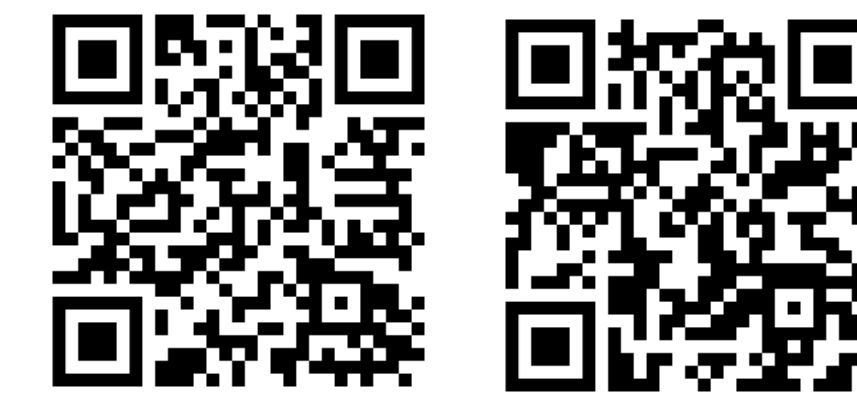# Pseudo-LiDAR from Visual Depth Estimation:
# Bridging the Gap in 3D Object Detection for Autonomous Driving

Yan Wang, Wei-Lun (Harry) Chao, Divyansh Garg, Bharath Hariharan, Mark Campbell, Kilian Q. Weinberger
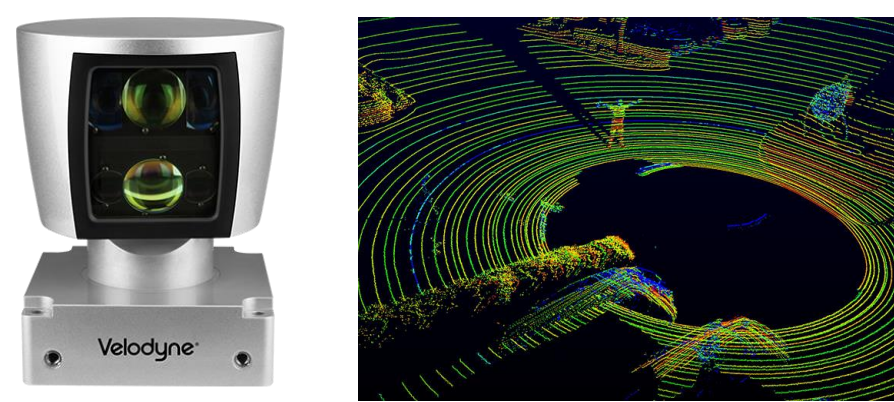
CVPR — LONG BEACH CALIFORNIA — June 16-20, 2019

