



TOYOTA  
RESEARCH INSTITUTE

# Scaling up ML for Autonomy

**Adrien Gaidon** (twitter: @adnothing)

**Head of Machine Learning Research**  
**Toyota Research Institute (TRI), CA, USA**

## Collaborators

### ML-Research team

(V. Guizilini, Jie Li,  
R. Ambrus, W. Kehl, et al)

### ML-Engineering team

(S. Pillai, A. Raventos, A.  
Bhargava, KH. Lee, et al)

### **Wolfram Burgard**

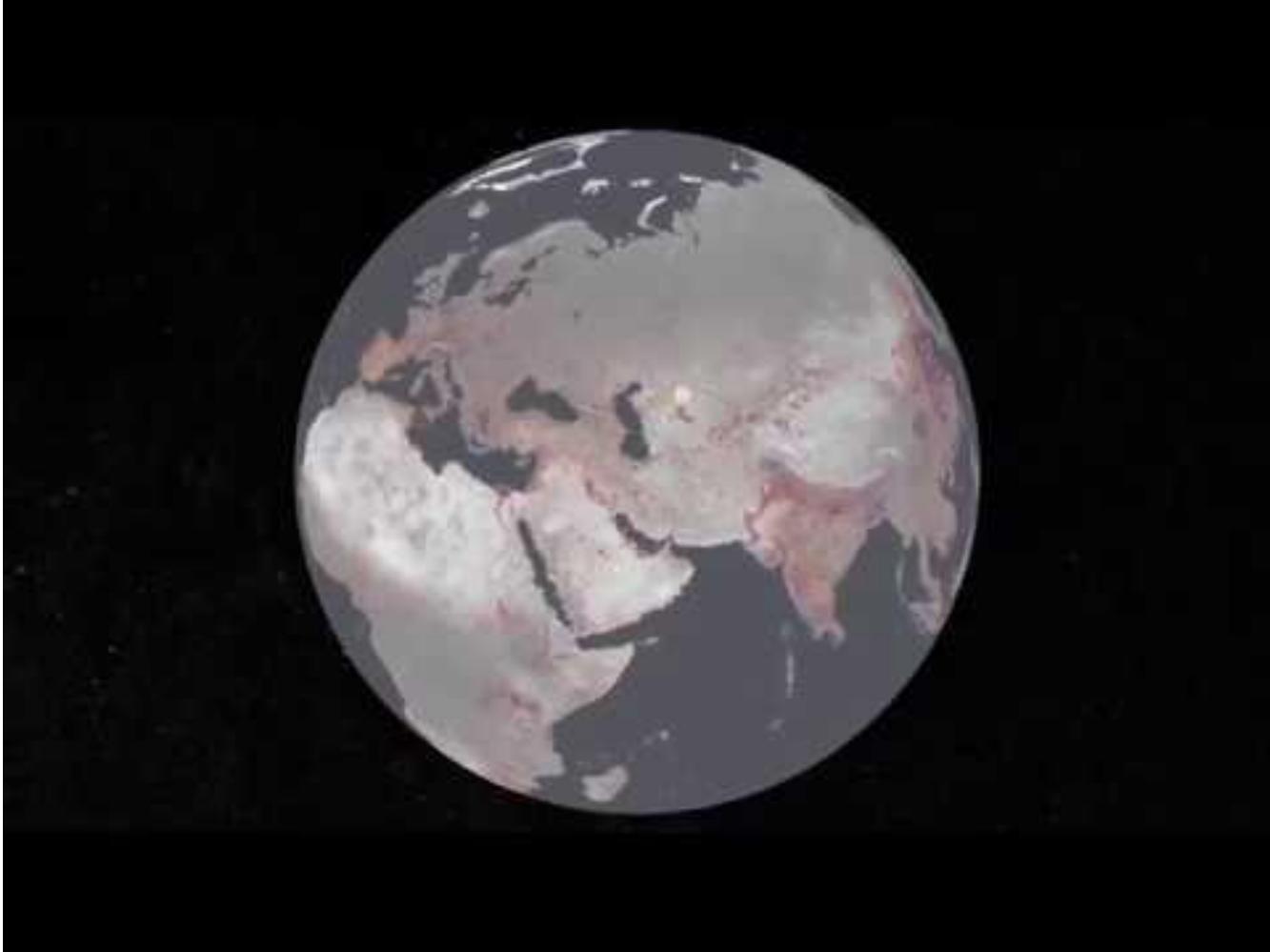
**Stanford** (SVL, MRL, ASL)



# 1.35 MILLION

---

ROAD TRAFFIC DEATHS PER YEAR



Planet ©

[youtu.be/ZqjKDpbtVn0](https://youtu.be/ZqjKDpbtVn0)



**TOYOTA**

100M Cars, 95% Parked

---

**~10s PB**

---

AMOUNT OF DATA PER DAY



# World-scale Autonomy?

THE SWE

PROGRAM

EVERYTHING

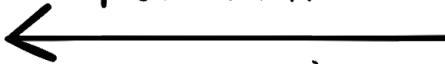
MAPS ?



# World-scale Autonomy?

THE SWE

PROGRAM



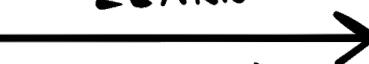
EVERYTHING

MAPS ?



THE SCIENTIST

LEARN



EVERYTHING

sim ?



# World-scale Autonomy?

THE SWE

THE MLE

THE SCIENTIST

PROGRAM

LABEL

LEARN

EVERYTHING

EVERYTHING

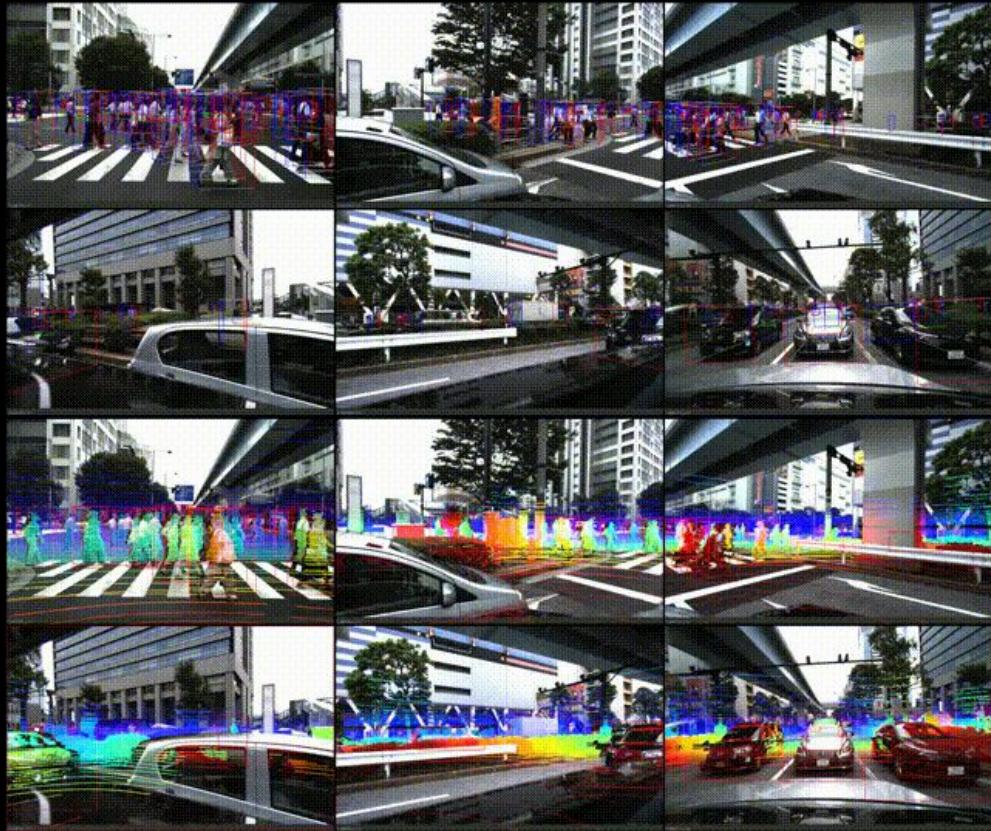
EVERYTHING

MAPS ?

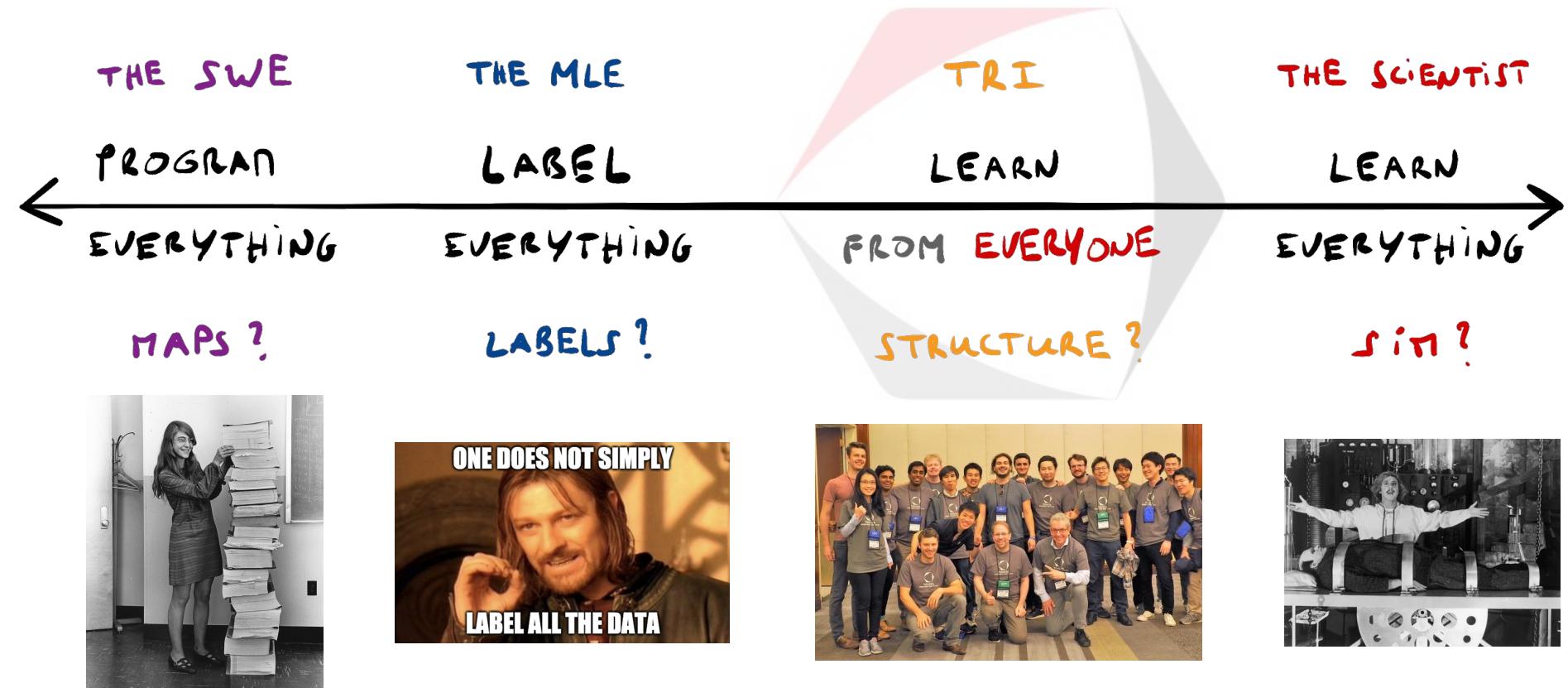
LABELS ?

sim ?





# World-scale Autonomy?



## Behavior: leverage large scale Demonstrations

**Exploring the Limitations of Behavior Cloning for Autonomous Driving, ICCV'19 (oral)**

Spatiotemporal Relationship Reasoning for Pedestrian Intent Prediction, RA-L & ICRA'20

It Is Not the Journey but the Destination: Endpoint Conditioned Trajectory Prediction, arXiv:2004.02025

Reinforcement Learning based Control of Imitative Policies for Near-Accident Driving, coming soon

Risk-Sensitive Sequential Action Control with Multi-Modal Human Trajectory Forecasting [...], coming soon

Driving Through Ghosts: Behavioral Cloning with False Positives, coming soon

## Supervised Learning: efficiently use available Labels

**ROI-10D: Monocular Lifting of 2D Detection to 6D Pose and Metric Shape, CVPR'19**

Learning Imbalanced Datasets with Label-Distribution-Aware Margin Loss, NeurIPS'19

Learning to Fuse Things and Stuff, arXiv:1812.01192

Spatio-Temporal Graph for Video Captioning with Knowledge Distillation, CVPR'20

**Real-Time Panoptic Segmentation from Dense Detections, CVPR'20 (oral)**

Hierarchical Lovász Embeddings for Proposal-free Panoptic Segmentation, coming soon

Unsupervised Estimation of Segmentation Difficulty, coming soon

## Geometry: Self / Semi-Supervised Pseudo-LiDAR and SfM

**SuperDepth: Self-Supervised, Super-Resolved Monocular Depth Estimation, ICRA'19**

Robust Semi-Supervised Monocular Depth Estimation with Reprojected Distances, CoRL'19

Two Stream Networks for Self-Supervised Ego-Motion Estimation, CoRL'19

Semantically-Guided Representation Learning for Self-Supervised Monocular Depth, ICLR'20

Neural Outlier Rejection for Self-Supervised Keypoint Learning, ICLR'20

Self-Supervised 3D Keypoint Learning for Ego-motion Estimation, arxiv:1912.03426

**3D Packing for Self-Supervised Monocular Depth Estimation, CVPR'20 (oral)**

Self-Supervised Neural Camera Models, coming soon

## Simulation: Domain Adaptation, Differentiable Rendering, RL

**SPIGAN: Privileged Adversarial Learning from Simulation, ICLR'19**

DeceptionNet: Network-Driven Domain Randomization, ICCV'19

Generating Human Action Videos by Coupling 3D Game Engines and Probabilistic Graphical Models, IJCV'20

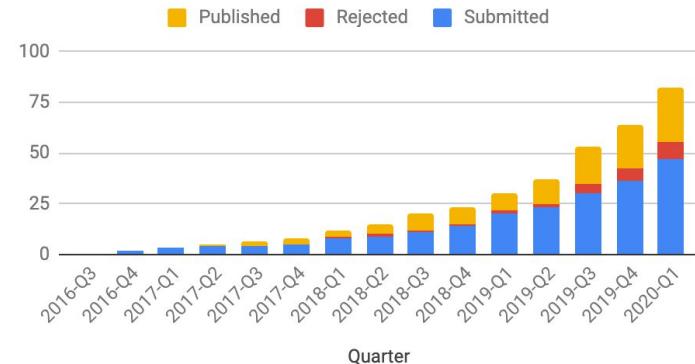
**Autolabeling 3D Objects with Differentiable Rendering of SDF Shape Priors, CVPR'20 (oral)**

Self-Supervised Differentiable Rendering for Monocular 3D Object Detection, coming soon

Behaviorally Diverse Traffic Simulation via Reinforcement Learning, coming soon

Discovering Avoidable Planner Failures [...] in Behaviorally Diverse Simulation, coming soon

