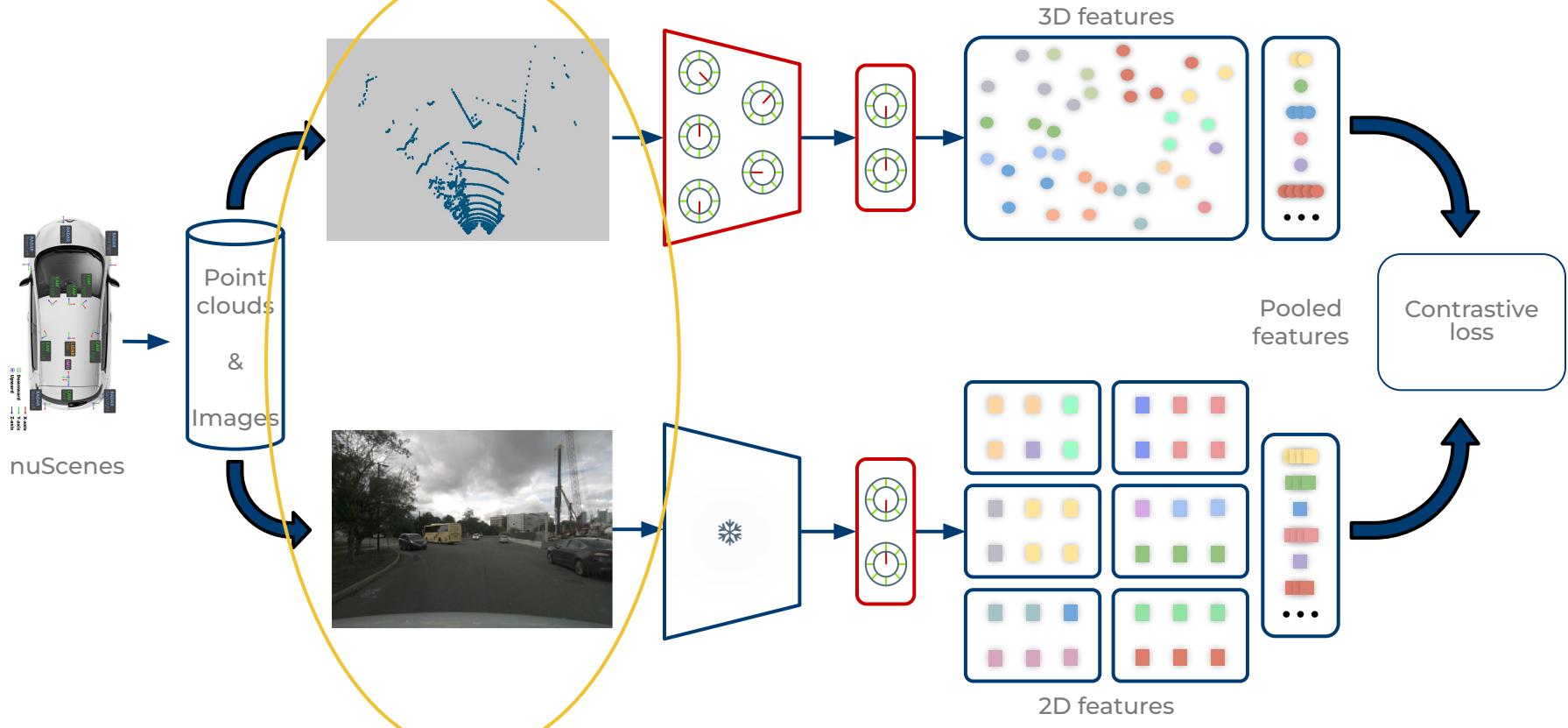


# Simplifications done in ScaLR



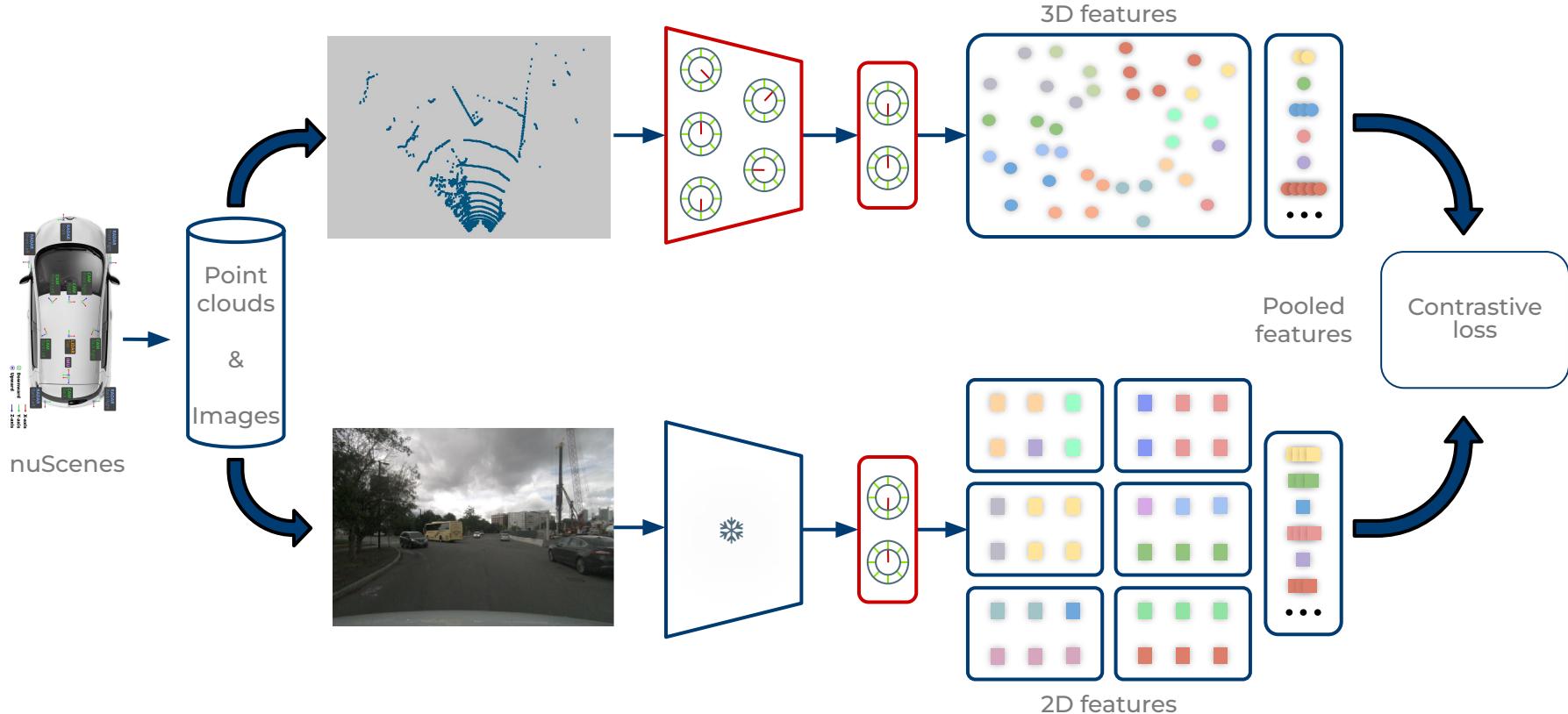


# Simplifications done in ScaLR



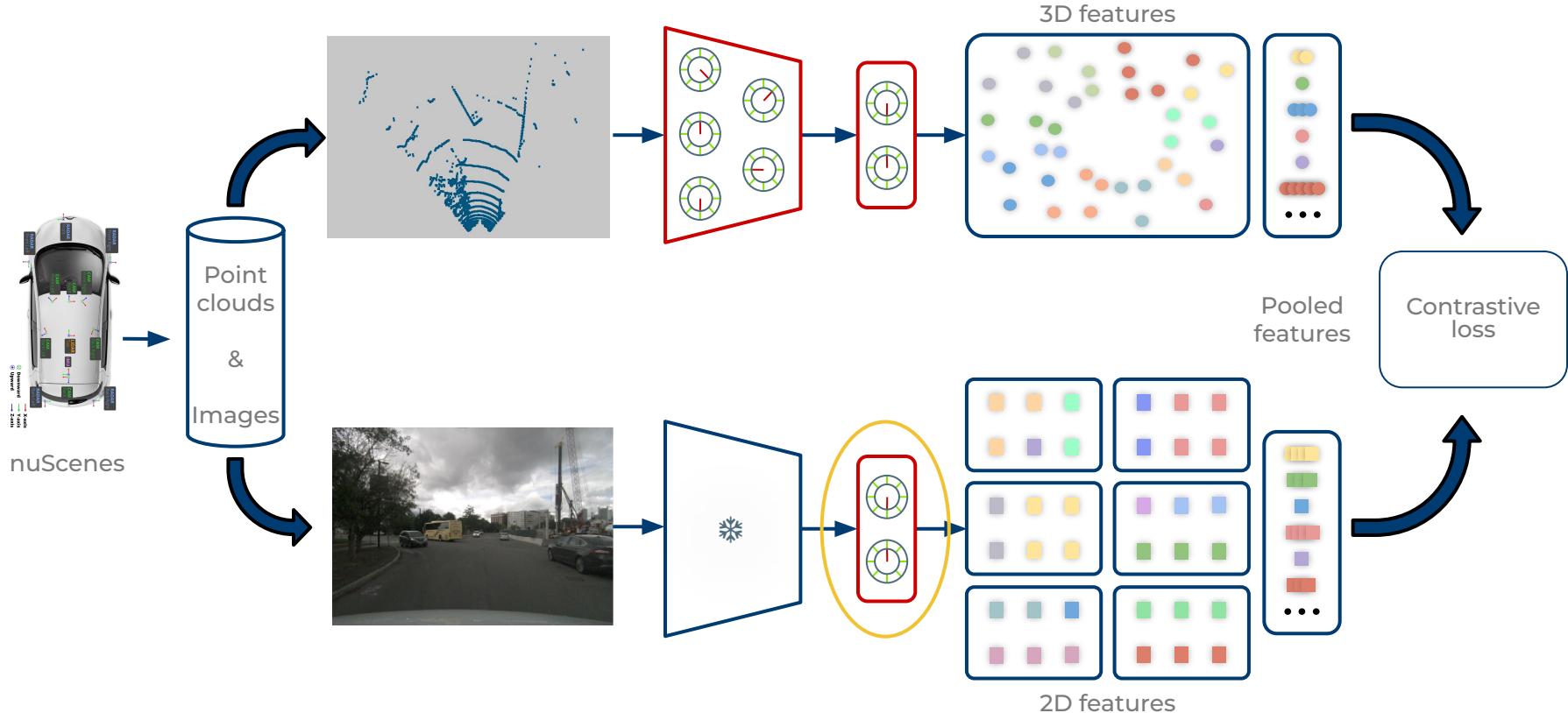


# Simplifications done in ScaLR



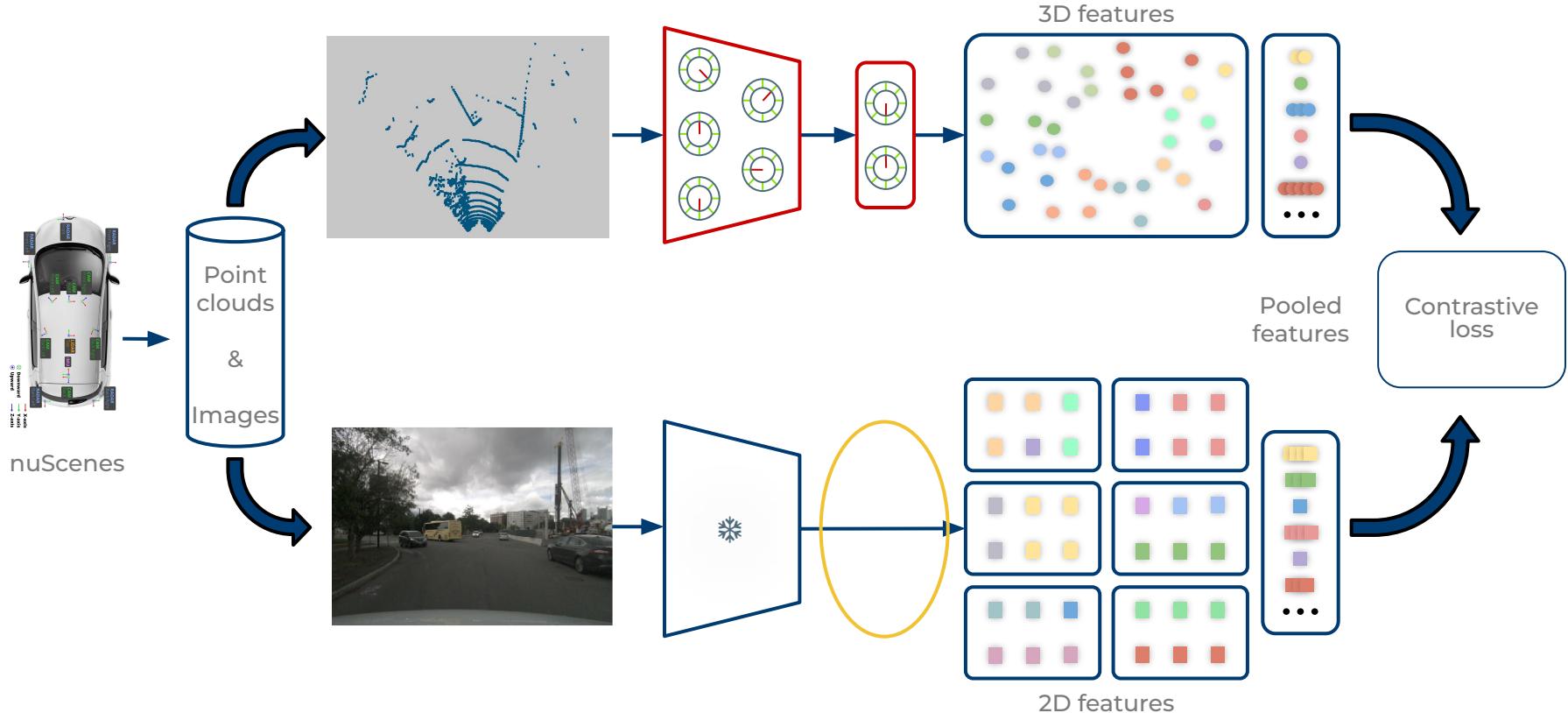