## Highlights

- **Propose an image-based 3D detection framework:** converting **image-based depth maps** to **pseudo-LiDAR representation** enables existing LiDAR-based 3D object detectors
- **Achieve a 45% AP$_{3D}$ on the KITTI benchmark, almost a 350% improvement over the previous SOTA**
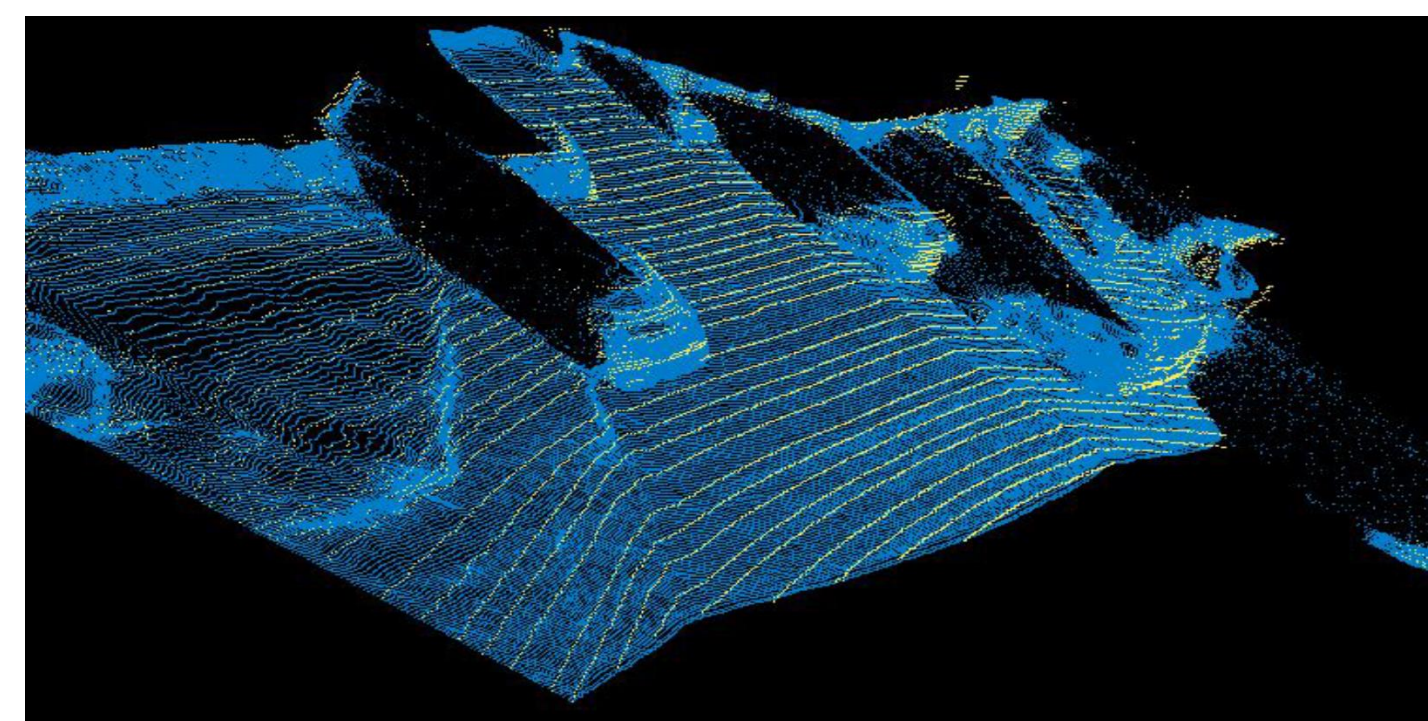
## Introduction

- 3D object detection is essential for autonomous driving.
- Most approaches rely on LiDAR for precise depths, but:
  - ➤ Expensive (64-line = $75K USD)
  - ➤ Over-reliance is risky.
  - ➤ Alternatives are needed.
- Image-based approaches fall far behind (10% vs. 74% AP$_{3D}$), commonly attributed to *poor image-base depth estimation*.

## Is image-based depth accurate?

- Image-based depth maps $Z$ can be transformed to 3D points

(depth) $z = Z(u,v)$

(width) $x = \frac{(u-c_U) \times z}{f_U}$

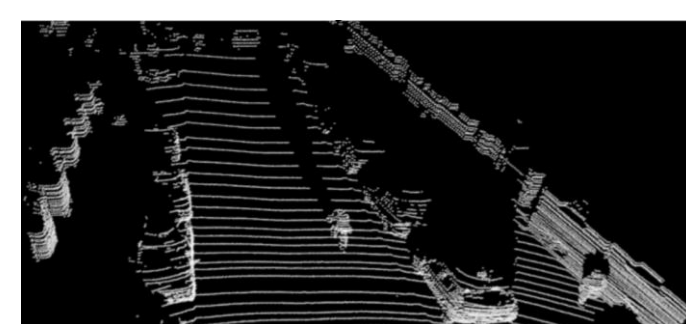(height) $y = \frac{(v-c_V) \times z}{f_V}$

$c_U, c_V$: image center
$f_U, f_V$: focal lengths

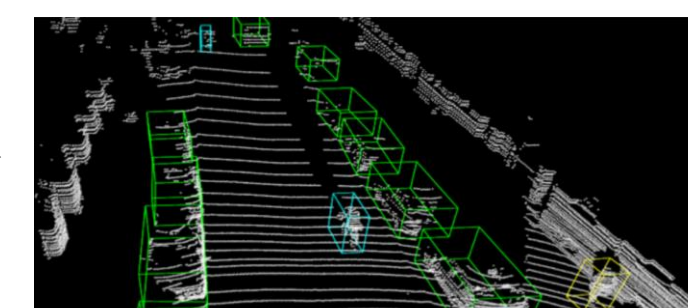- **Stereo depth vs. LiDAR: points are surprisingly consistent!**

