# Implicit Neural Representations: From Objects to 3D Scenes

Andreas Geiger

Autonomous Vision Group
University of Tübingen / MPI for Intelligent Systems Tübingen

June 19, 2020

University of Tübingen
MPI for Intelligent Systems
**Autonomous Vision Group**

# Collaborators

Songyou Peng

Michael Oechsle

Carolin Schmitt

Michael Niemeyer

Lars Mescheder

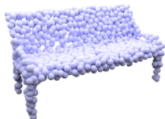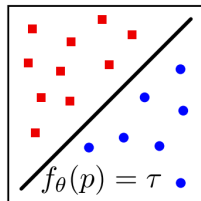Simon Donne
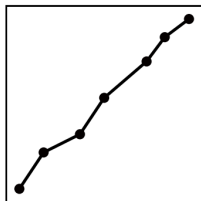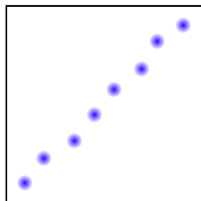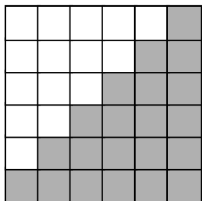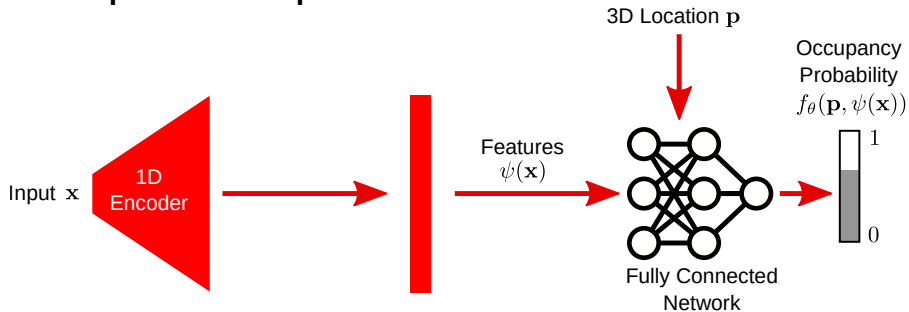
Gernot Riegler

Vladlen Koltun

Marc Pollefeys

Andreas Geiger

# 3D Representations



► Traditional Explicit Representations ⇒ **Discrete**

► Implicit Neural Representation ⇒ **Continuous**

Mescheder, Oechsle, Niemeyer, Nowozin and Geiger: Occupancy Networks: Learning 3D Reconstruction in Function Space. CVPR, 2019.

# Limitations

**Structure of implicit neural representations:**


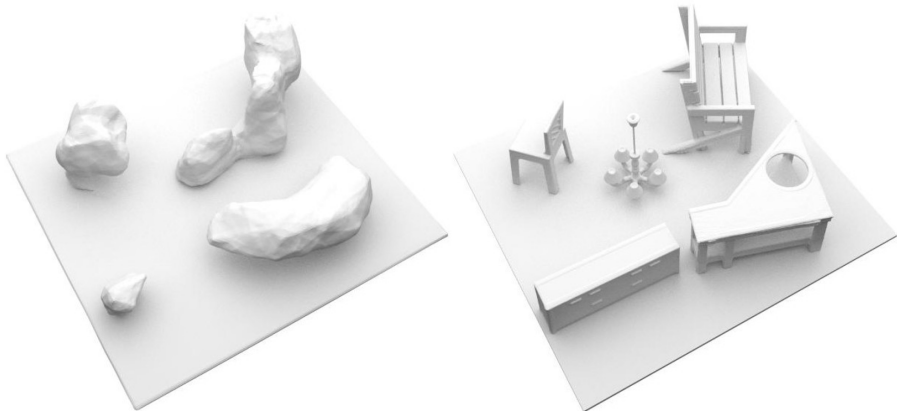
- Global latent code ⇒ no local information, overly smooth geometry
- Fully connected architecture ⇒ does not exploit translation equivariance