# Simplifications done in ScaLR



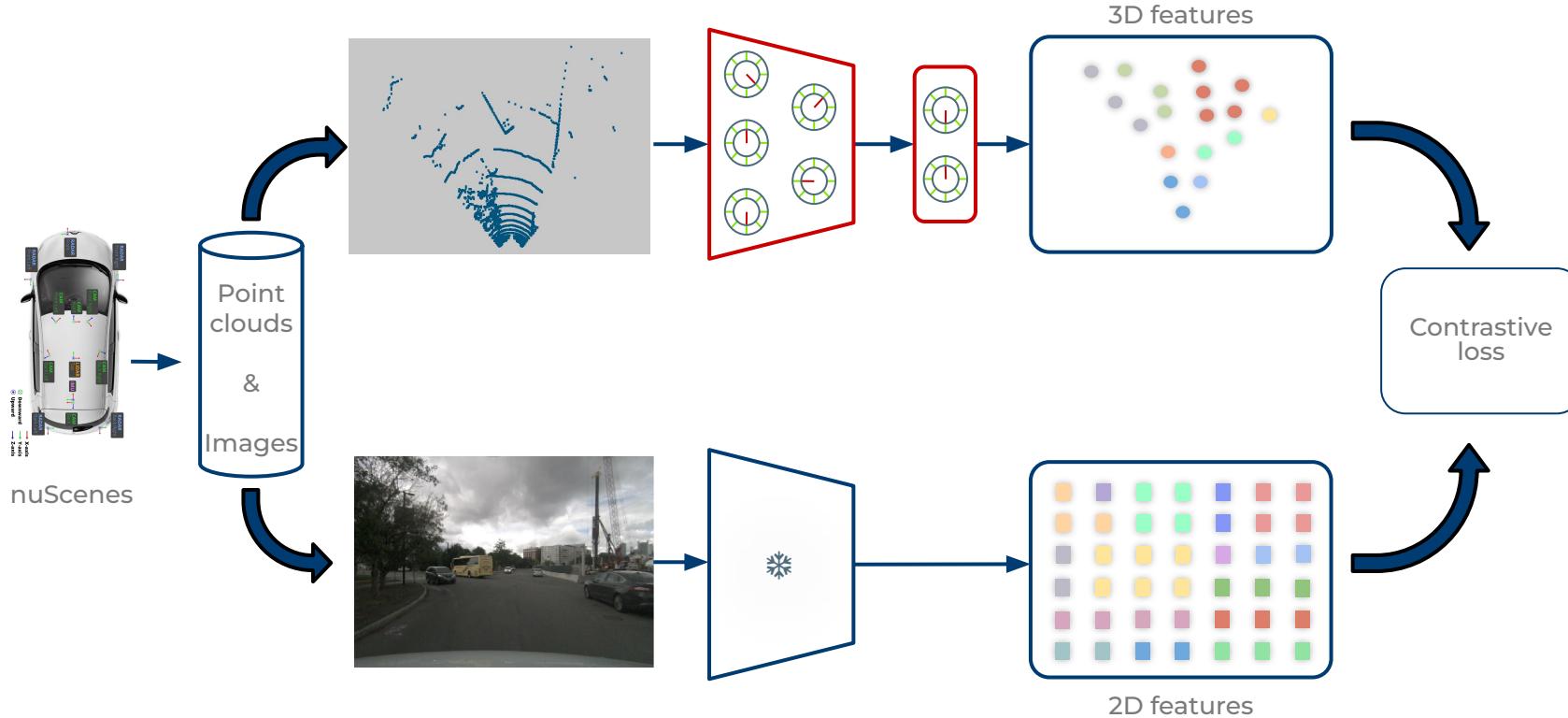


# Simplifications done in ScaLR



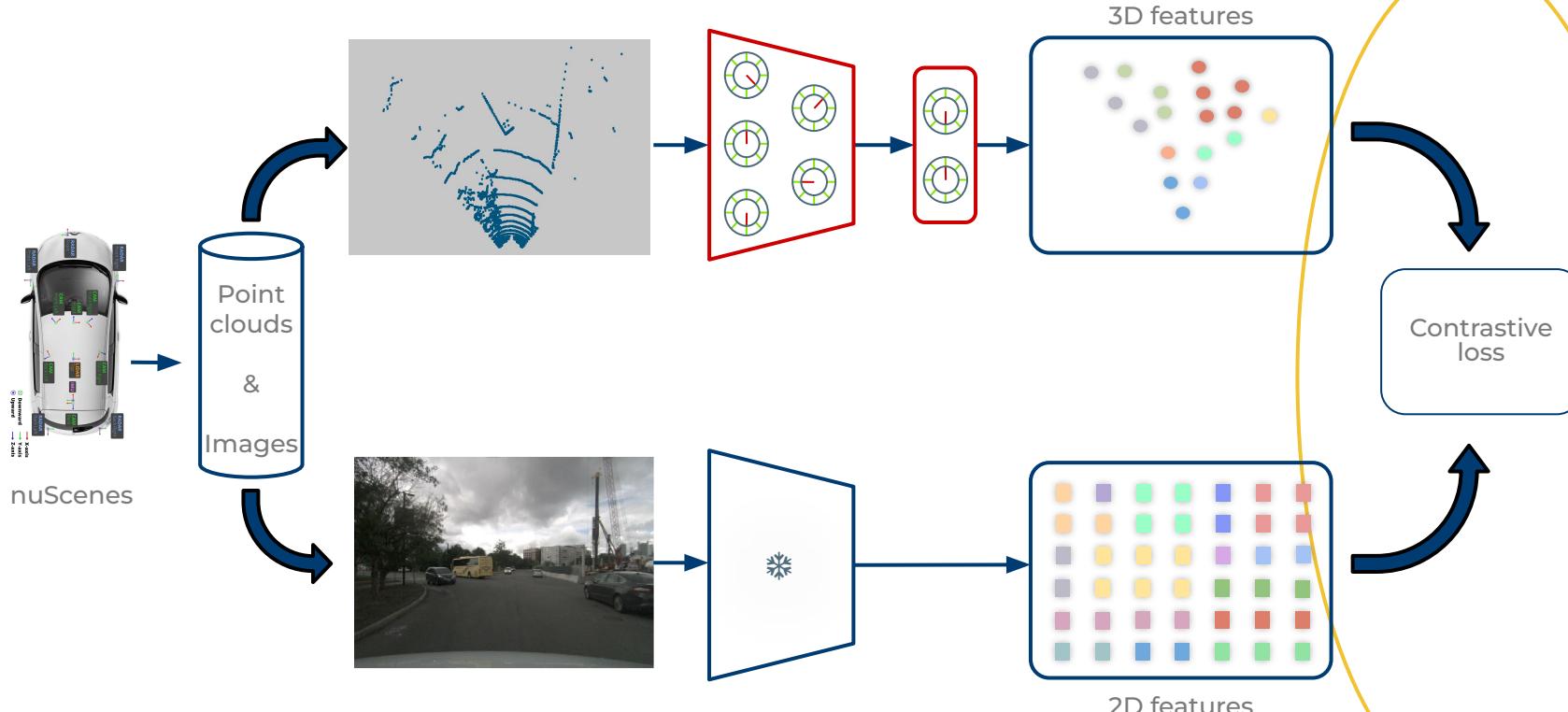


# Simplifications done in ScaLR



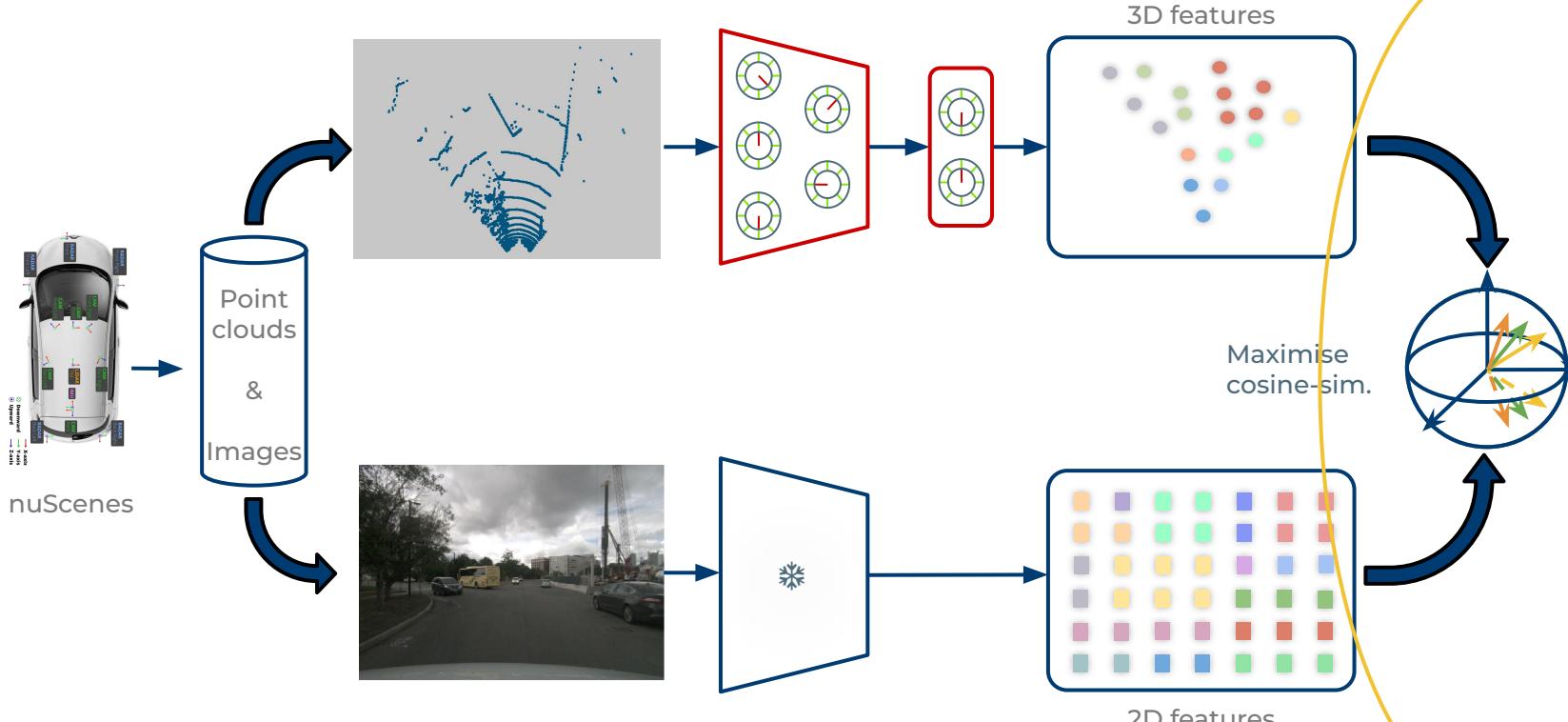


# Simplifications done in ScaLR



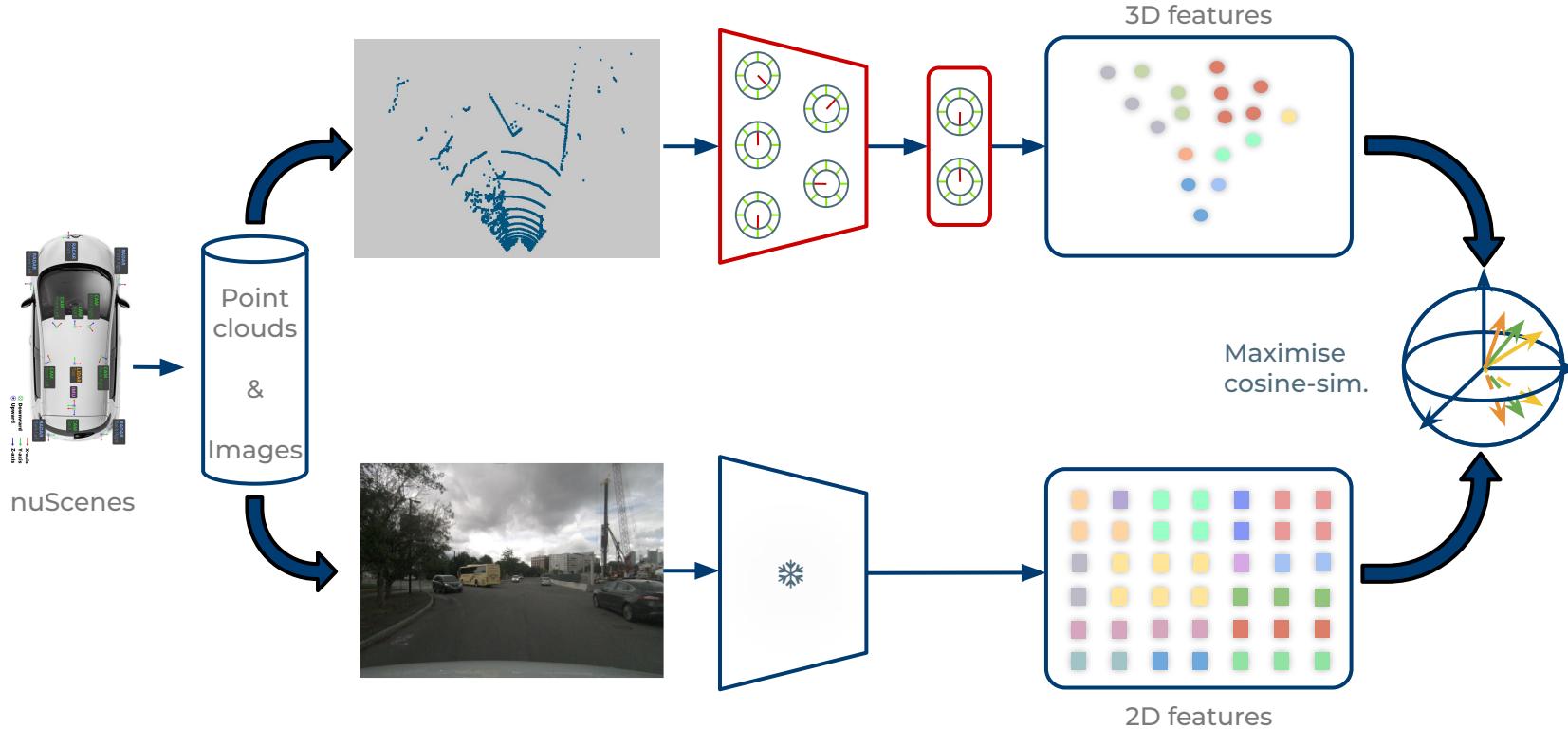


# Simplifications done in ScaLR





# Simplifications done in ScaLR





## Impact of these simplifications



# Impact of simplifications

	2D head	Contrast. loss	Cosine loss	No. cam.	Mem. / Time ↓	mIoU% ↑
MinkUNet	✓	✓	-	6	1.0 / 1.0	39.3
	✗	✓	-	6	1.6 / 1.0	43.4
	✗	-	✓	6	1.6 / 1.1	43.6
	✗	-	✓	1	0.8 / 0.3	42.1
WI-48-256	✓	✓	-	1	-	50.0
	✗	✓	-	1	-	54.1
	✗	-	✓	1	-	56.7



# Impact of simplifications

	2D head	Contrast. loss	Cosine loss	No. cam.	Mem. / Time ↓	mIoU% ↑
MinkUNet	✓	✓	-	6	1.0 / 1.0	39.3
	✗	✓	-	6	1.6 / 1.0	43.4
	✗	-	✓	6	1.6 / 1.1	43.6
	✗	-	✓	1	0.8 / 0.3	42.1
WI-48-256	✓	✓	-	1	-	50.0
	✗	✓	-	1	-	54.1
	✗	-	✓	1	-	56.7