## Data representation matters!

- LiDAR-based 3D detectors

Point cloud or bird's-eye view (BEV)
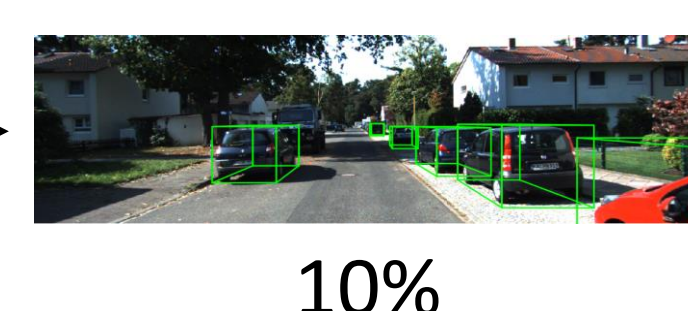
[VoxelNet, 2018]
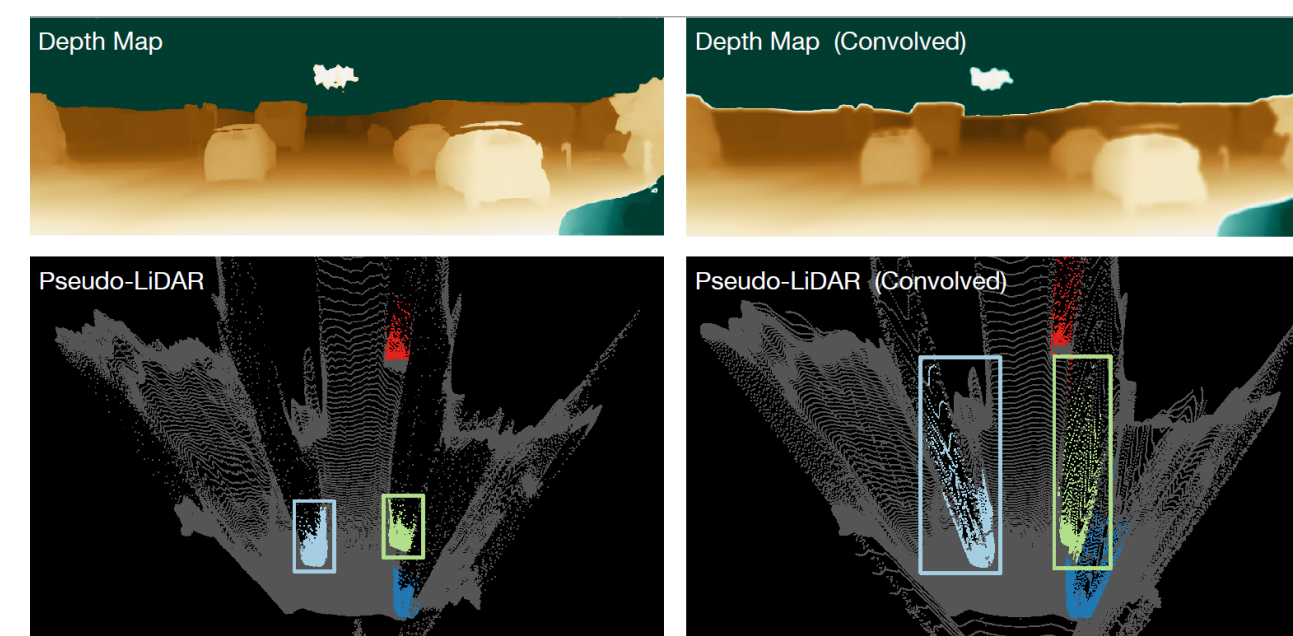
74%