## ML Publications History (cumulative)



## Upcoming workshops co-organized by TRI

### ICML: AI for Autonomous Driving (AIAD)

<https://sites.google.com/view/aiad2020>

### ECCV: Perception for Autonomous Driving (PAD)

<https://sites.google.com/view/pad2020>

## Upcoming TRI Dataset Releases

### STIP: Stanford-TRI Intent Prediction

<http://stip.stanford.edu/>

### DDAD: Dense Depth for Autonomous Driving

<https://github.com/TRI-ML/DDAD>

# Scaling up ML for Autonomy

## Behavior Cloning and its Limitations

*Exploring the Limitations of Behavior Cloning  
for Autonomous Driving, Codevilla et al, ICCV'19 (oral)*

Real-time Panoptic Segmentation

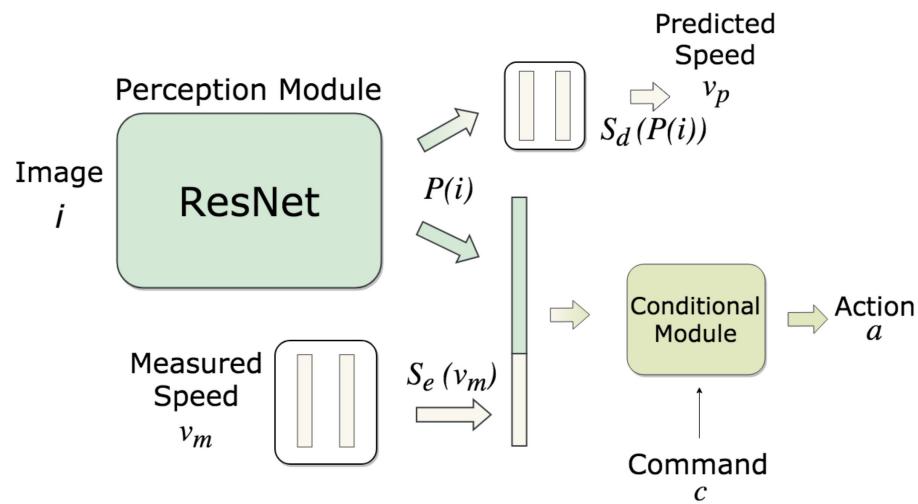
Self-Supervised Pseudo-Lidar Networks

Auto-labeling via Differentiable Rendering

# Trillions of Miles driven yearly (3.2T in US alone)

## Behavior Cloning: simplest Imitation Learning

## Modernized: deeper ResNet, ImageNet, Data++



| module           | input dimension        | channels |
|------------------|------------------------|----------|
| Perception       | ResNet 34 [15] outputs | 512      |
|                  | 1                      | 128      |
| Measured Speed   | 128                    | 128      |
|                  | 128                    | 128      |
| Speed Prediction | 512                    | 256      |
|                  | 256                    | 256      |
|                  | 256                    | 1        |
| Joint input      | $512 + 128$            | 512      |
|                  | 512                    | 256      |
| Control          | 256                    | 256      |
|                  | 256                    | 1        |

| Task         | Training conditions |           |            |        |      |           | New town & weather |           |         |        |           |           |
|--------------|---------------------|-----------|------------|--------|------|-----------|--------------------|-----------|---------|--------|-----------|-----------|
|              | CIL[10]             | CIRL[26]  | CAL[36]    | MT[25] | CILR | CILRS     | CIL[10]            | CIRL[26]  | CAL[36] | MT[25] | CILR      | CILRS     |
| Straight     | 98                  | 98        | <b>100</b> | 96     | 94   | 96        | 80                 | <b>98</b> | 94      | 96     | 92        | 96        |
| One Turn     | 89                  | <b>97</b> | <b>97</b>  | 87     | 92   | 92        | 48                 | 80        | 72      | 82     | <b>92</b> | <b>92</b> |
| Navigation   | 86                  | 93        | 92         | 81     | 88   | <b>95</b> | 44                 | 68        | 68      | 78     | 88        | <b>92</b> |
| Nav. Dynamic | 83                  | 82        | 83         | 81     | 85   | <b>92</b> | 42                 | 62        | 64      | 62     | 82        | <b>90</b> |

Table 1. Comparison with the state of the art on the original CARLA benchmark. The “CILRS” version corresponds to our CIL-based ResNet using the speed prediction branch, whereas “CILR” is without this speed prediction. These two models and CIL are the only ones that do not use any extra supervision or online interaction with the environment during training. The table reports the percentage of successfully completed episodes in each condition, selecting the best seed out of five runs.

| Task    | Training conditions |                              |            |            |                              |            | New Town & Weather |            |            |                              |  |
|---------|---------------------|------------------------------|------------|------------|------------------------------|------------|--------------------|------------|------------|------------------------------|--|
|         | CIL[10]             | CAL[36]                      | MT[25]     | CILR       | CILRS                        | CIL[10]    | CAL[36]            | MT[25]     | CILR       | CILRS                        |  |
| Empty   | $79 \pm 1$          | $81 \pm 1$                   | $84 \pm 1$ | $92 \pm 1$ | <b><math>97 \pm 2</math></b> | $24 \pm 1$ | $25 \pm 3$         | $57 \pm 0$ | $66 \pm 2$ | <b><math>90 \pm 2</math></b> |  |
| Regular | $60 \pm 1$          | $73 \pm 2$                   | $54 \pm 2$ | $72 \pm 5$ | <b><math>83 \pm 0</math></b> | $13 \pm 2$ | $14 \pm 2$         | $32 \pm 2$ | $54 \pm 2$ | <b><math>56 \pm 2</math></b> |  |
| Dense   | $21 \pm 2$          | <b><math>42 \pm 3</math></b> | $13 \pm 4$ | $28 \pm 1$ | <b><math>42 \pm 2</math></b> | $2 \pm 0$  | $10 \pm 0$         | $14 \pm 2$ | $13 \pm 4$ | <b><math>24 \pm 8</math></b> |  |

Table 2. Results on our *NoCrash* benchmark. Mean and standard deviation on three runs, as CARLA 0.8.4 has significant non-determinism.

# Motivation

Trillion miles driven yearly - Requires expensive data annotation



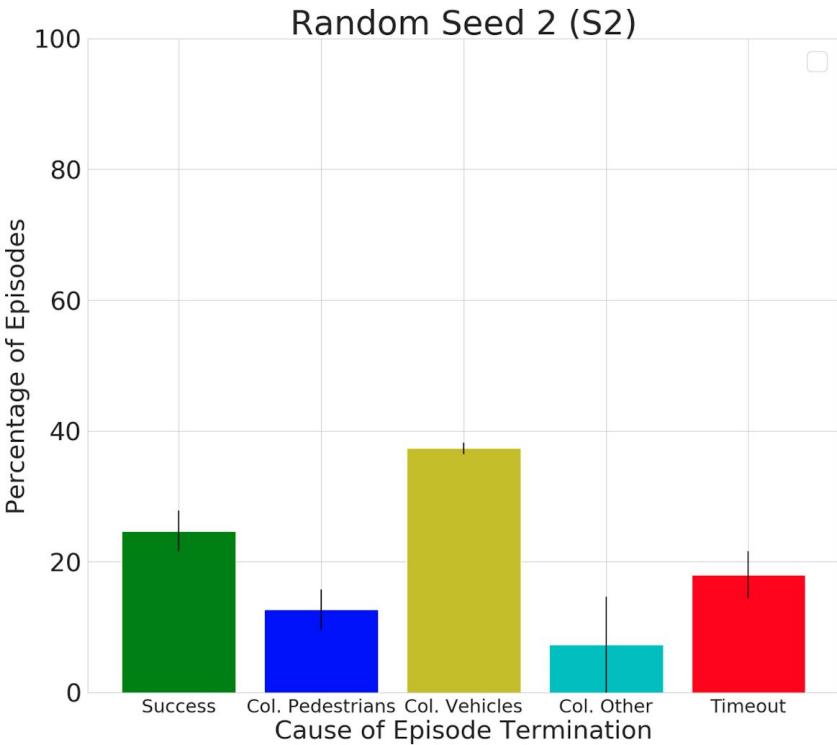
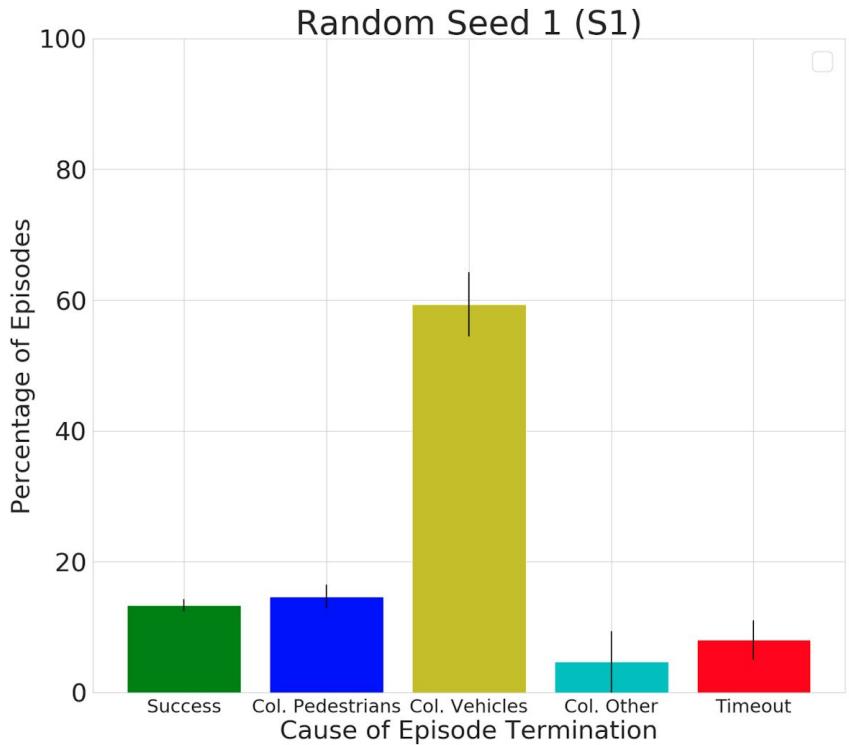


Figure 6. Cause of episode termination on *NoCrash* for two CILRS models (trained on 10 hours with ImageNet initialization) with identical parameters but different random seeds. The episodes were ran under “New Weather & Town” conditions of the “Dense Traffic” task.

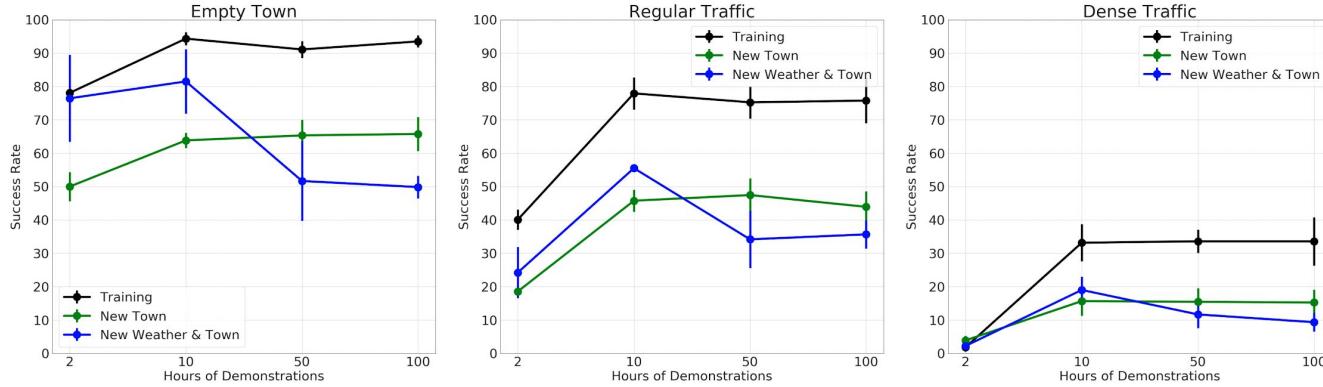


Figure 3. Due to biases in the data, the results may get either saturated or worse with increasing amounts of training data.

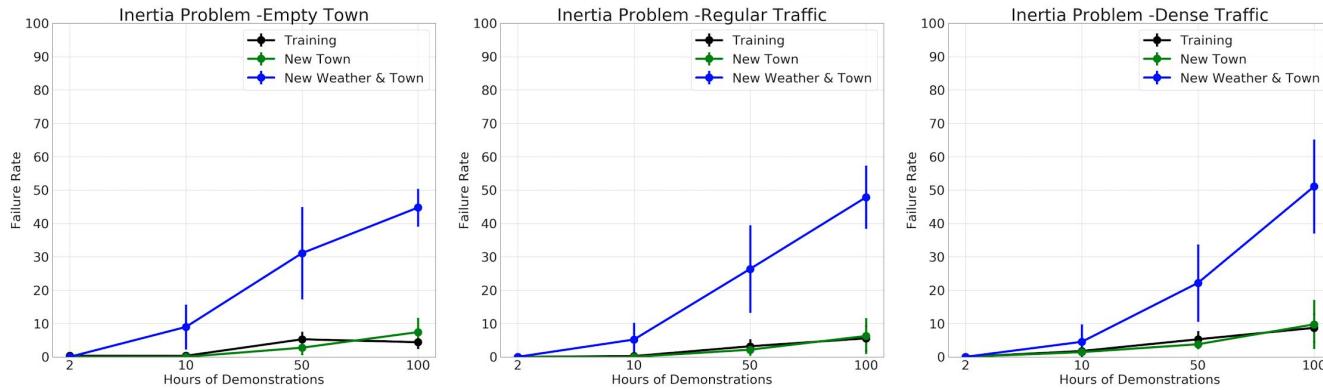


Figure 4. The percentage of episodes that failed due to the inertia problem. We can see that by increasing the amount of data, this bias may further degrade the generalization capabilities of the models.