# Impact of simplifications

	2D head	Contrast. loss	Cosine loss	No. cam.	Mem. / Time ↓	mIoU% ↑
MinkUNet	✓	✓	-	6	1.0 / 1.0	39.3
	✗	✓	-	6	1.6 / 1.0	43.4
	✗	-	✓	6	1.6 / 1.1	43.6
	✗	-	✓	1	0.8 / 0.3	42.1
WI-48-256	✓	✓	-	1	-	50.0
	✗	✓	-	1	-	54.1
	✗	-	✓	1	-	56.7



# Impact of simplifications

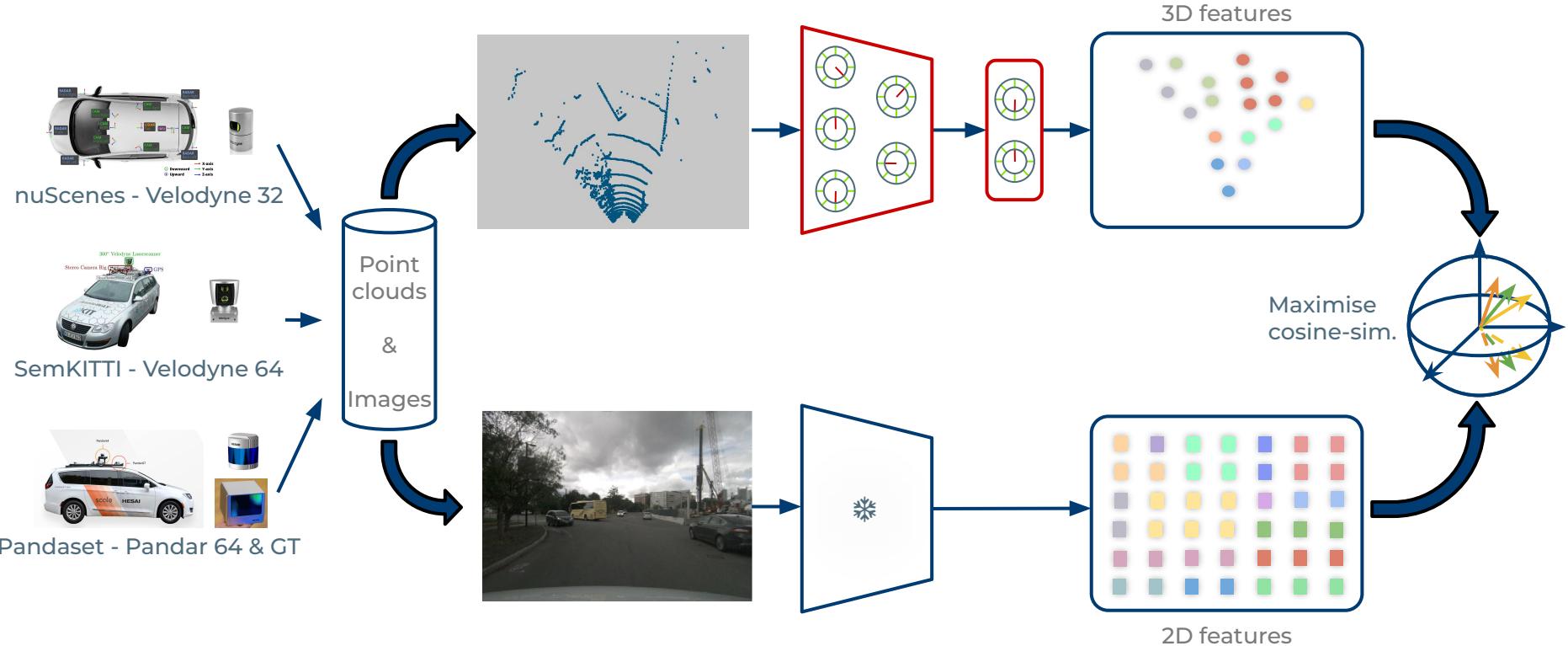
	2D head	Contrast. loss	Cosine loss	No. cam.	Mem. / Time ↓	mIoU% ↑
MinkUNet	✓	✓	-	6	1.0 / 1.0	39.3
	✗	✓	-	6	1.6 / 1.0	43.4
	✗	-	✓	6	1.6 / 1.1	43.6
	✗	-	✓	1	0.8 / 0.3	42.1
WI-48-256	✓	✓	-	1	-	50.0
	✗	✓	-	1	-	54.1
	✗	-	✓	1	-	56.7