- Image-based 3D detectors

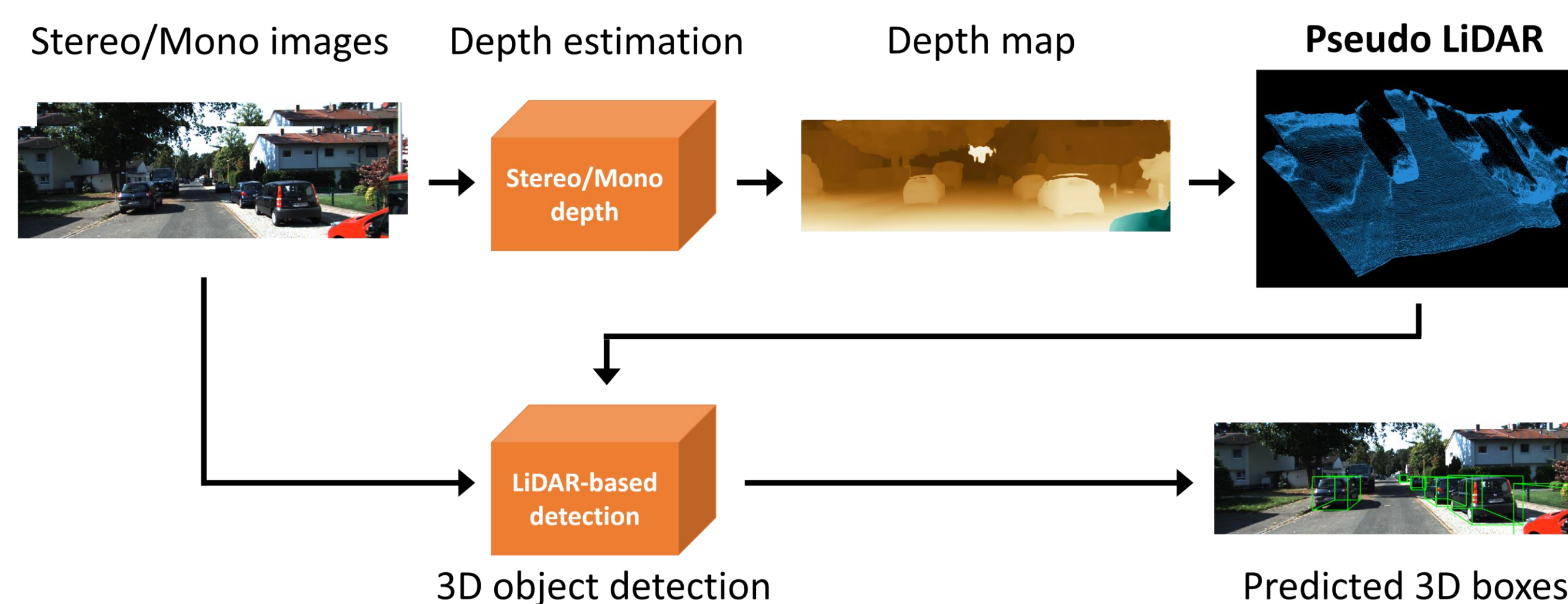Frontal view, followed by the 2D detection pipeline