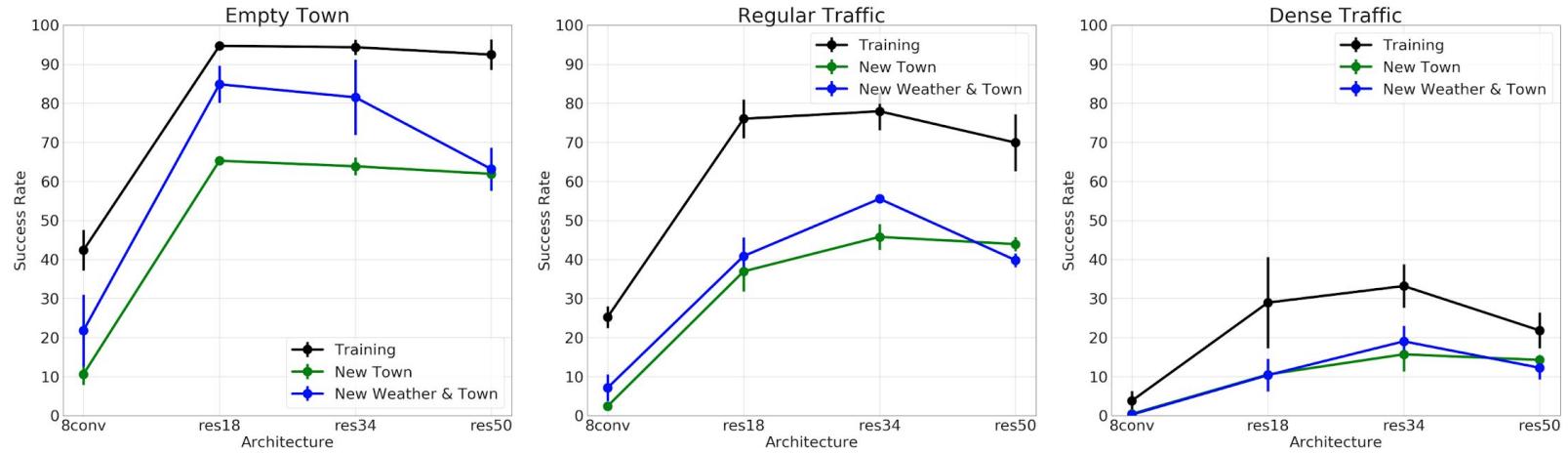
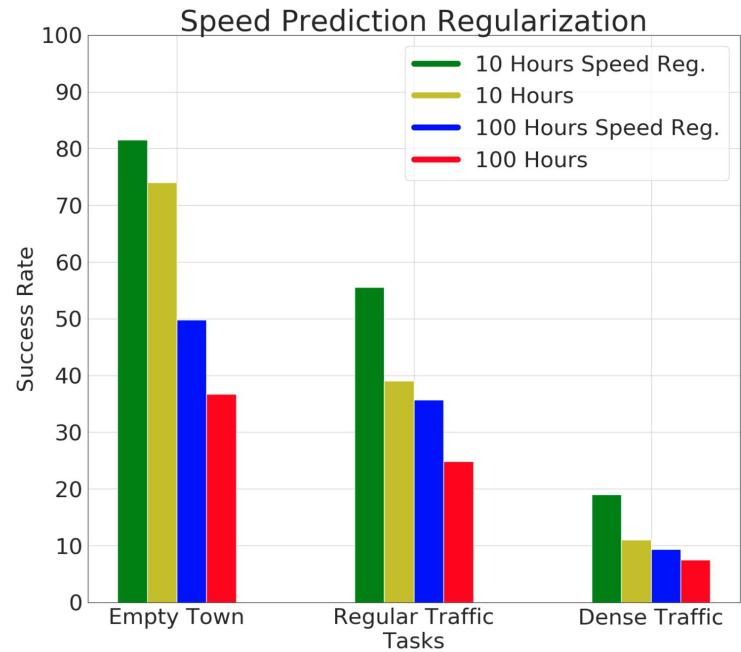


Figure 12. Ablative analysis between different architectures. The eight convolutions architecture, “8conv”, proposed by Codevilla [10] obtained poor results on the more complex CARLA100 benchmark. ResNet based deeper architectures, “res18” and “res34”, were able to improve the results. However, when testing ResNet 50 we notice a significant drop in the quality of the results.

|                  | Task    | Variance |
|------------------|---------|----------|
| CILRS            | Empty   | 23%      |
|                  | Regular | 26%      |
|                  | Dense   | 42%      |
| CILRS (ImageNet) | Empty   | 4%       |
|                  | Regular | 12%      |
|                  | Dense   | 38%      |

Table 3. Estimated variance of the success rate of CILRS on *NoCrash* computed by training 12 times the same model with different random seeds. The variance is reduced by fixing part of the initial weights with ImageNet pre-training.



# Scaling up ML for Autonomy

## Behavior Cloning and its Limitations

*Exploring the Limitations of Behavior Cloning  
for Autonomous Driving, Codevilla et al, ICCV'19 (oral)*

Real-time Panoptic Segmentation

Self-Supervised Pseudo-Lidar Networks

Auto-labeling via Differentiable Rendering

# Scaling up ML for Autonomy

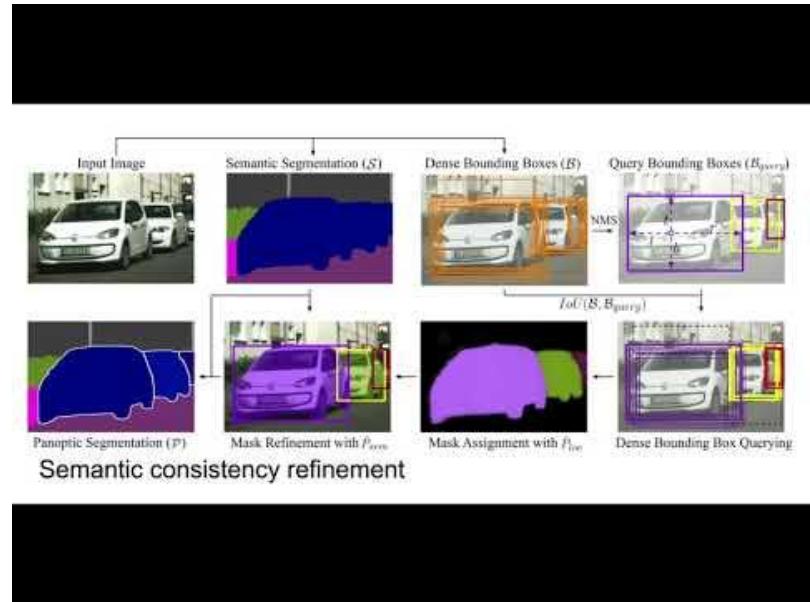
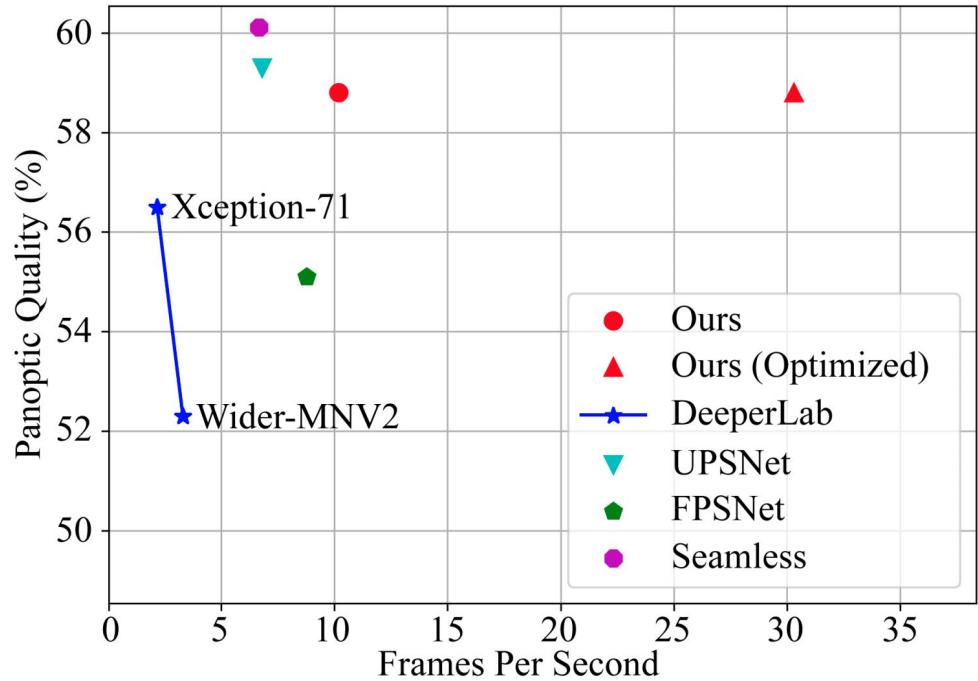
## Behavior Cloning and its Limitations

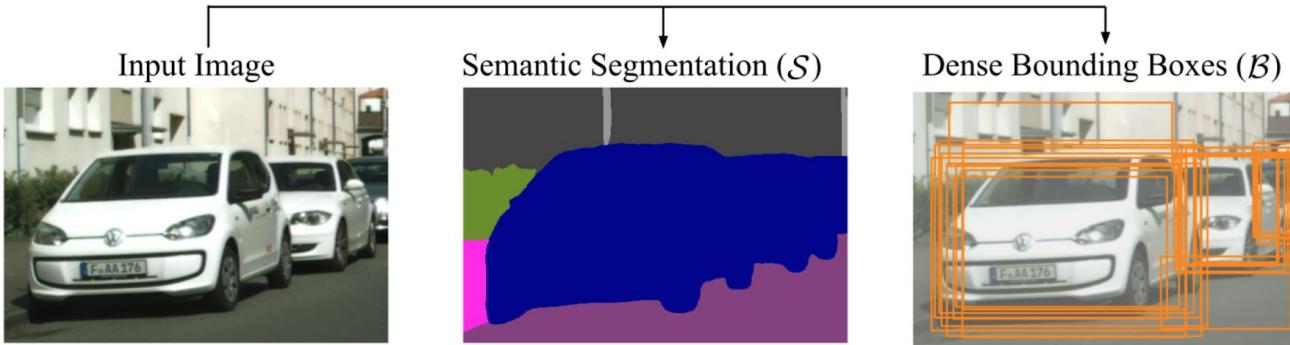
## Real-time Panoptic Segmentation

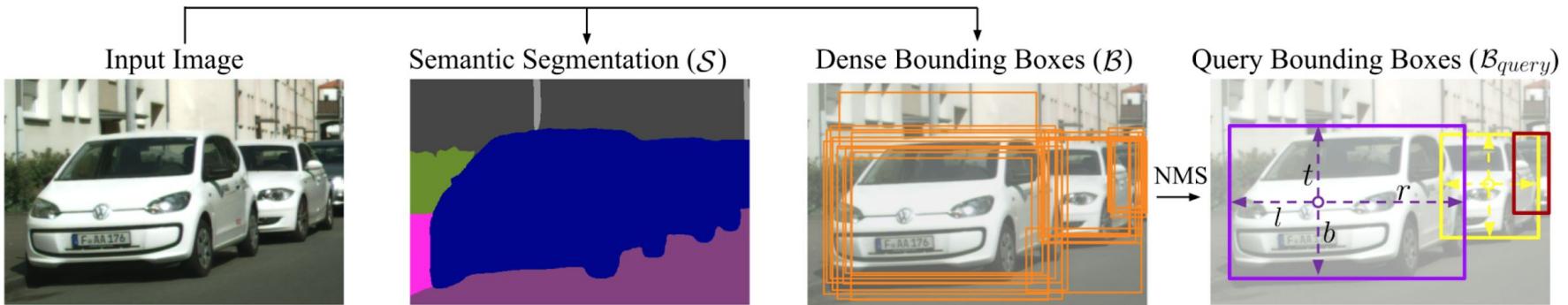
*Real-Time Panoptic Segmentation from Dense Detections,  
R. Hou\*, J. Li\*, A. Bhargava, A. Raventos et al, CVPR'20 (oral)*

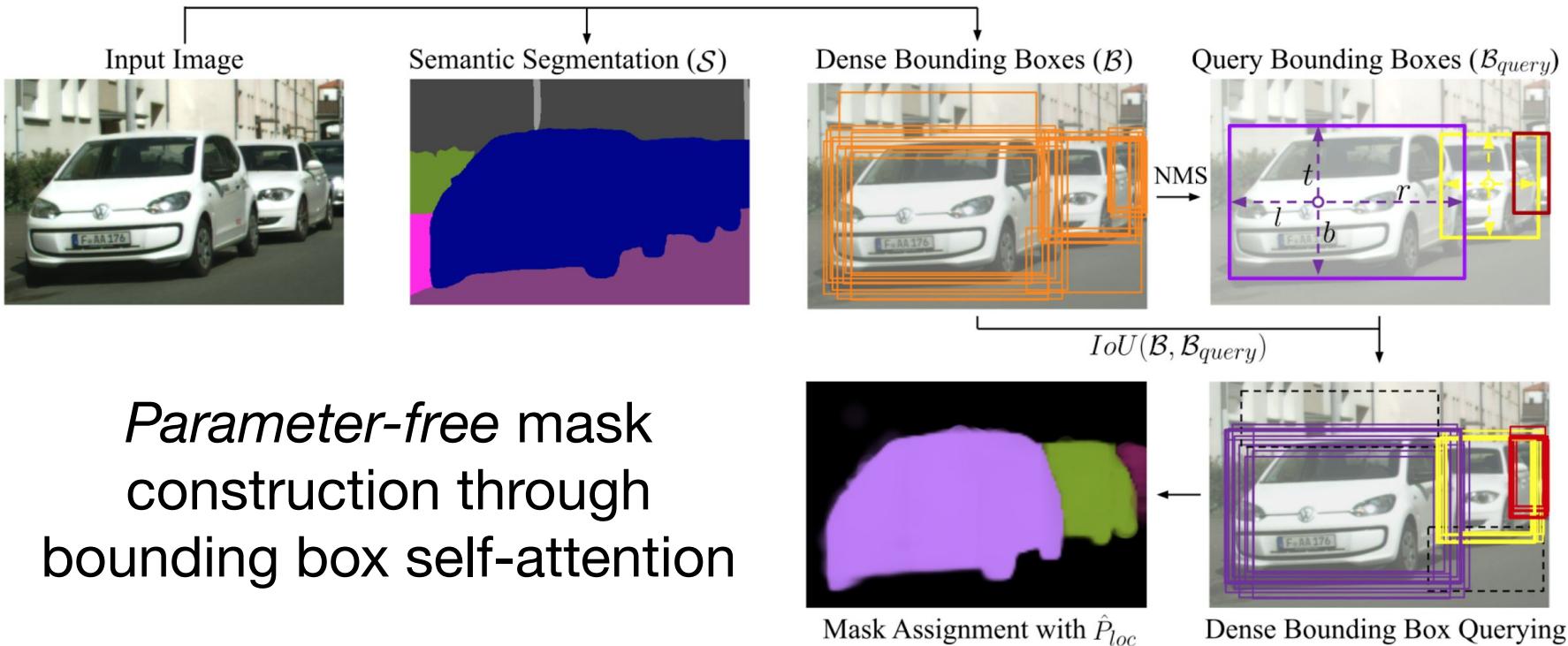
## Self-Supervised Pseudo-Lidar Networks