Mescheder, Oechsle, Niemeyer, Nowozin and Geiger: Occupancy Networks: Learning 3D Reconstruction in Function Space. CVPR, 2019.
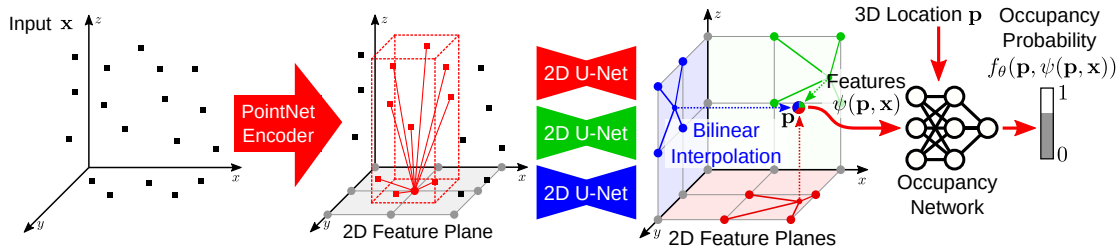
# Limitations

Implicit models work well for simple objects but **poorly on complex scenes:**

How to reconstruct large-scale 3D scenes with implicit neural representations?

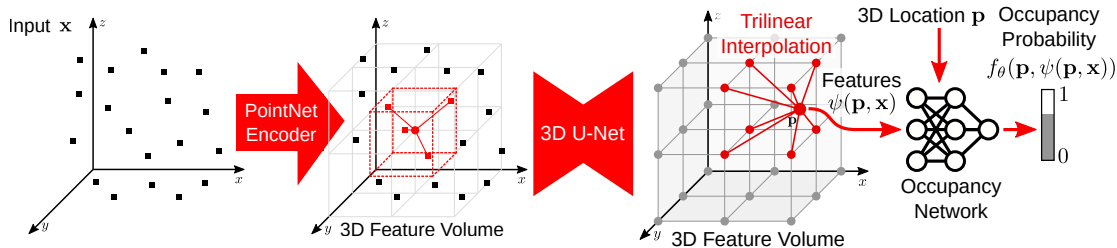**Convolutional Occupancy Networks**

# Convolutional Occupancy Networks



▶ **2D Plane Encoder:** Local PointNet processes input, project onto canonical plane

▶ **2D Plane Decoder:** Processed by U-Net, query features via bilinear interpolation

▶ **Occupancy Readout:** Shallow occupancy network $f_\theta(\cdot)$

Peng, Niemeyer, Mescheder, Pollefeys and Geiger: Convolutional Occupancy Networks. arXiv, 2020.

# Convolutional Occupancy Networks



- ▶ **3D Volume Encoder:** Local PointNet processes input, volumetric feature encoding
- ▶ **3D Volume Decoder:** Processed by 3D U-Net, query features via trilinear interp.
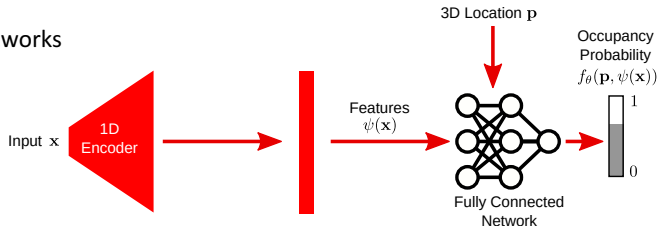- ▶ **Occupancy Readout:** Shallow occupancy network $f_\theta(\cdot)$

# Comparison

## Occupancy Networks



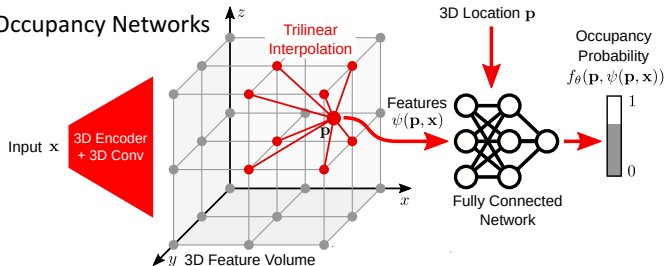## Convolutional Occupancy Networks



Peng, Niemeyer, Mescheder, Pollefeys and Geiger: Convolutional Occupancy Networks. arXiv, 2020.

Results

# Object-Level Reconstruction



Input    ONet    Ours    GT

# Training Speed

# Training Speed

Peng, Niemeyer, Mescheder, Pollefeys and Geiger: Convolutional Occupancy Networks. arXiv, 2020.

# Scene-Level Reconstruction



Input      ONet      SPSR      Ours      GT

► Trained and evaluated on synthetic rooms

Peng, Niemeyer, Mescheder, Pollefeys and Geiger: Convolutional Occupancy Networks. arXiv, 2020.