10%

- Issues with convolution from the frontal view:
  - ➤ Object sizes vary with depth.
  - ➤ Neighboring pixels may be far-away in 3D, making it hard for convolutional networks to precisely localize objects in 3D.
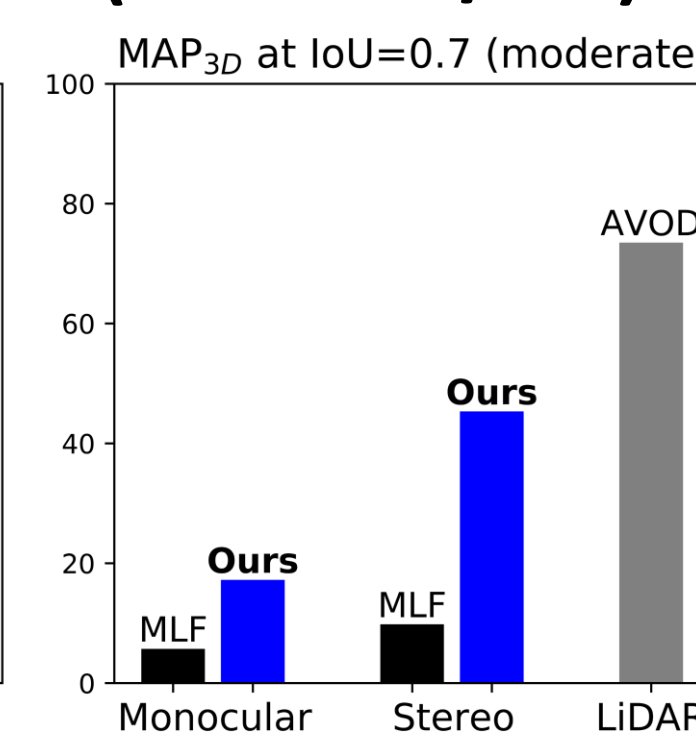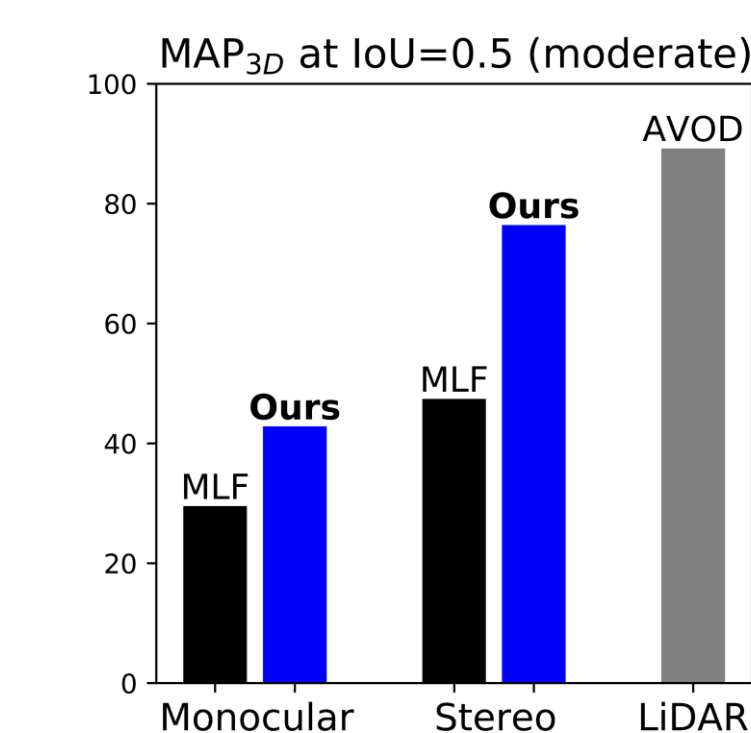
Depth Map — Depth Map (Convolved)
Pseudo-LiDAR — Pseudo-LiDAR (Convolved)