## Auto-labeling via Differentiable Rendering

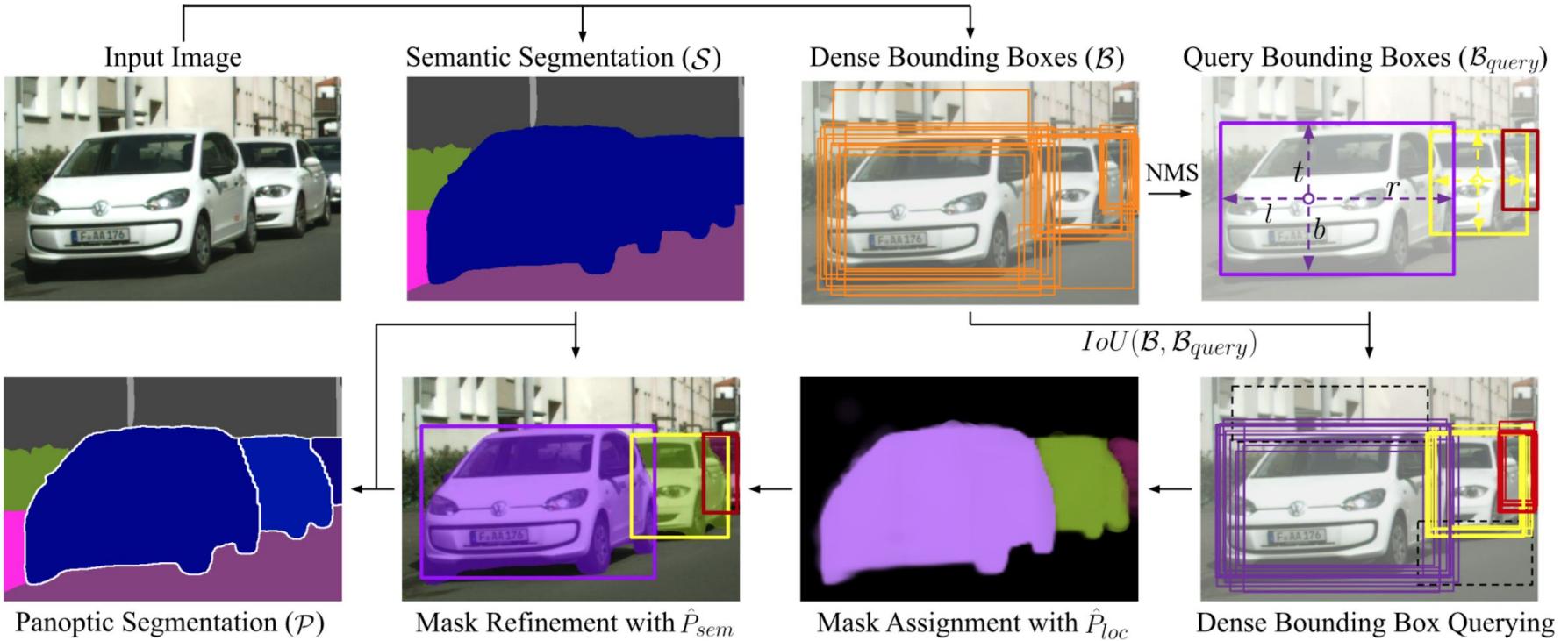


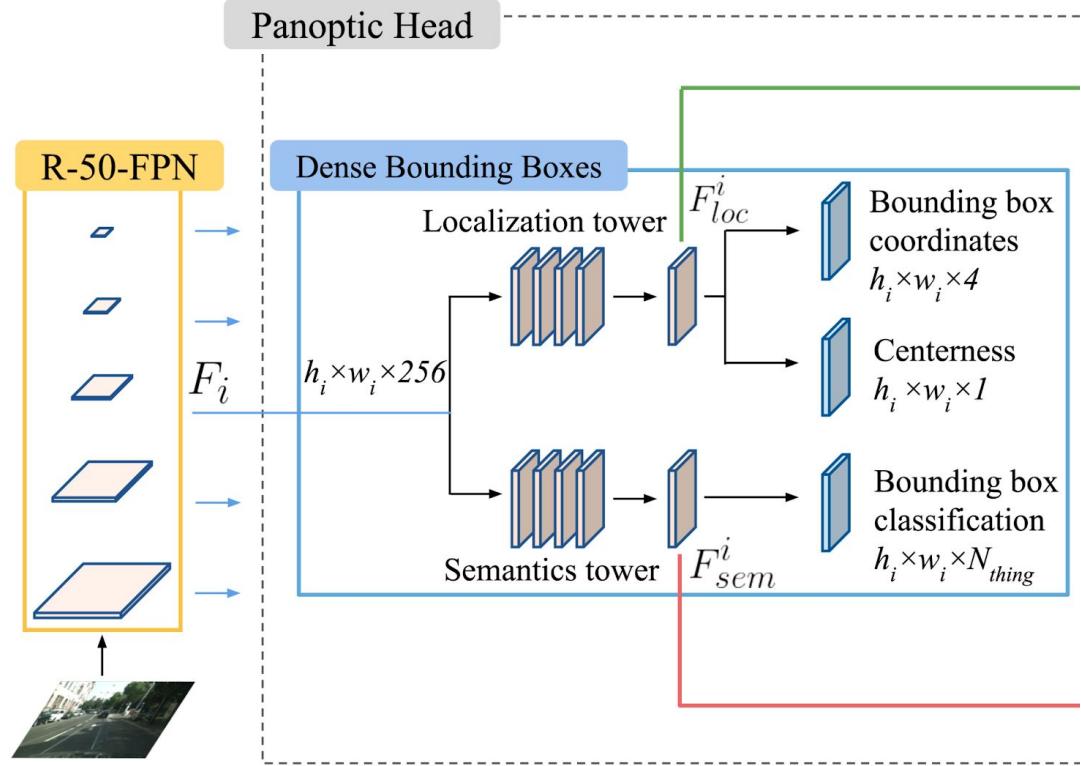


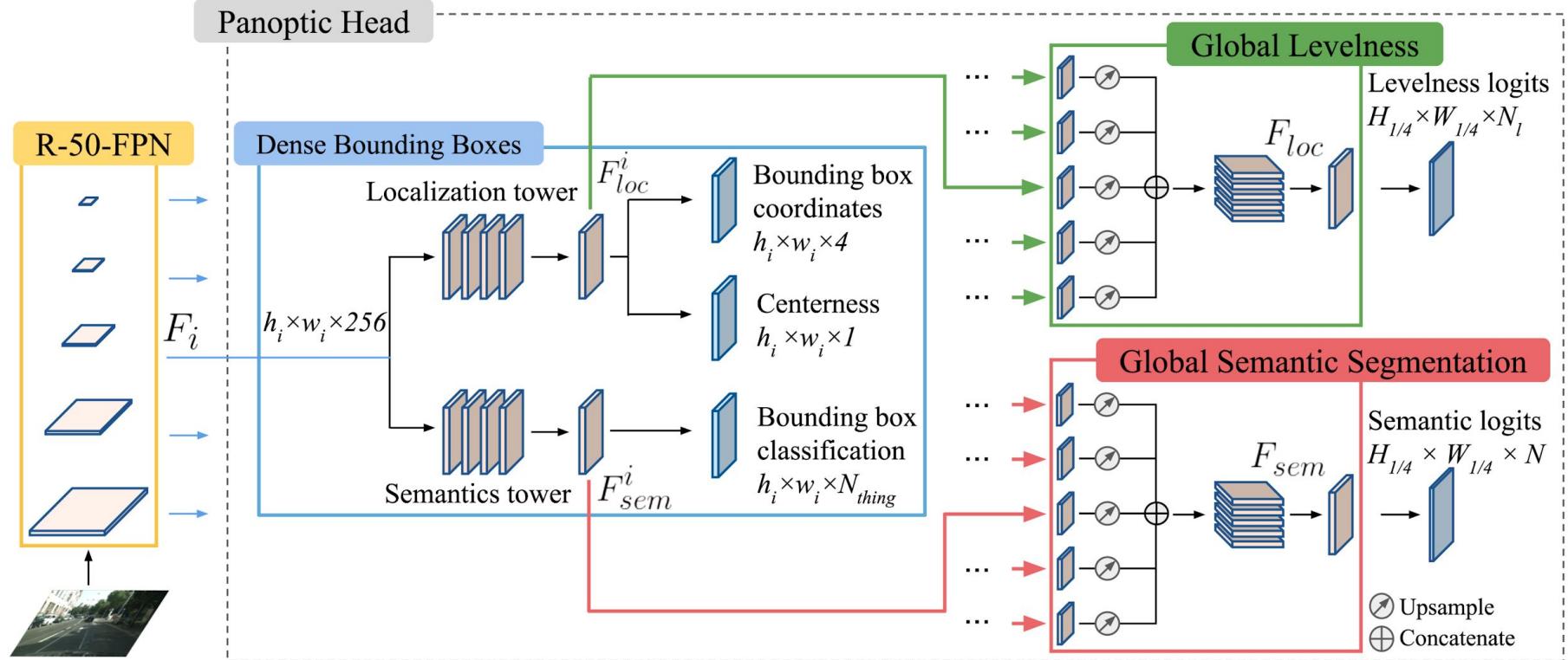




*Parameter-free mask construction through  
bounding box self-attention*







## Cityscapes (val)

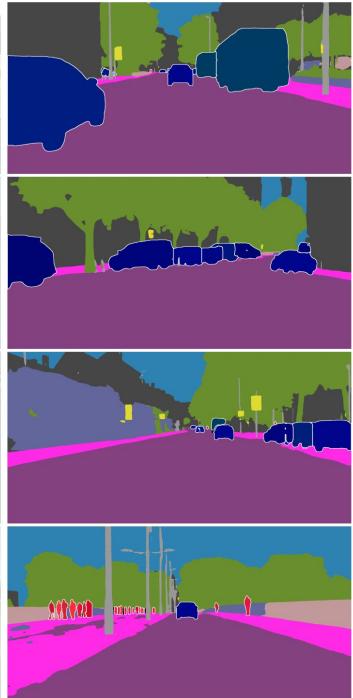
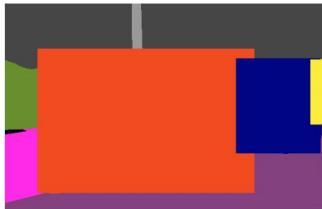
| Method                    | Backbone      | PQ          | $PQ^{th}$   | $PQ^{st}$   | mIoU        | AP          | GPU       | Inference Time |
|---------------------------|---------------|-------------|-------------|-------------|-------------|-------------|-----------|----------------|
| <b>Two-Stage</b>          |               |             |             |             |             |             |           |                |
| TASCNet [15]              | ResNet-50-FPN | 55.9        | 50.5        | 59.8        | -           | -           | V100      | 160ms          |
| AUNet[16]                 | ResNet-50-FPN | 56.4        | 52.7        | 59.0        | 73.6        | <b>33.6</b> | -         | -              |
| Panoptic-FPN [13]         | ResNet-50-FPN | 57.7        | 51.6        | 62.2        | 75.0        | 32.0        | -         | -              |
| AdaptIS <sup>†</sup> [30] | ResNet-50     | 59.0        | <b>55.8</b> | 61.3        | 75.3        | 32.3        | -         | -              |
| UPSNet [36]               | ResNet-50-FPN | 59.3        | 54.6        | 62.7        | 75.2        | 33.3        | V100      | 140ms*         |
| Seamless Panoptic [28]    | ResNet-50-FPN | 60.2        | 55.6        | <b>63.6</b> | 74.9        | 33.3        | V100      | 150ms*         |
| <b>Single-Stage</b>       |               |             |             |             |             |             |           |                |
| Deeplab [38]              | Wider MNV2    | 52.3        | -           | -           | -           | -           | V100      | 251ms          |
| FPSNet [7]                | ResNet-50-FPN | 55.1        | 48.3        | 60.1        | -           | -           | TITAN RTX | 114ms          |
| SSAP [8]                  | ResNet-50     | 56.6        | 49.2        | -           | -           | <b>31.5</b> | 1080Ti    | >260ms         |
| Deeplab [38]              | Xception-71   | 56.5        | -           | -           | -           | -           | V100      | 312ms          |
| Ours                      | ResNet-50-FPN | <b>58.8</b> | <b>52.1</b> | <b>63.7</b> | <b>77.0</b> | 29.8        | V100      | <b>99ms</b>    |



## COCO (val)

| Method                    | Backbone      | PQ          | $PQ^{th}$   | $PQ^{st}$   | Inf. Time   |
|---------------------------|---------------|-------------|-------------|-------------|-------------|
| <b>Two-Stage</b>          |               |             |             |             |             |
| Panoptic-FPN [13]         | ResNet-50-FPN | 33.3        | 45.9        | 28.7        | -           |
| AdaptIS <sup>†</sup> [30] | ResNet-50     | 35.9        | 40.3        | 29.3        | -           |
| AUNet [16]                | ResNet-50-FPN | 39.6        | <b>49.1</b> | 25.2        | -           |
| UPSNet [36]               | ResNet-50-FPN | 42.5        | 48.5        | <b>33.4</b> | 110ms*      |
| <b>Single-Stage</b>       |               |             |             |             |             |
| Deeplab [38]              | Xcep-71       | 33.8        | -           | -           | 94ms        |
| SSAP [8]                  | ResNet-50     | 36.5        | -           | -           | -           |
| Ours                      | ResNet-50-FPN | <b>37.1</b> | <b>41.0</b> | <b>31.3</b> | <b>63ms</b> |

# Supervision: Weak = 95% Strong



| Two towers                               | Levelness | Mask loss | PQ          | PQ <sup>th</sup> | PQ <sup>st</sup> |
|--|-----------|-----------|-------------|------------------|------------------|
| <b>Fully Supervised</b>                  |           |           |             |                  |                  |
|  |           |           | 56.8        | 48.1             | 63.1             |
| ✓  |           |           | 57.1        | 47.8             | <b>63.8</b>      |
| ✓  | ✓         |           | 58.1        | 50.4             | 63.7             |
| ✓  | ✓         | ✓         | <b>58.8</b> | <b>52.1</b>      | 63.7             |
| <b>Weakly Supervised (No mask label)</b> |           |           |             |                  |                  |
| ✓  | ✓         |           | 55.7        | 45.2             | 63.3             |

# Scaling up ML for Autonomy

## Behavior Cloning and its Limitations

## Real-time Panoptic Segmentation

*Real-Time Panoptic Segmentation from Dense Detections,  
R. Hou\*, J. Li\*, A. Bhargava, A. Raventos et al, CVPR'20 (oral)*

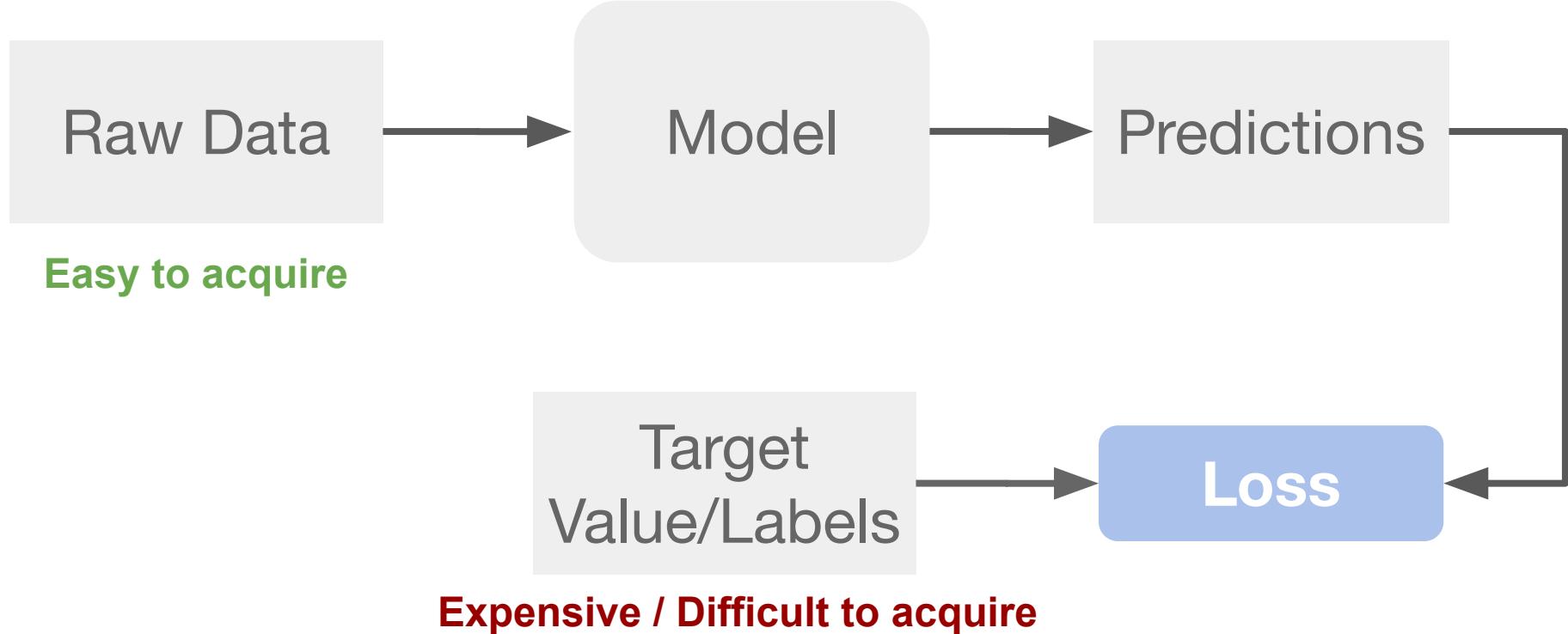
## Self-Supervised Pseudo-Lidar Networks