## **Study influence of three pillars**



# ScaLR





# The 3 pillars

## Pillar 3: Mix of datasets



nuScenes - Velodyne 32



SemKITTI - Velodyne 64

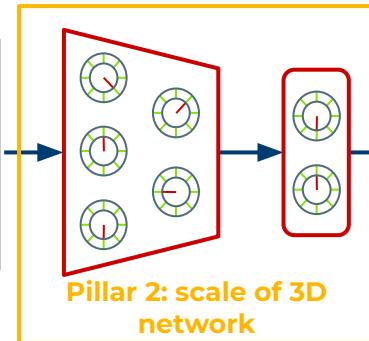


Pandaset - Pandar 64 & GT

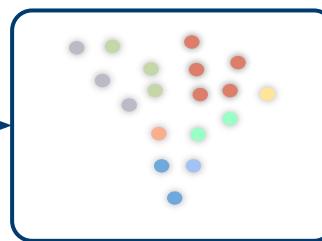
Point clouds  
&  
Images



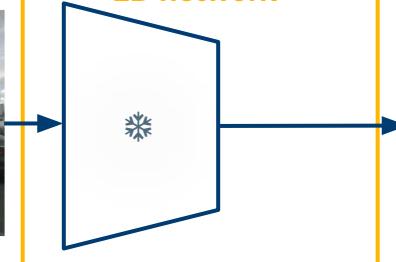
## Pillar 2: scale of 3D network



3D features



## Pillar 1: choice & scale of 2D network



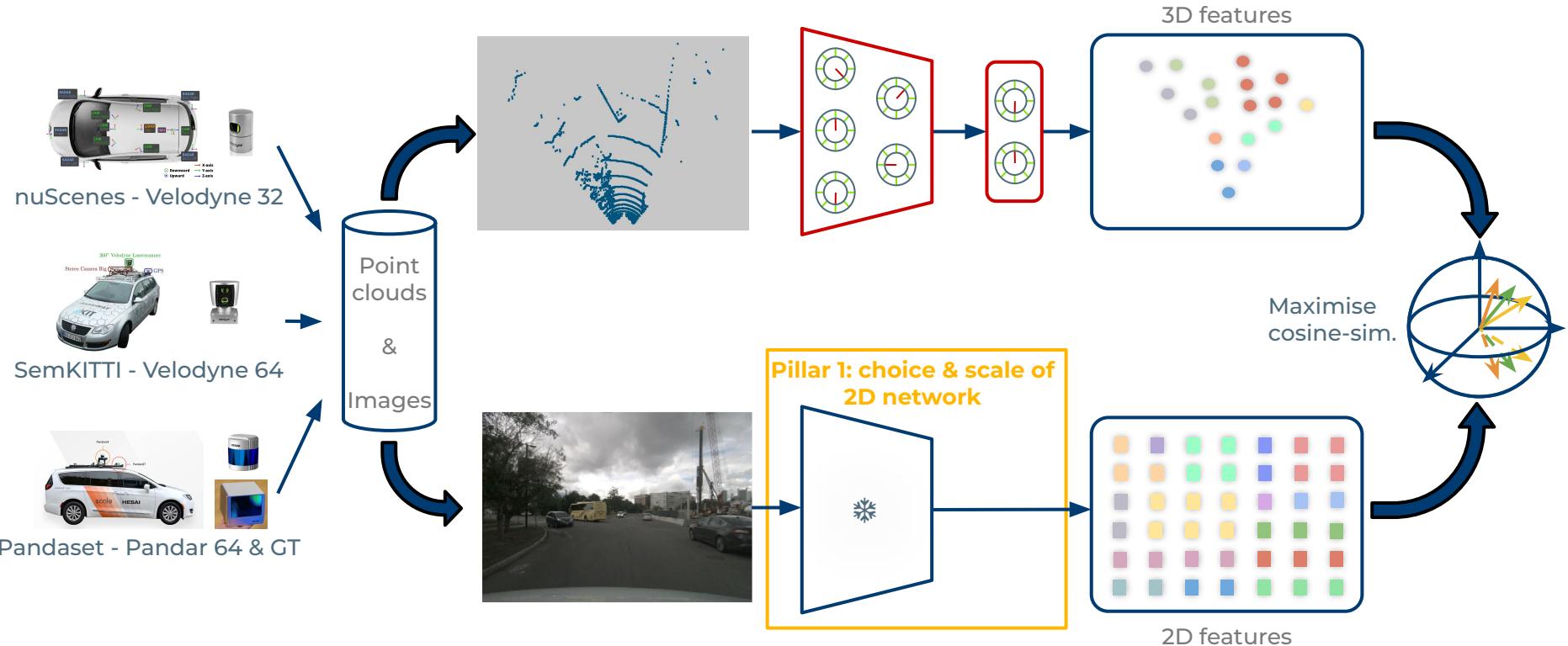
2D features

Maximise  
cosine-sim.



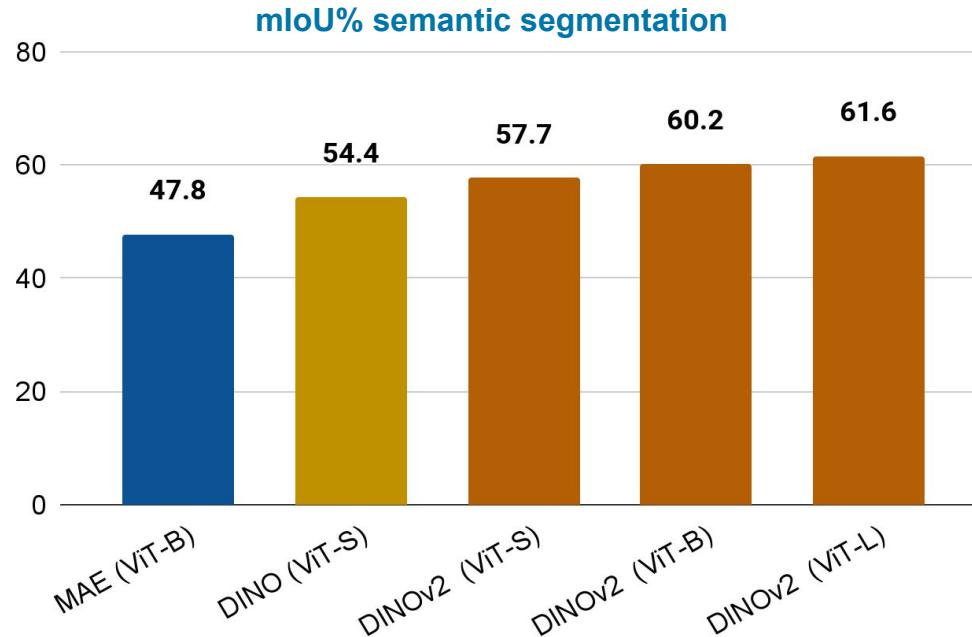


# Pillar 1: choice & scale of 2D network





## Pillar 1: choice & scale of 2D network

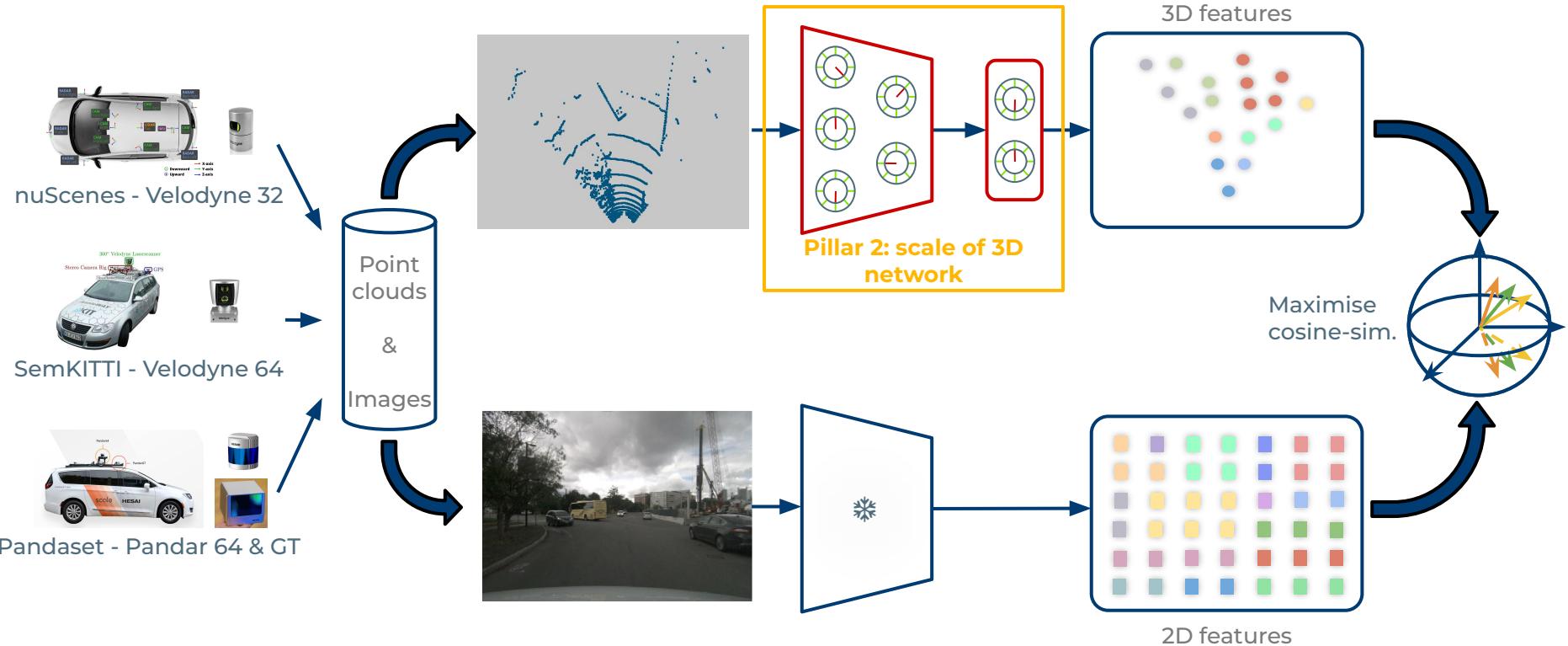


Pretraining & linear probing on nuScenes

- Distilled into WaffleIron-256 (features of size 256)