## Proposed pseudo-LiDAR framework

Stereo/Mono images — Depth estimation — Depth map — Pseudo LiDAR

Stereo/Mono depth

LiDAR-based detection

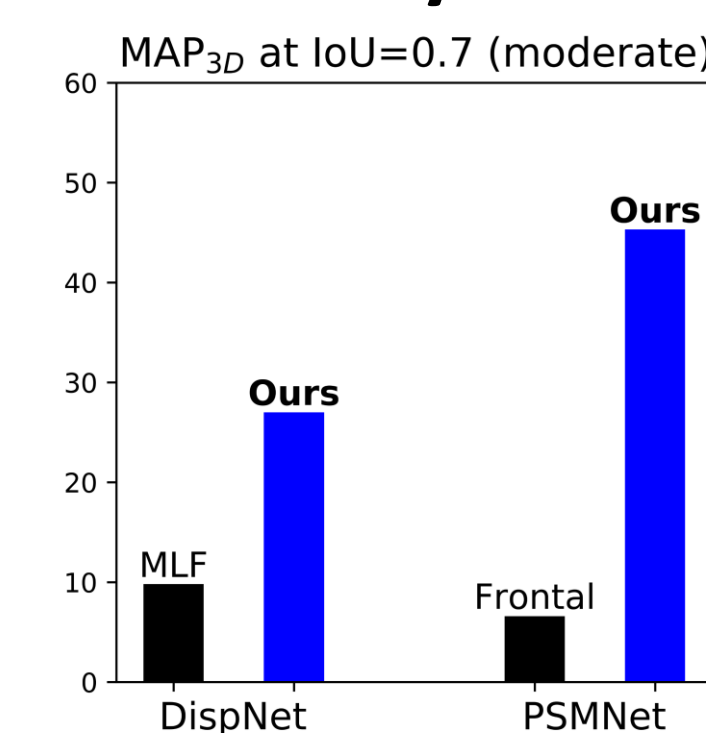3D object detection — Predicted 3D boxes

## Experiments

- **Dataset:** KITTI object detection (4K/4K/8K images for train/val/test), focusing on "car"
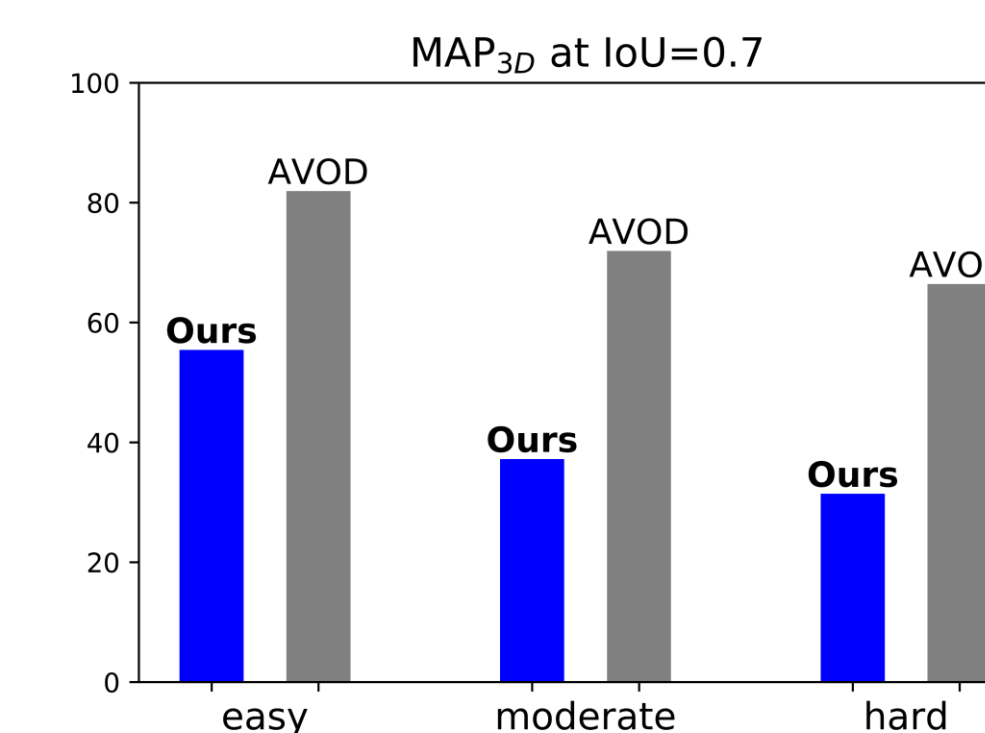- **Our approach:** PSMNet [1]/Dorn [2] for Stereo/monocular depth + AVOD detector [3]

**Validation results (IoU = 0.5/0.7)**

MAP$_{3D}$ at IoU=0.5 (moderate)
(bars: Monocular MLF/Ours, Stereo MLF/Ours, LiDAR AVOD)

MAP$_{3D}$ at IoU=0.7 (moderate)
(bars: Monocular MLF/Ours, Stereo MLF/Ours, LiDAR AVOD)

**Analysis**

MAP$_{3D}$ at IoU=0.7 (moderate)
(DispNet MLF/Ours, PSMNet Frontal/Ours)

**Test results**

MAP$_{3D}$ at IoU=0.7
(easy Ours/AVOD, moderate Ours/AVOD, hard Ours/AVOD)

AVOD — Ours — Frontal

## Discussion, conclusion, and future work

- The historic performance gap between image- and LiDAR-based approaches may be more due to differences in processing rather than data quality.
- Pseudo-LiDAR largely improves image-based 3D detection, and may be a promising alternative (or complimentary) to LiDAR.
- **Future directions:** improve stereo depth for far-away objects and computational efficiency
- **Current progress:**
  - ➤ Novel stereo depth network: 45.3% → **50.4%**
  - ➤ Fuse stereo with 4-line LiDAR: 50.4% → **63.4%**
- **Code:** https://github.com/mileyan/pseudo_lidar

[1] Pyramid stereo matching network. In CVPR, 2018.
[2] Deep ordinal regression network for monocular depth estimation. In CVPR, 2018.
[3] Joint 3d proposal generation and object detection from view aggregation. In IROS, 2018.
[4] Multi-level fusion based 3d object detection from monocular images. In CVPR, 2018.
[5] A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In CVPR, 2016.