## Auto-labeling via Differentiable Rendering

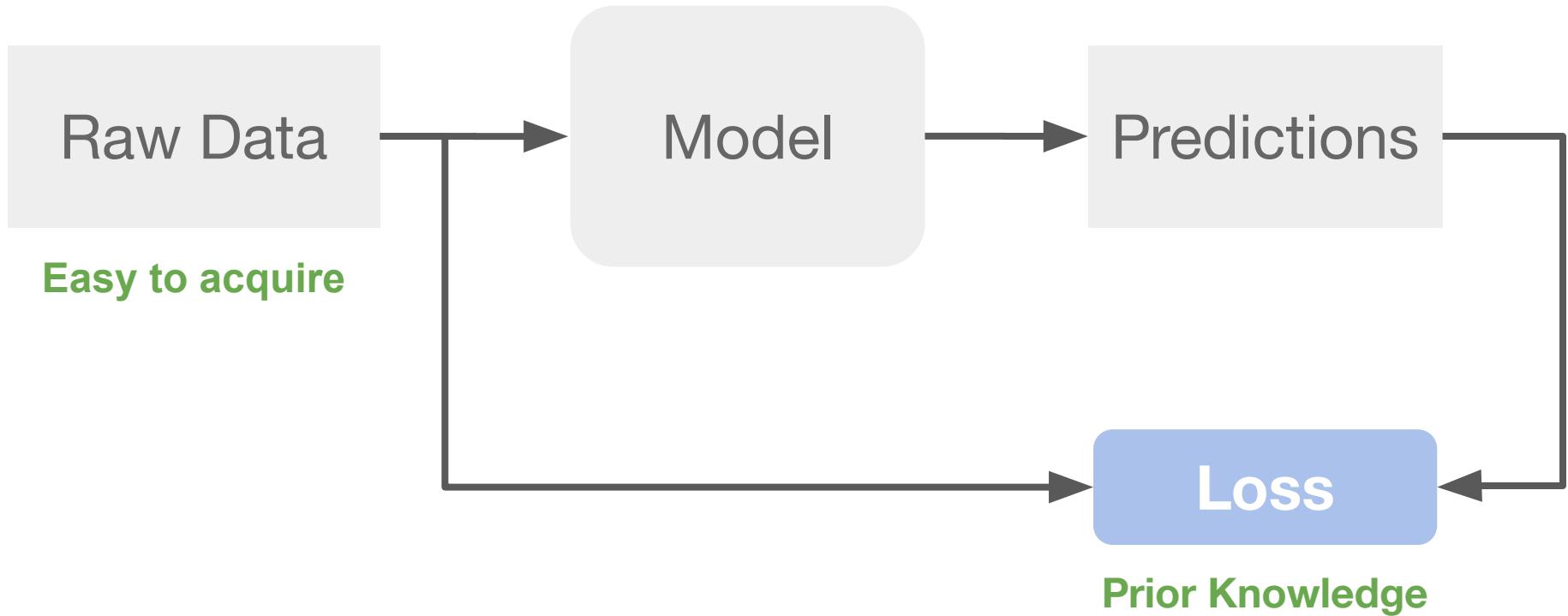
# Scaling up ML for Autonomy

- Behavior Cloning and its Limitations
- Real-time Panoptic Segmentation
- Self-Supervised Pseudo-Lidar Networks**
  - 3D Packing for Self-Supervised Monocular Depth Estimation*
  - V. Guizilini, R. Ambrus, S. Pillai et al, CVPR'20 (oral)*
- Auto-labeling via Differentiable Rendering

# Supervised Learning

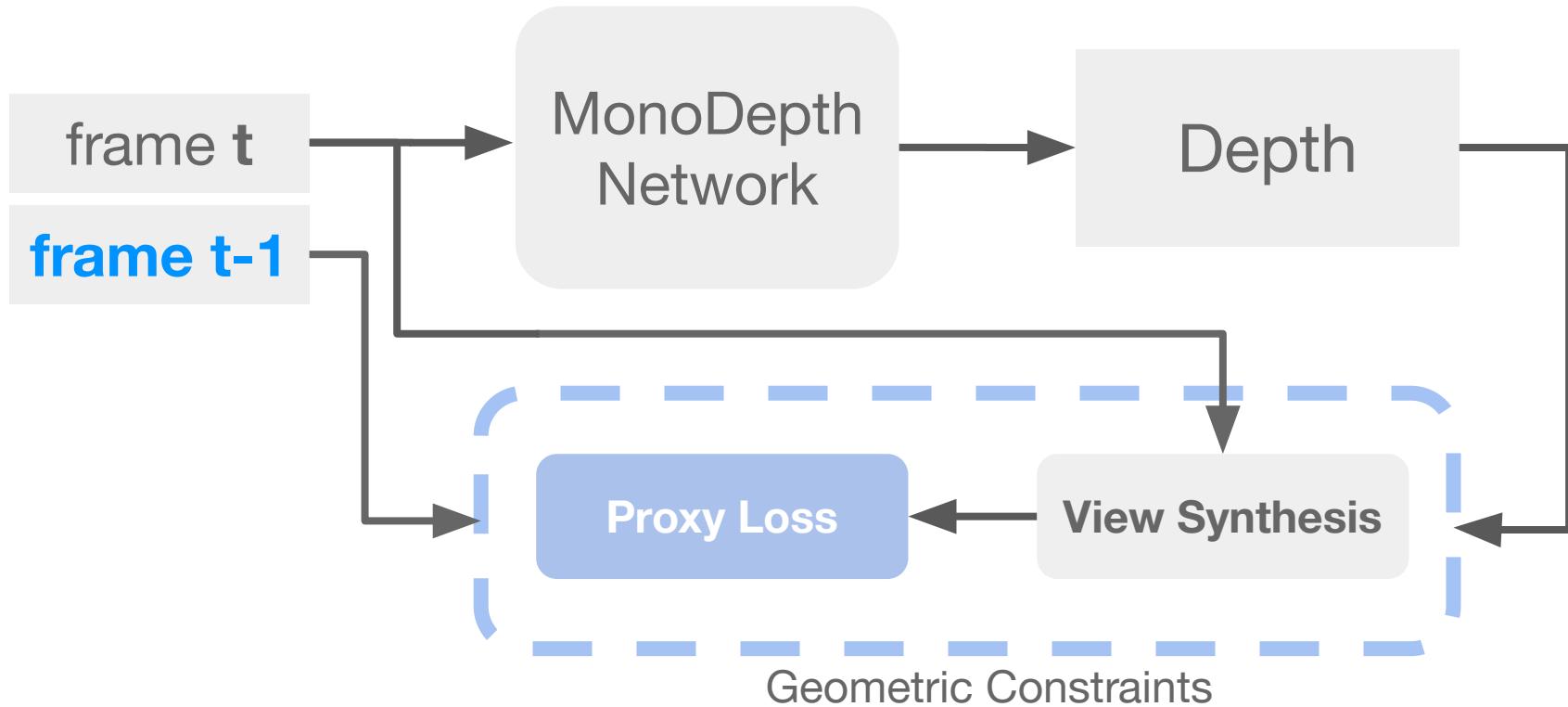


# Self-Supervised Learning

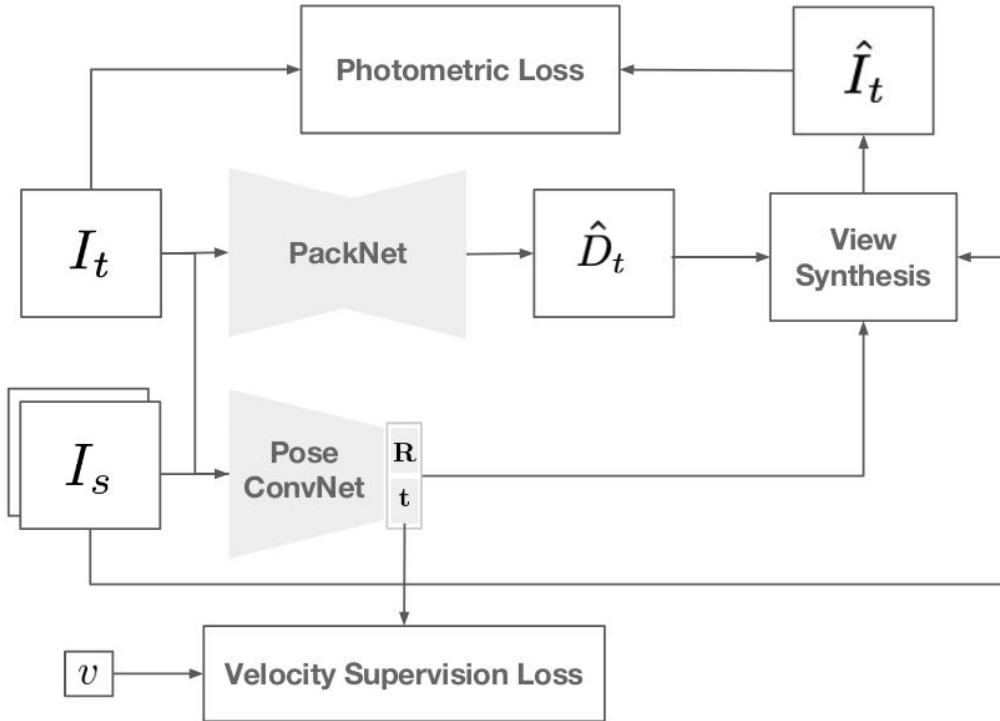


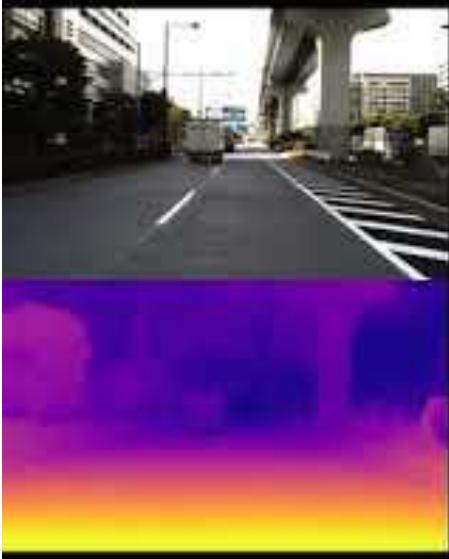
# Self-Supervised Structure-from-Motion (SfM)

Monocular  
Video



# Self-Supervised Structure-from-Motion (SfM)





- No LiDAR information is used at training or test time
- Samples shown were not seen during training

# PackNet: Pack it, don't pool it

|                        | Layer Description   | K | Output Tensor Dim.            |
|------------------------|---|---|-------------------------------|
| #0                     | Input RGB image   |   | $3 \times H \times W$         |
| <b>Encoding Layers</b> |   |   |                               |
| #1                     | Conv2d  | 5 | $64 \times H \times W$        |
| #2                     | Conv2d → <b>Packing</b>   | 7 | $64 \times H/2 \times W/2$    |
| #3                     | ResidualBlock (x2) → <b>Packing</b>                                   | 3 | $64 \times H/4 \times W/4$    |
| #4                     | ResidualBlock (x2) → <b>Packing</b>                                   | 3 | $128 \times H/8 \times W/8$   |
| #5                     | ResidualBlock (x3) → <b>Packing</b>                                   | 3 | $256 \times H/16 \times W/16$ |
| #6                     | ResidualBlock (x3) → <b>Packing</b>                                   | 3 | $512 \times H/32 \times W/32$ |
| <b>Decoding Layers</b> |   |   |                               |
| #7                     | <b>Unpacking</b> (#6) → Conv2d ( $\oplus$ #5)                         | 3 | $512 \times H/16 \times W/16$ |
| #8                     | <b>Unpacking</b> (#7) → Conv2d ( $\oplus$ #4)                         | 3 | $256 \times H/8 \times W/8$   |
| #9                     | InvDepth (#8)   | 3 | $1 \times H/8 \times W/8$     |
| #10                    | <b>Unpacking</b> (#8) → Conv2d ( $\oplus$ #3 $\oplus$ Upsample(#9))   | 3 | $128 \times H/4 \times W/4$   |
| #11                    | InvDepth (#10)  | 3 | $1 \times H/4 \times W/4$     |
| #12                    | <b>Unpacking</b> (#10) → Conv2d ( $\oplus$ #2 $\oplus$ Upsample(#11)) | 3 | $64 \times H/2 \times W/2$    |
| #13                    | InvDepth (#12)  | 3 | $1 \times H/2 \times W/2$     |
| #14                    | <b>Unpacking</b> (#12) → Conv2d ( $\oplus$ #1 $\oplus$ Upsample(#13)) | 3 | $64 \times H \times W$        |
| #15                    | InvDepth (#14)  | 3 | $1 \times H \times W$         |