# Scene-Level Reconstruction



| Input | ONet | SPSR | Ours |

► Trained on synthetic rooms, evaluated on **ScanNet**

Peng, Niemeyer, Mescheder, Pollefeys and Geiger: Convolutional Occupancy Networks. arXiv, 2020.

# Large-Scale Reconstruction



**Results on Matterport3D**

► Fully convolutional model

► Trained on synthetic crops

► Sliding window evaluation

► Scales to any scene size

Room-Level Reconstruction

Peng, Niemeyer, Mescheder, Pollefeys and Geiger: Convolutional Occupancy Networks. arXiv, 2020.

# Key Insights

**Key Insights:**

- ► Convolutional models allow for scaling implicit models to larger scenes
- ► Convolutional models train faster than fully implicit models
- ► Convolutional models allow for incorporating local feature information
- ► For objects, the 3-plane model has the best accuracy/memory trade-off
- ► For scenes, the volumetric representation performs best
- ► Models transfer from synthetic to real scenes

How to capturing the visual appearance of objects?

**Conditional Surface Light Fields**

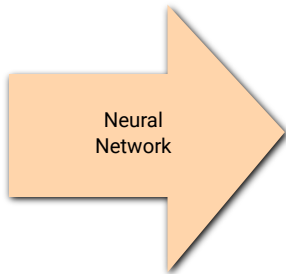# Problem Definition

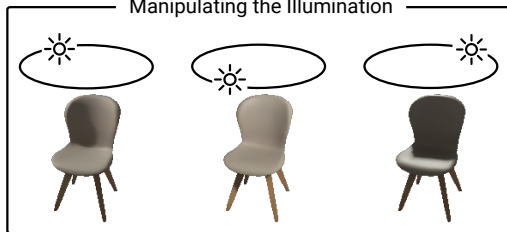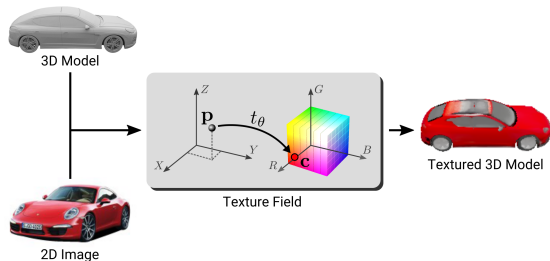# Existing Representation

**Texture Fields**

► 3D consistent

► Generalize across objects

► View-point independent

► Do not model lighting



3D Model

2D Image

Texture Field

Textured 3D Model

[Oechsle et al., ICCV 2019]

# Conditional Surface Light Field

**Rendering equation:**

$$L(\mathbf{p}, \mathbf{v}, \mathbf{l}, \mathbf{n}) = \int_\Omega \mathsf{svBRDF}(\mathbf{p}, \mathbf{r}, \mathbf{v}) \cdot \mathbf{l}(\mathbf{r}) \cdot (\mathbf{n}^T \mathbf{r}) \, d\mathbf{r}$$

**Conditional surface light field:**

$$L_{\mathsf{cSLF}}(\mathbf{p}, \mathbf{v}, \mathbf{l}) : \mathbb{R}^3 \times \mathbb{R}^3 \times \mathbb{R}^M \to \mathbb{R}^3$$



Oechsle, Niemeyer, Mescheder, Strauss and Geiger: Learning Implicit Surface Light Fields. arXiv, 2020.

# Overfitting to Single Objects

# Single-Image Appearance Prediction



Oechsle, Niemeyer, Mescheder, Strauss and Geiger: Learning Implicit Surface Light Fields. arXiv, 2020.

# Generative Model

Oechsle, Niemeyer, Mescheder, Strauss and Geiger: Learning Implicit Surface Light Fields. arXiv, 2020.

How to obtain training data with materials?

**Joint Estimation of Pose, Geometry and svBRDF**

# Joint Estimation of Pose, Geometry and svBRDF

**Goal: Dataset of 3D indoor scenes**
captured with high accuracy
from a handheld mobile sensor.

**Custom built sensor rig:**

► Custom IR depth sensor
   similar to Microsoft Kinect

► Active illumination + RGB camera
   for material estimation



Schmitt, Donné, Riegler, Koltun and Geiger: On Joint Estimation of Pose, Geometry and svBRDF From a Handheld Scanner. CVPR, 2020.          23

# Joint Estimation of Pose, Geometry and svBRDF



Schmitt, Donné, Riegler, Koltun and Geiger: On Joint Estimation of Pose, Geometry and svBRDF From a Handheld Scanner. CVPR, 2020.