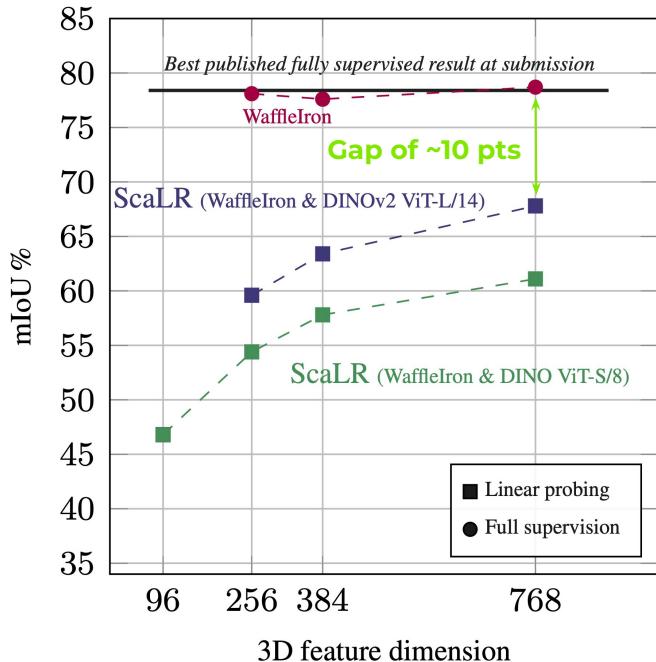


## Pillar 2: scale of 3D network





## Pillar 2: scale of 3D network



Pretraining & linear probing on nuScenes

- Distilled into WaffleIron with varying feature sizes



## Pillar 3: Mix of datasets

### Pillar 3: Mix of datasets



nuScenes - Velodyne 32

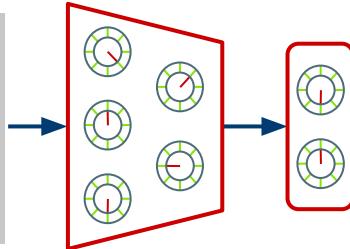


SemKITTI - Velodyne 64

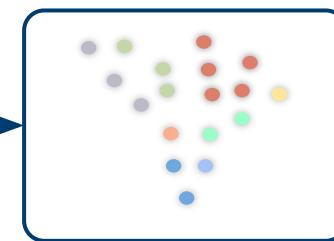


Pandaset - Pandar 64 & GT

Point clouds  
&  
Images



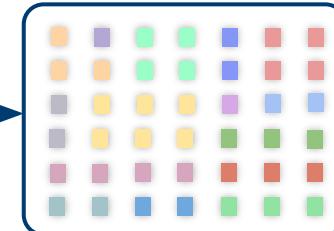
3D features



Maximise  
cosine-sim.



2D features





## Pillar 3: Mix of datasets

Pretrain. Dataset	Downstream & Test Dataset				
	nuScenes	SemKITTI	Pand. 64	Pand. GT	
<i>Pretraining with DINO-ViT-S/8 and linear probing</i>					
WI-256	nuScenes	<b>54.4</b>	28.8	26.9	25.2
	SemKITTI	39.5	<b>46.6</b>	25.3	25.7
	Pand. 64	39.6	25.6	<b>30.0</b>	24.7
	Pand. GT	29.9	26.9	23.5	<b>28.5</b>
	All datasets	<b>54.6</b>	<b>50.6</b>	<b>33.1</b>	<b>32.3</b>
<i>Pretraining with DINoV2-ViT-L/14 and linear probing</i>					
WI-768	nuScenes	<b>67.8</b>	43.1	33.9	29.9
	All datasets	<b>67.8</b>	<b>55.8</b>	<b>37.9</b>	<b>34.5</b>

Pretraining on individual datasets or mix of all - Linear probing on each datasets



## Pillar 3: Mix of datasets

Dataset	Pretrain.	Downstream & Test Dataset			
		nuScenes	SemKITTI	Pand. 64	Pand. GT
<i>Pretraining with DINO-ViT-S/8 and linear probing</i>					
WI-256	nuScenes	<b>54.4</b>	28.8	26.9	25.2
	SemKITTI	39.5	<b>46.6</b>	25.3	25.7
	Pand. 64	39.6	25.6	<b>30.0</b>	24.7
	Pand. GT	29.9	26.9	23.5	<b>28.5</b>
	All datasets	<b>54.6</b>	<b>50.6</b>	<b>33.1</b>	<b>32.3</b>
<i>Pretraining with DINoV2-ViT-L/14 and linear probing</i>					
WI-768	nuScenes	<b>67.8</b>	43.1	33.9	29.9
	All datasets	<b>67.8</b>	<b>55.8</b>	<b>37.9</b>	<b>34.5</b>

Pretraining on individual datasets or mix of all - Linear probing on each datasets



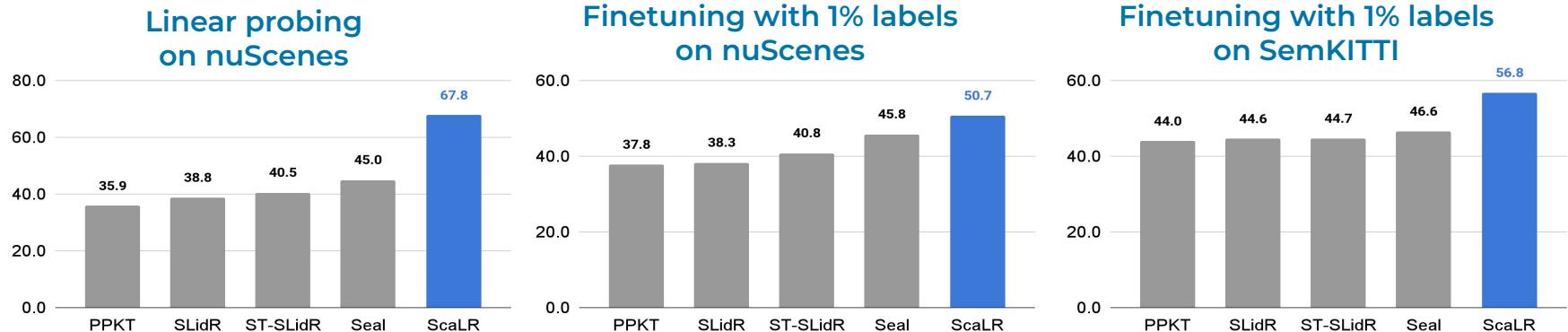
## Pillar 3: Mix of datasets

Pretrain. Dataset	Downstream & Test Dataset				
	nuScenes	SemKITTI	Pand. 64	Pand. GT	
<i>Pretraining with DINO-ViT-S/8 and linear probing</i>					
WI-256	nuScenes	<b>54.4</b>	28.8	26.9	25.2
	SemKITTI	39.5	<b>46.6</b>	25.3	25.7
	Pand. 64	39.6	25.6	<b>30.0</b>	24.7
	Pand. GT	29.9	26.9	23.5	28.5
	All datasets	<b>54.6</b>	<b>50.6</b>	<b>33.1</b>	<b>32.3</b>
<i>Pretraining with DINoV2-ViT-L/14 and linear probing</i>					
WI-768	nuScenes	<b>67.8</b>	43.1	33.9	29.9
	All datasets	<b>67.8</b>	<b>55.8</b>	<b>37.9</b>	<b>34.5</b>

Pretraining on individual datasets or mix of all - Linear probing on each datasets