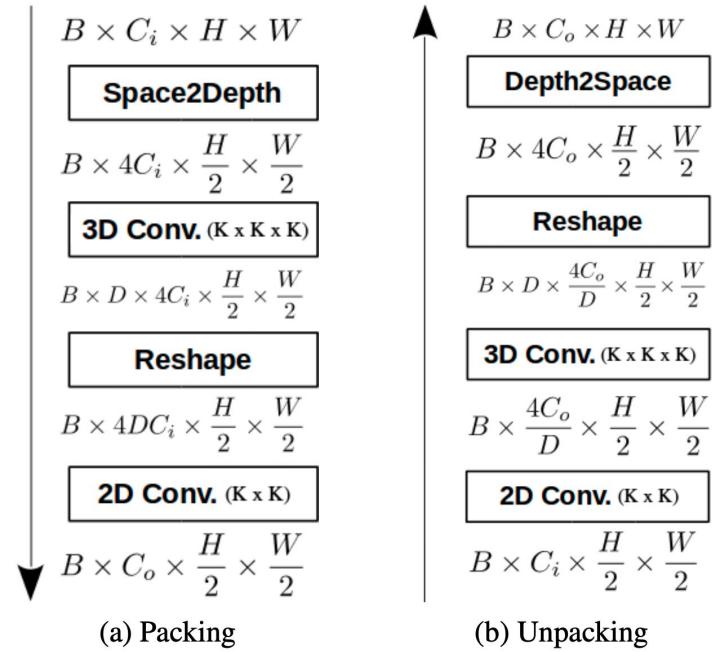
(a) Input Image



(b) Max Pooling +  
Bilinear Upsample

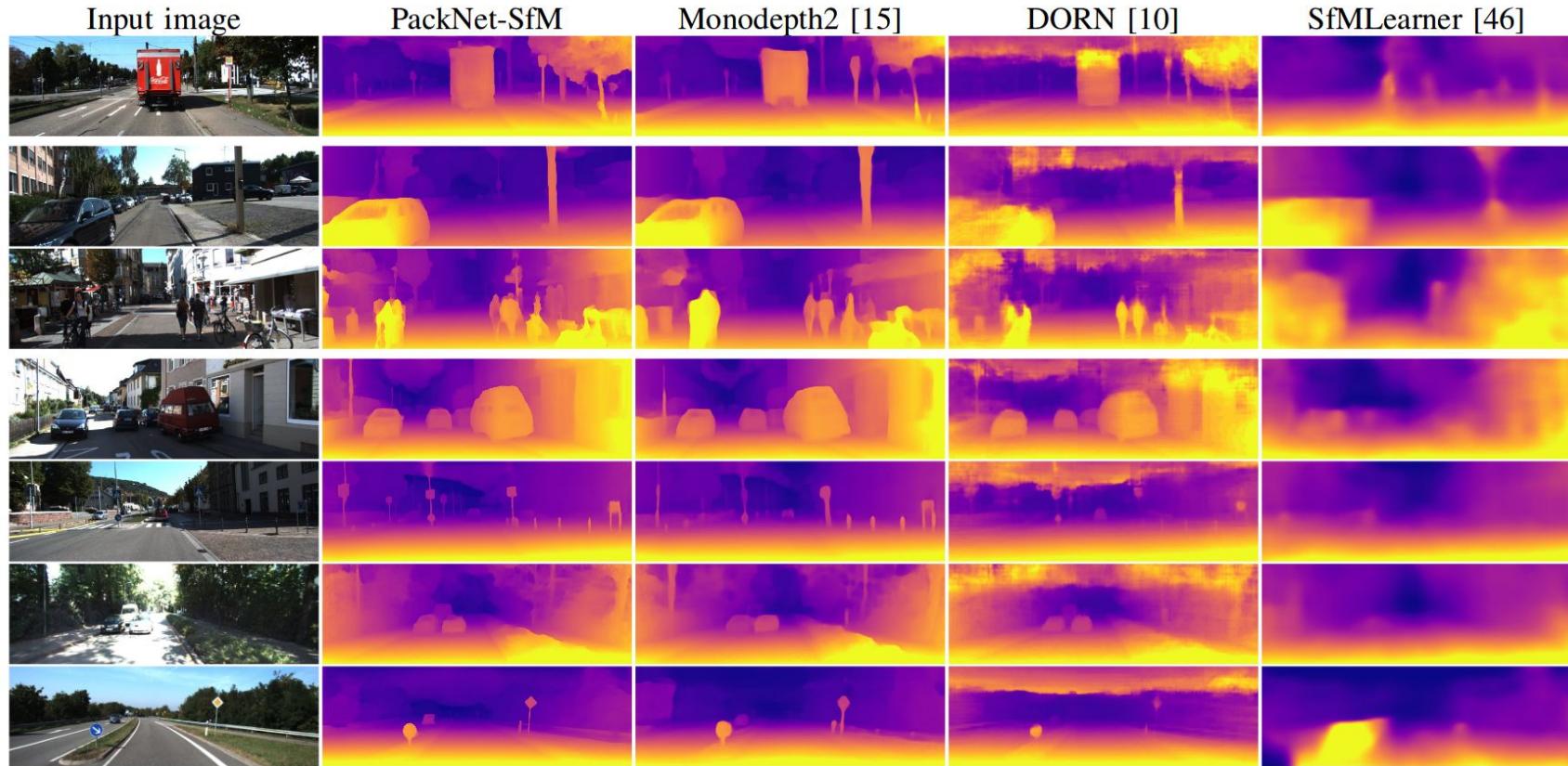


(c) Pack + Unpack



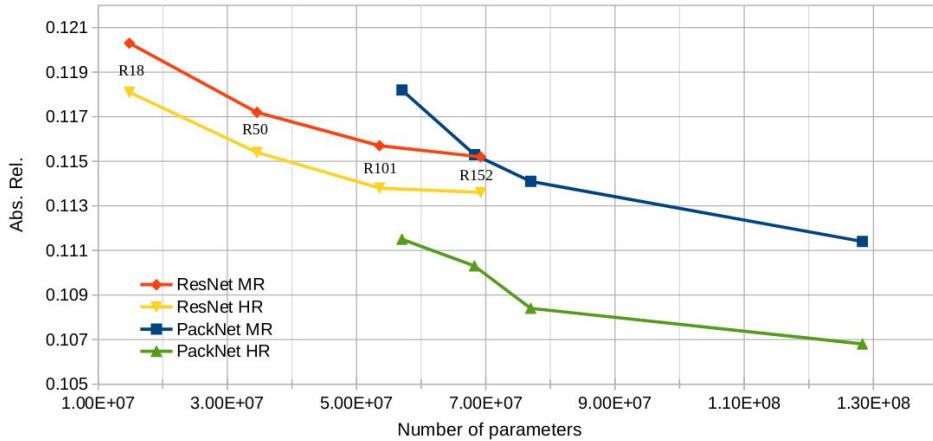


# Experimental Results (KITTI)



# Experimental Results

*Better use of network capacity...*

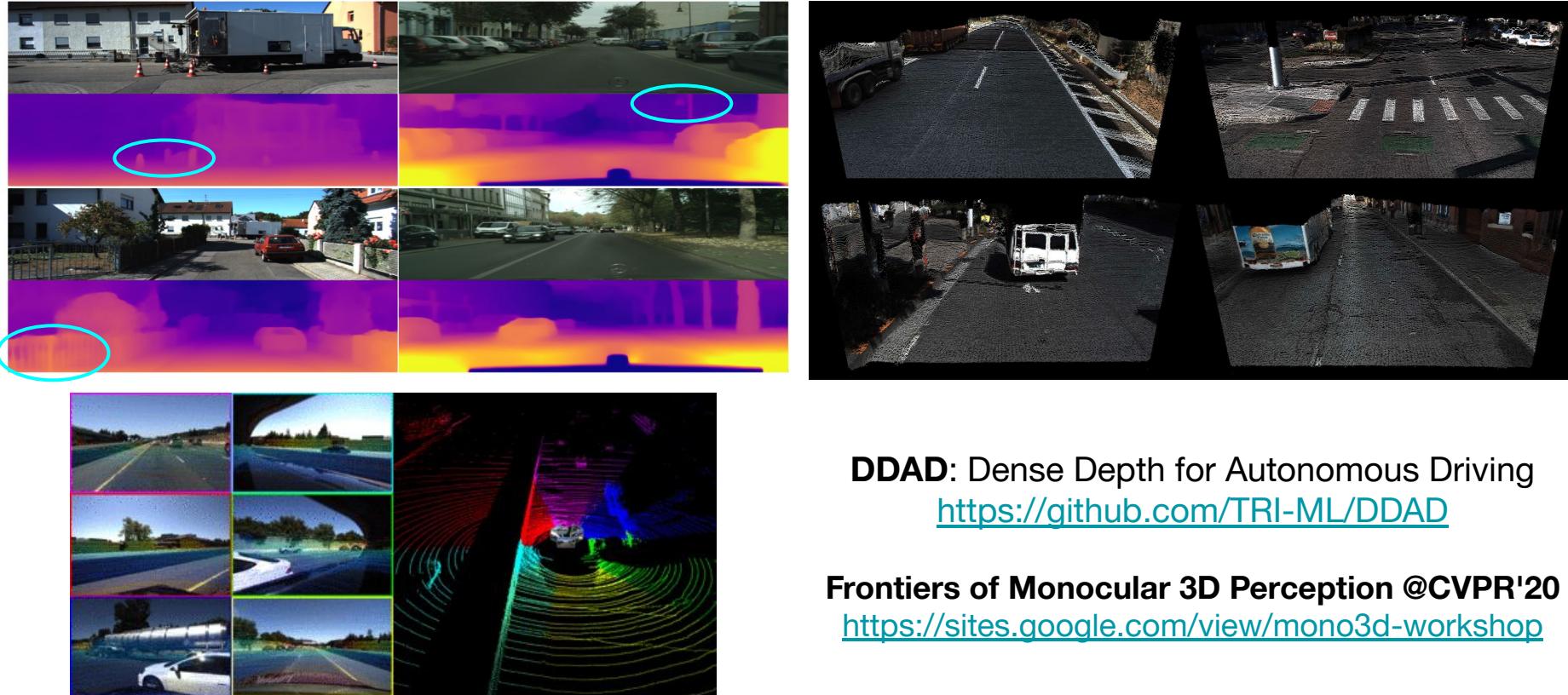


**And better generalization!**  
(KITTI → NuScenes)

| Depth Network                    | Abs          | Rel          | Sq Rel       | RMSE         | RMSE <sub>log</sub> | $\delta < 1.25$ |
|----------------------------------|--------------|--------------|--------------|--------------|---------------------|-----------------|
| ResNet18                         | 0.133        | 1.023        | 5.123        | 0.211        | 0.845               |                 |
| ResNet18 <sup>‡</sup>            | 0.120        | 0.896        | 4.869        | 0.198        | 0.868               |                 |
| ResNet50                         | 0.127        | 0.977        | 5.023        | 0.205        | 0.856               |                 |
| ResNet50 <sup>‡</sup>            | 0.117        | 0.900        | 4.826        | 0.196        | 0.873               |                 |
| PackNet18                        | 0.118        | 0.802        | 4.656        | 0.194        | 0.868               |                 |
| PackNet50                        | 0.114        | 0.818        | 4.621        | 0.190        | 0.875               |                 |
| PackNet-SfM<br>(w/o pack/unpack) | 0.122        | 0.880        | 4.816        | 0.198        | 0.864               |                 |
| PackNet-SfM<br>(w/o 3D convs.)   | 0.118        | 0.922        | 4.831        | 0.195        | 0.872               |                 |
| <b>PackNet-SfM</b>               | <b>0.111</b> | <b>0.785</b> | <b>4.601</b> | <b>0.189</b> | <b>0.878</b>        |                 |

| Method                | Abs          | Rel          | Sq Rel       | RMSE         | RMSE <sub>log</sub> | $\delta < 1.25$ |
|-----------------------|--------------|--------------|--------------|--------------|---------------------|-----------------|
| ResNet18              | 0.218        | 2.053        | 8.154        | 0.355        | 0.650               |                 |
| ResNet18 <sup>‡</sup> | 0.212        | 1.918        | 7.958        | 0.323        | 0.674               |                 |
| ResNet50              | 0.216        | 2.165        | 8.477        | 0.371        | 0.637               |                 |
| ResNet50 <sup>‡</sup> | 0.210        | 2.017        | 8.111        | 0.328        | 0.697               |                 |
| <b>PackNet-SfM</b>    | <b>0.187</b> | <b>1.852</b> | <b>7.636</b> | <b>0.289</b> | <b>0.742</b>        |                 |

# Experimental Results



**DDAD:** Dense Depth for Autonomous Driving  
<https://github.com/TRI-ML/DDAD>

**Frontiers of Monocular 3D Perception @CVPR'20**  
<https://sites.google.com/view/mono3d-workshop>

# Scaling up ML for Autonomy

- Behavior Cloning and its Limitations
- Real-time Panoptic Segmentation
- Self-Supervised Pseudo-Lidar Networks**
  - 3D Packing for Self-Supervised Monocular Depth Estimation*
  - V. Guizilini, R. Ambrus, S. Pillai et al, CVPR'20 (oral)*
- Auto-labeling via Differentiable Rendering

# Scaling up ML for Autonomy

Behavior Cloning and its Limitations  
Real-time Panoptic Segmentation  
Self-Supervised Pseudo-Lidar Networks  
**Auto-labeling via Differentiable Rendering**

*Autolabeling 3D Objects with Differentiable Rendering of  
SDF Shape Priors, S. Zakharov, W. Kehl\* et al, CVPR'20 (oral)*

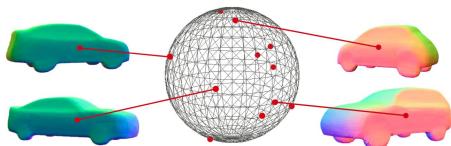
# Auto-labeling in 3D

Input: image, point cloud, 2d bounding boxes

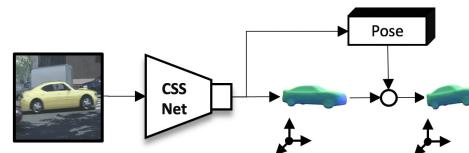
Output: 3d boxes with pose + shape

Goal: use auto-labels instead of manual ones

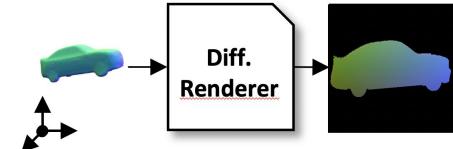
Shape Representation

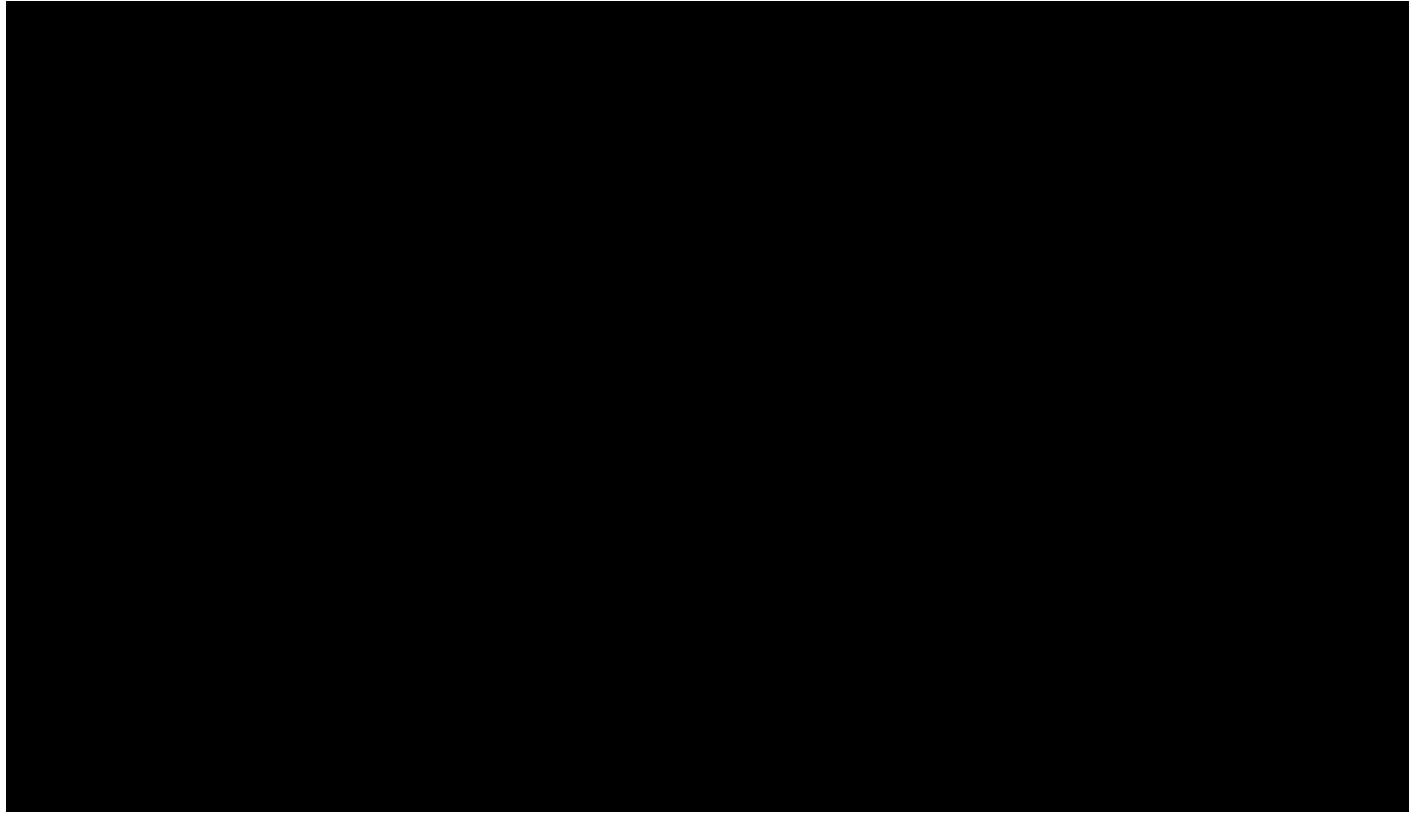


Pose/Shape Estimator



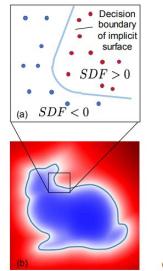
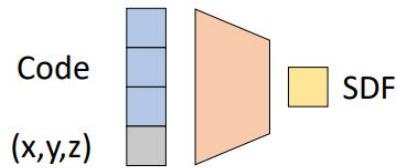
Differentiable Renderer



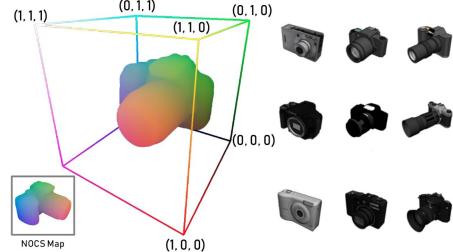


# Coordinate Shape Space (CSS): DeepSDF + NOCS

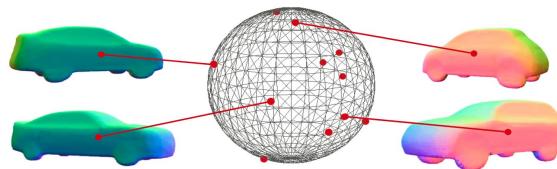
DeepSDF network  
maps  $[x,y,z,\text{code}]$  to  
an SDF value



Normalized Object  
Coordinate Space  
(NOCS)

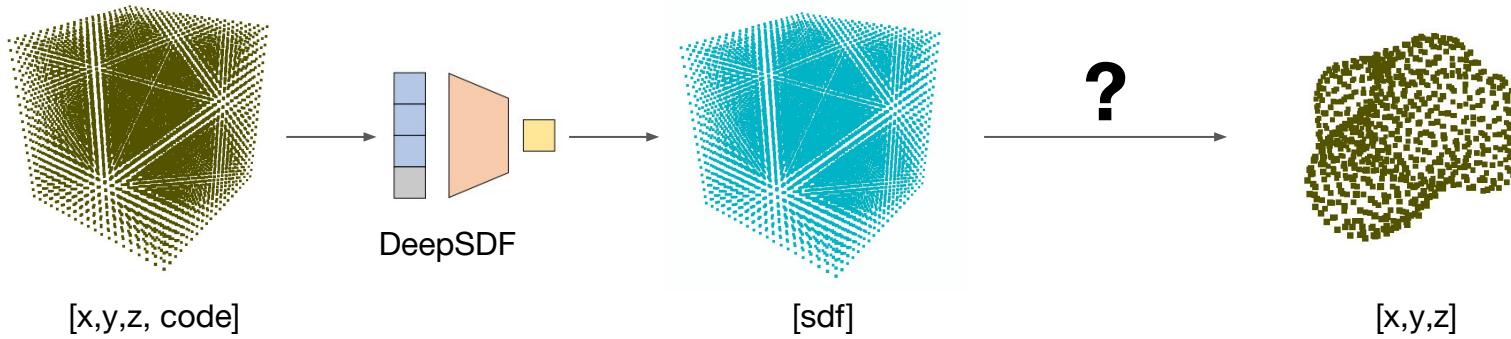


3D models represented in the shape space



3D latent code represents a unique car shape

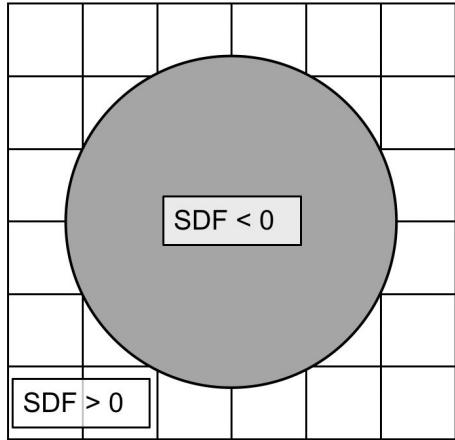
# From SDF to Coordinates?



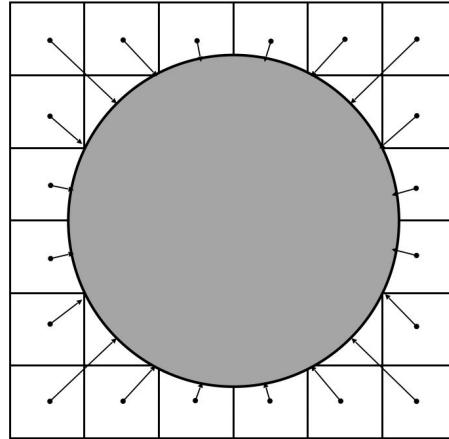
## How to render SDF?

- Marching Cubes (differentiable?)
- Raycasting (slow)

# Zero-Isosurface Projection



$$\mathbf{s} = f(\mathit{grid})$$



$$\mathbf{n} = \frac{\partial f(\mathit{grid})}{\partial \mathit{grid}}$$



$$\mathit{grid} - \frac{\partial f(\mathit{grid})}{\partial \mathit{grid}} f(\mathit{grid})$$

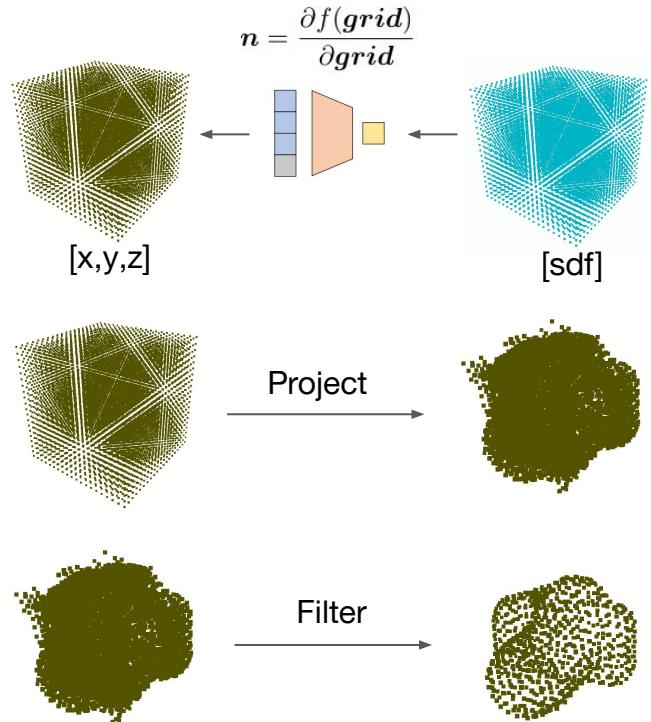
# Zero-Isosurface Projection

1. Project the grid points to the surface using the *SDF values* and the *analytically estimated normals*:

$$p = \text{grid} - \frac{\partial f(\text{grid})}{\partial \text{grid}} f(\text{grid})$$

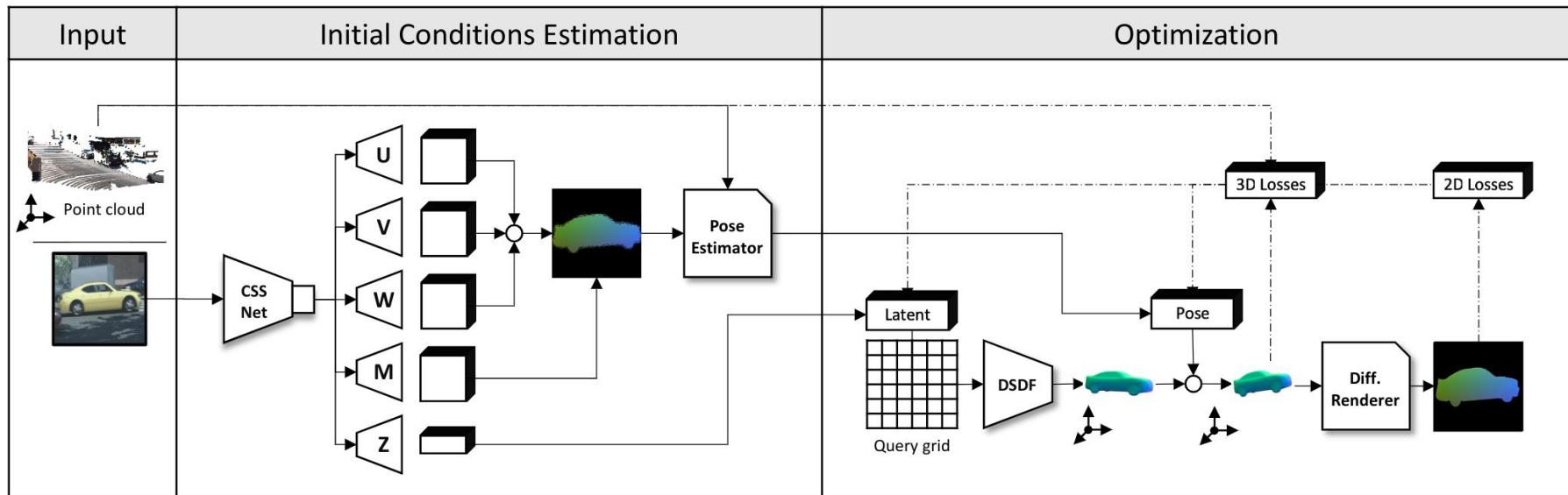
2. Mask the points that are far from the surface:

$$p_{\text{masked}} = p, \text{ where } |f(\text{grid})| < 0.1$$

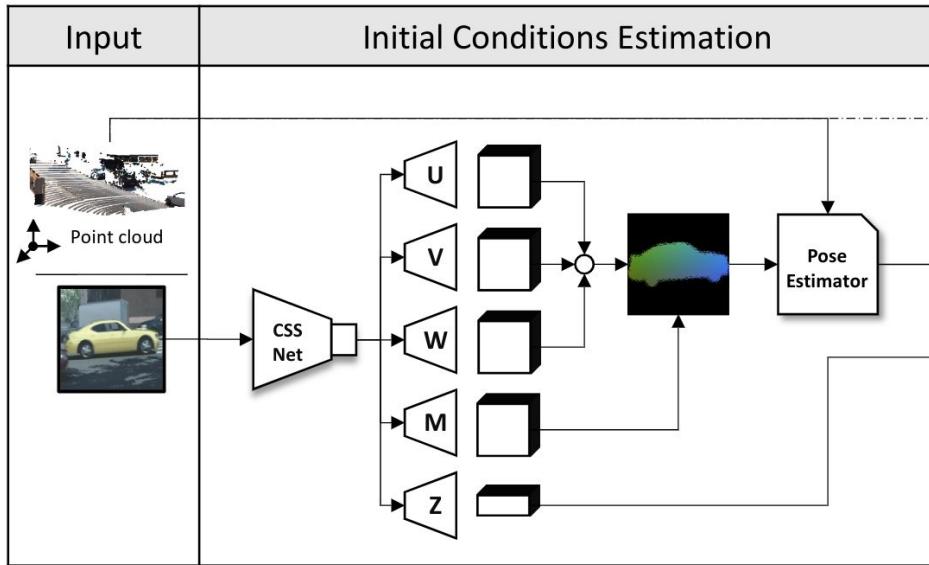


**Differentiable all the way back to the latent vector!**

# Auto-labeling Pipeline



# Auto-labeling Pipeline



## CSS Network (R18-based):

**U, V, W:** Normalized Object Coordinates (NOCS) for each pixel  
**M:** Object mask  
**z:** Latent shape vector

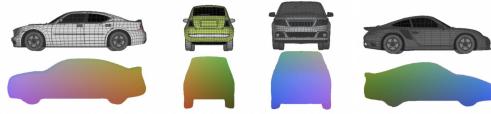


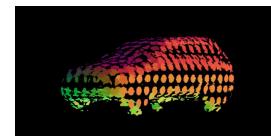
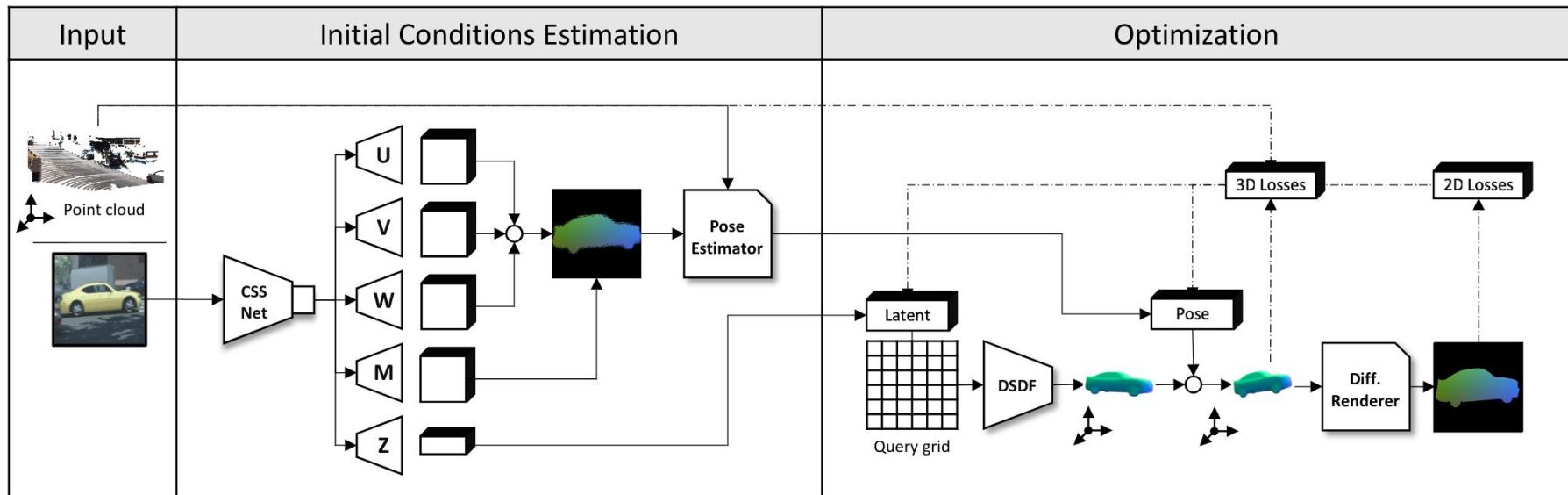
Figure 2: CSS representation. Top: Car models from the PD dataset [1]. Bottom: The same cars in the CSS representation: decoded shape vector  $z$  colored with NOCS.