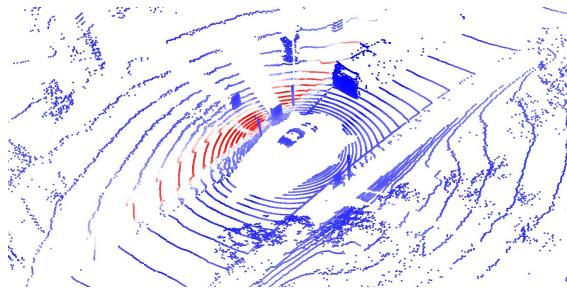
# Progress after optimizing the 3 pillars



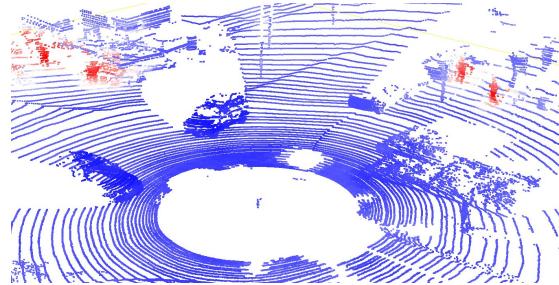


# Qualitative results

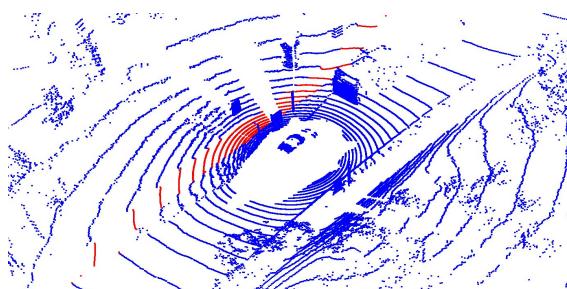
Correlation maps with class prototypes



*nuScenes - Sidewalk*



*Sem.KITTI - Pedestrian*

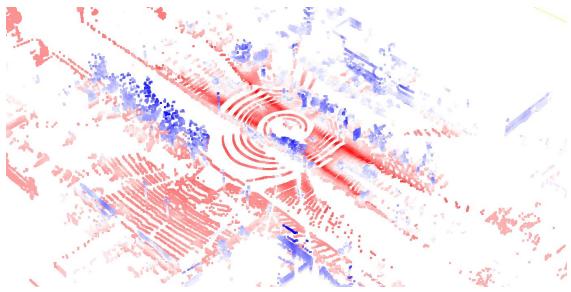


*Ground truth*

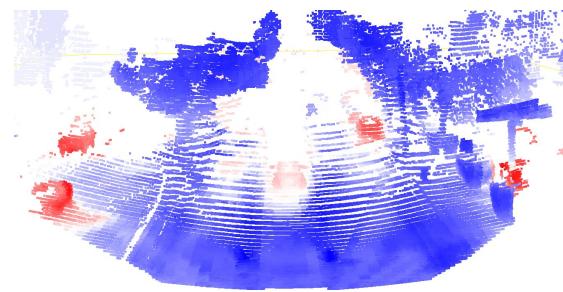


## Qualitative results

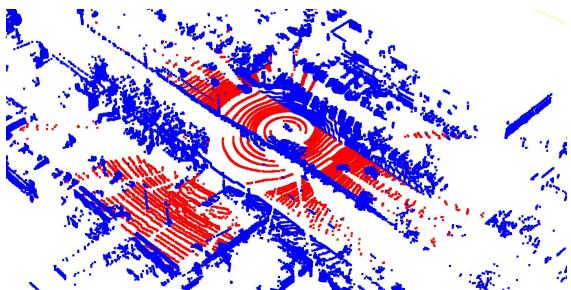
*Correlation maps with class prototypes*



*PandaSet 64 - Road*



*PandaSet GT - Car*

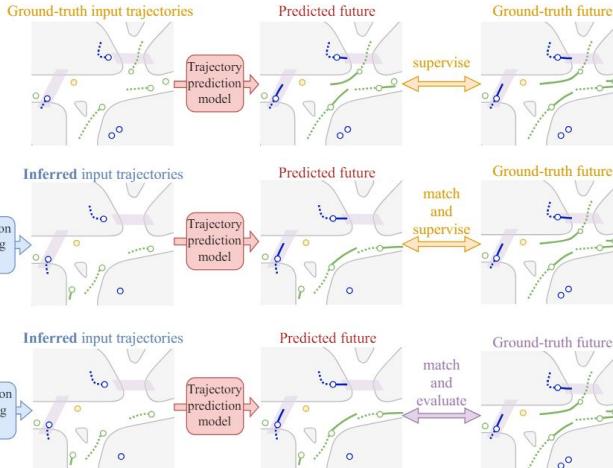


*Ground truth*

# **Conclusion and perspectives**

# Conclusion

- Scalable and simple pretraining strategies become first class citizens in AD
- Data outside the AD domain (internet, robotics) is key for camera perception
- Training over multiple datasets and sensor configurations shows promising results (ScaLR, UniTraj)
- Increased reliance on previously trained modules



Valeo4Cast approach at the Argoverse 2  
end-to-end forecasting challenge 2024

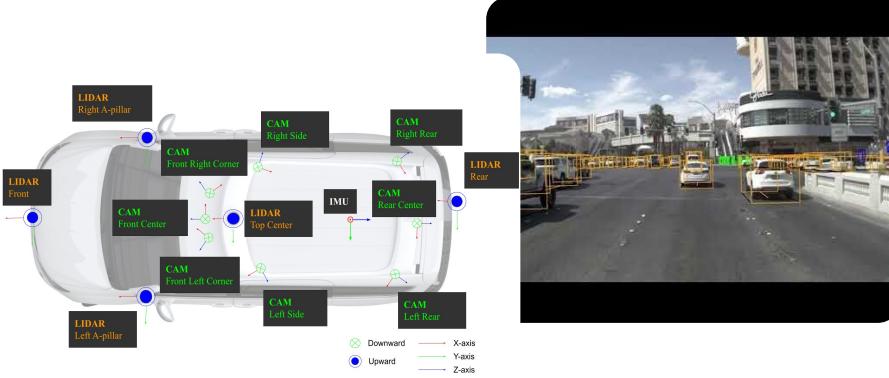
		Forecasting			Detection			Tracking
		mAP <sub>F</sub> (↑)	ADE (↓)	FDE (↓)	Inputs	Training range	mCDS (↑)	HOTA (↑)
2023	CenterPoint [28]	-	-	-	L	150	14	-
	BEVFusion [13]	-	-	-	L+C	150	37	-
	Anony_3D (v0) [11]	-	-	-	L	150	31	44.36
	Host_4626_Team [17]	14.51	5.10	7.32	-	-	-	39.98
	Dgist-cvlab [25]	42.91	4.11	4.59	L+C	150	34	41.49
	Le3DE2E [22]	46.70	3.22	3.76	L+C	150	39	56.19
2024	Dgist-cvlab	45.83	4.09	4.53	L+C	150	34	41.49
	Le3DE2E	50.53	4.07	4.60	L+C	150	<b>43</b>	<b>64.60</b>
	<b>Valeo4Cast</b>	<b>63.82</b>	<b>2.14</b>	<b>2.43</b>	L	75	31	61.28

Valeo4Cast results across sub-challenges

# Perspectives

- Different forms of world models taking over paving the path to embodied AI
- More downstream use-cases around foundation models
- Arrival of new larger scale datasets in the community
- Need for more open research in the field (data, code, models, collabs)

nuPlan  
(calibrated multi-cam + LiDAR)  
~ 100h of driving data @10Hz



OpenDV-Youtube  
(only non-calib. front-cam)  
~ 1700h of driving data @10Hz



**<https://valeoai.github.io>**