Pose estimator uses either:

PnP (2D-based)

Kabsch / Procrustes (3D-based)

# Auto-labeling Pipeline



# Training: Domain Adaptation



Parallel Domain (Training data)



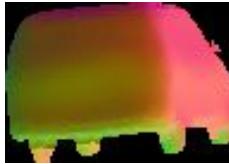
KITTI (Testing data)

# Training: Domain Randomization

- 2D:
  - Transforms:
    - Random rotation
    - Random horizontal flip
    - Random resized crop
  - Noise:
    - Color jitter: brightness, contrast, saturation, hue
- 3D:
  - Phong Lighting:
    - Random ambient, diffuse and specular lights from pre-computed normals



RGB



Normals

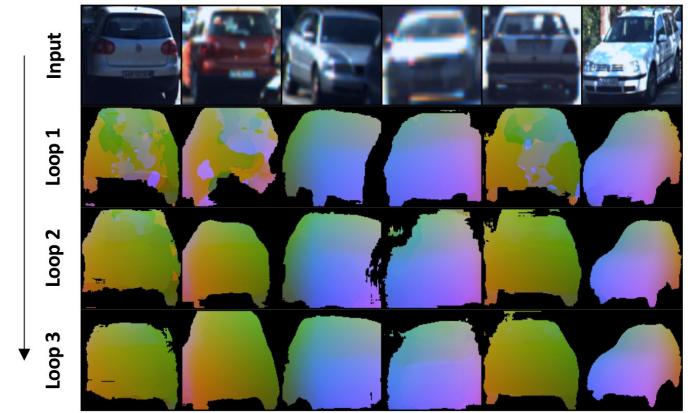
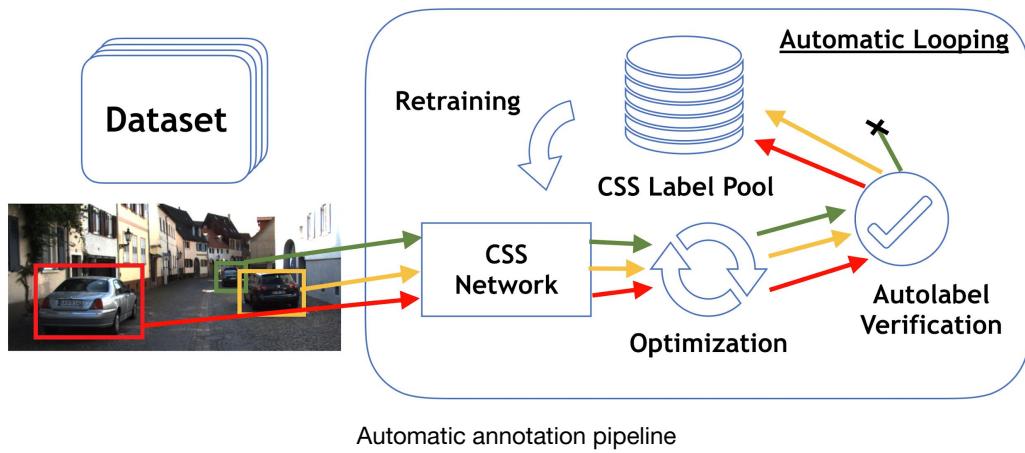


Color jitter



+Random  
transforms

# Curriculum Learning



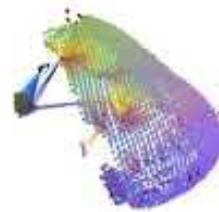
CSS Network prediction quality of our network over consecutive loops for the same patch

# Qualitative Results



2D alignment

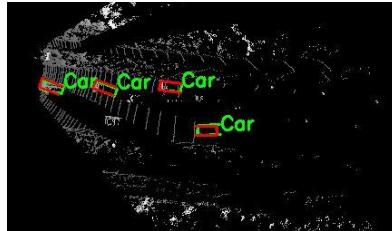
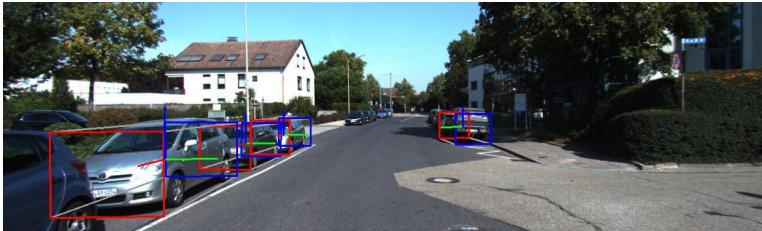
3D alignment



# Quantitative Results

| Method                              | 2D AP @ 0.5/0.7 |             | 3D AP @ 0.5/0.7 |             | BEV AP @ 0.5/0.7 |              |
|-------------------------------------|-----------------|-------------|-----------------|-------------|------------------|--------------|
|                                     | Easy            | Moderate    | Easy            | Moderate    | Easy             | Moderate     |
| PointPillars [20] (Original Labels) | - / -           | - / -       | 94.8 / 81.1     | 92.4 / 68.2 | 95.1 / 92.1      | 95.1 / 84.7  |
| PointPillars [20] (Autolabels)      | - / -           | - / -       | 90.7 / 22.4     | 71.1 / 13.3 | 94.9 / 81.0      | 88.5 / 59.8  |
| MonoDIS [35] (Original Labels)      | 96.1 / 95.5     | 92.6 / 86.5 | 45.7 / 11.0     | 32.9 / 7.1  | 52.4 / 17.7      | 37.2 / 11.9  |
| MonoDIS [35] (Autolabels)           | 96.7 / 85.8     | 86.2 / 67.6 | 32.9 / 1.23     | 22.1 / 0.54 | 51.1 / 15.7      | 34.5 / 10.52 |

Table 2: We compare the performance of 3D object detectors trained on true KITTI labels vs. our autolabels. On the BEV metric, **the detectors trained on autolabels alone achieve results equal to the current state-of-the-art**. On the 3D AP metric, both autolabel trained detectors achieve competitive results at the IoU 0.5 threshold.

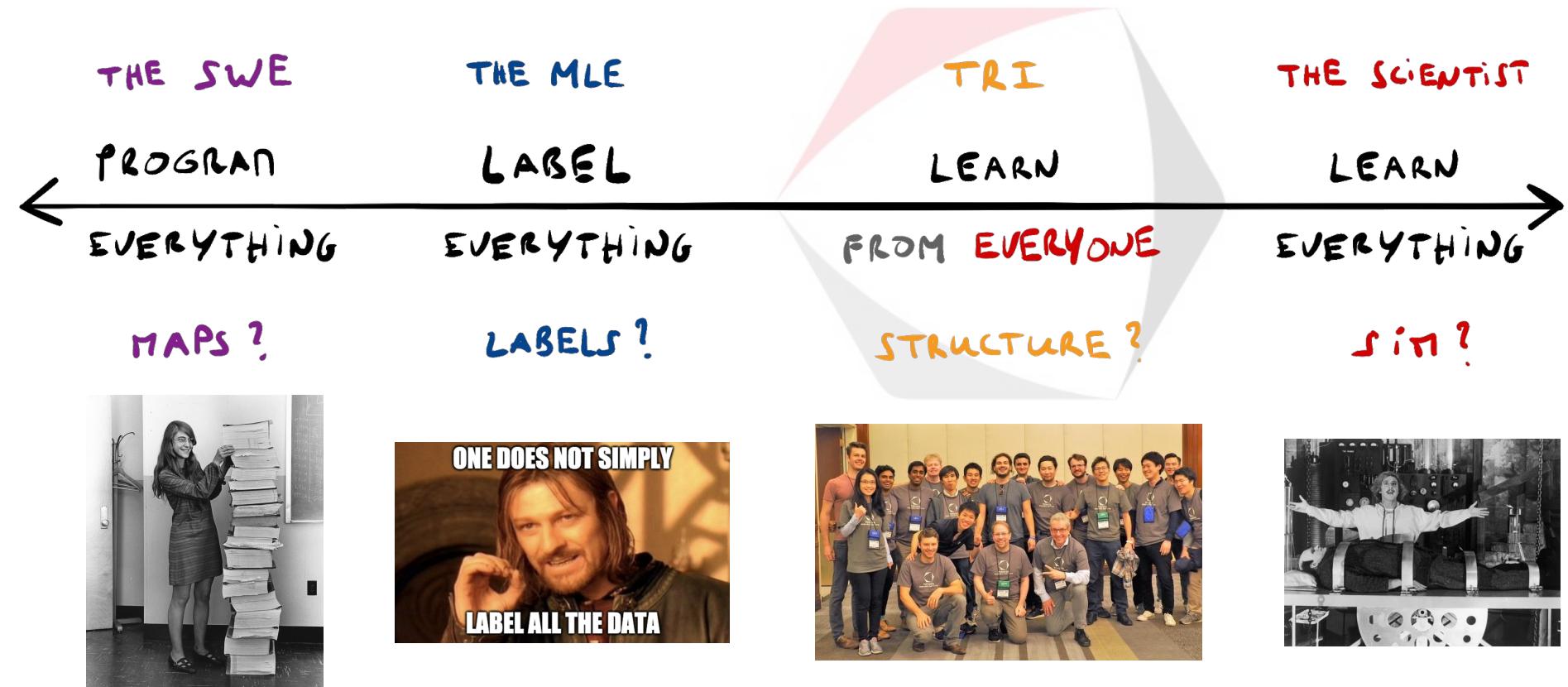


# Scaling up ML for Autonomy

Behavior Cloning and its Limitations  
Real-time Panoptic Segmentation  
Self-Supervised Pseudo-Lidar Networks  
**Auto-labeling via Differentiable Rendering**

*Autolabeling 3D Objects with Differentiable Rendering of  
SDF Shape Priors, S. Zakharov, W. Kehl\* et al, CVPR'20 (oral)*

# World-scale Autonomy?



## Behavior: leverage large scale Demonstrations

*Exploring the Limitations of Behavior Cloning for Autonomous Driving, ICCV'19 (oral)*

*Spatiotemporal Relationship Reasoning for Pedestrian Intent Prediction, RA-L & ICRA'20*

*It Is Not the Journey but the Destination: Endpoint Conditioned Trajectory Prediction, arXiv:2004.02025*

*Reinforcement Learning based Control of Imitative Policies for Near-Accident Driving, coming soon*

*Risk-Sensitive Sequential Action Control with Multi-Modal Human Trajectory Forecasting [...], coming soon*

*Driving Through Ghosts: Behavioral Cloning with False Positives, coming soon*

## Supervised Learning: efficiently use available Labels

*ROI-10D: Monocular Lifting of 2D Detection to 6D Pose and Metric Shape, CVPR'19*

*Learning Imbalanced Datasets with Label-Distribution-Aware Margin Loss, NeurIPS'19*

*Learning to Fuse Things and Stuff, arXiv:1812.01192*

*Spatio-Temporal Graph for Video Captioning with Knowledge Distillation, CVPR'20*

## Real-Time Panoptic Segmentation from Dense Detections, CVPR'20 (oral)

*Hierarchical Lovász Embeddings for Proposal-free Panoptic Segmentation, coming soon*

*Unsupervised Estimation of Segmentation Difficulty, coming soon*

## Geometry: Self / Semi-Supervised Pseudo-LiDAR and SfM

*SuperDepth: Self-Supervised, Super-Resolved Monocular Depth Estimation, ICRA'19*

*Robust Semi-Supervised Monocular Depth Estimation with Reprojected Distances, CoRL'19*

*Two Stream Networks for Self-Supervised Ego-Motion Estimation, CoRL'19*

*Semantically-Guided Representation Learning for Self-Supervised Monocular Depth, ICLR'20*

*Neural Outlier Rejection for Self-Supervised Keypoint Learning, ICLR'20*

*Self-Supervised 3D Keypoint Learning for Ego-motion Estimation, arxiv:1912.03426*

## 3D Packing for Self-Supervised Monocular Depth Estimation, CVPR'20 (oral)

*Self-Supervised Neural Camera Models, coming soon*

## Simulation: Domain Adaptation, Differentiable Rendering, RL

*SPIGAN: Privileged Adversarial Learning from Simulation, ICLR'19*

*DeceptionNet: Network-Driven Domain Randomization, ICCV'19*

*Generating Human Action Videos by Coupling 3D Game Engines and Probabilistic Graphical Models, IJCV'20*

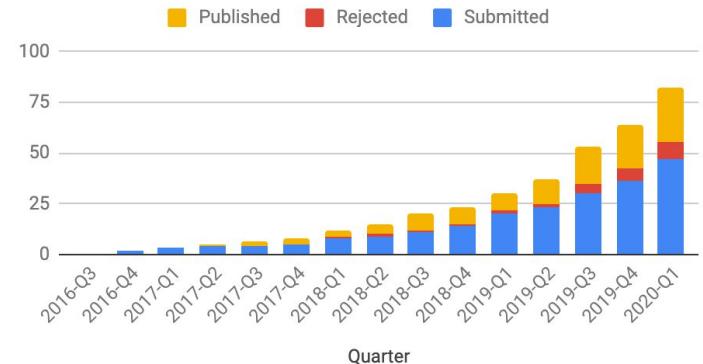
## Autolabeling 3D Objects with Differentiable Rendering of SDF Shape Priors, CVPR'20 (oral)

*Self-Supervised Differentiable Rendering for Monocular 3D Object Detection, coming soon*

*Behaviorally Diverse Traffic Simulation via Reinforcement Learning, coming soon*

*Discovering Avoidable Planner Failures [...] in Behaviorally Diverse Simulation, coming soon*

## ML Publications History (cumulative)



## Upcoming workshops co-organized by TRI

### ICML: AI for Autonomous Driving (AIAD)

<https://sites.google.com/view/aiad2020>

### ECCV: Perception for Autonomous Driving (PAD)

<https://sites.google.com/view/pad2020>

## Upcoming TRI Dataset Releases

### STIP: Stanford-TRI Intent Prediction

<http://stip.stanford.edu/>

### DDAD: Dense Depth for Autonomous Driving

<https://github.com/TRI-ML/DDAD>

# Scaling up ML for Autonomy

## Behavior Cloning and its Limitations

More data & params → SotA policy but...

## Real-time Panoptic from Bounding Boxes

SotA, 4x faster, weak sup. = 95% strong

## Self-Supervised Pseudo-Lidar Networks

Self sup > sup! Don't pool it: Pack it.

## Auto-labeling via Differentiable Rendering

Diff. shape priors + geometry ~ labels