# Materials $\longleftrightarrow$ Geometry

$\longrightarrow$   Accurate geometry reconstruction requires known appearance properties

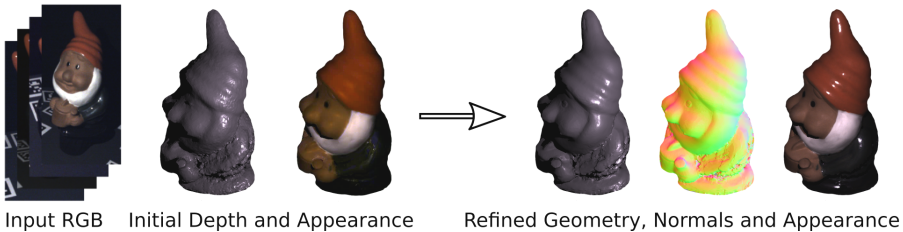$\longleftarrow$   Accurate appearance estimation requires very well known geometry

$\longleftrightarrow$   Joint estimation requires **only a rough initialization** for both
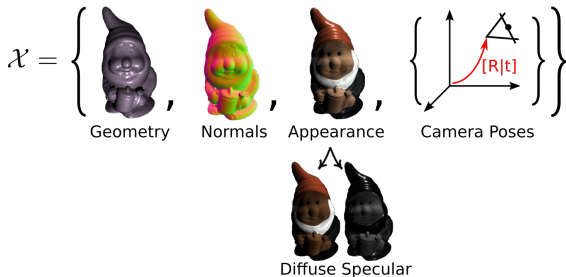


Input RGB   Initial Depth and Appearance      Refined Geometry, Normals and Appearance

# Joint Estimation of Pose, Geometry and svBRDF

**Contributions:**

- ► **Joint** formulation

- ► **Single objective function**
  minimized using off-the-shelf
  gradient-based solvers

- ► **Meaningful segmentation**
  differentiably part of the optimization

- ► **Accurate geometry**
  with very fine details



$$\mathcal{X} = \left\{ \, \text{Geometry}, \, \text{Normals}, \, \text{Appearance}, \, \left\{ \, \text{Camera Poses} \, \right\} \right\}$$

Geometry  Normals  Appearance  Camera Poses

Diffuse Specular

$$\mathcal{X}^* = \underset{\mathcal{X}}{\operatorname{argmin}} \quad \mathcal{L}(\mathcal{X})$$

Schmitt, Donné, Riegler, Koltun and Geiger: On Joint Estimation of Pose, Geometry and svBRDF From a Handheld Scanner. CVPR, 2020.

# Joint Estimation of Pose, Geometry and svBRDF

**Contributions:**

- ► **Joint** formulation
- ► **Single objective function**
  minimized using off-the-shelf
  gradient-based solvers
- ► **Meaningful segmentation**
  differentiably part of the optimization
- ► **Accurate geometry**
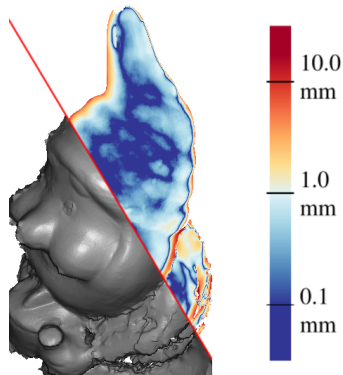  with very fine details



Reconstruction    Segmentation

Schmitt, Donné, Riegler, Koltun and Geiger: On Joint Estimation of Pose, Geometry and svBRDF From a Handheld Scanner. CVPR, 2020.

# Joint Estimation of Pose, Geometry and svBRDF

**Contributions:**

- ▶ **Joint** formulation
- ▶ **Single objective function**
  minimized using off-the-shelf
  gradient-based solvers
- ▶ **Meaningful segmentation**
  differentiably part of the optimization
- ▶ **Accurate geometry**
  with very fine details



Geometric error

Schmitt, Donné, Riegler, Koltun and Geiger: On Joint Estimation of Pose, Geometry and svBRDF From a Handheld Scanner. CVPR, 2020.

# Qualitative Results



Relighting          Novel Viewpoint

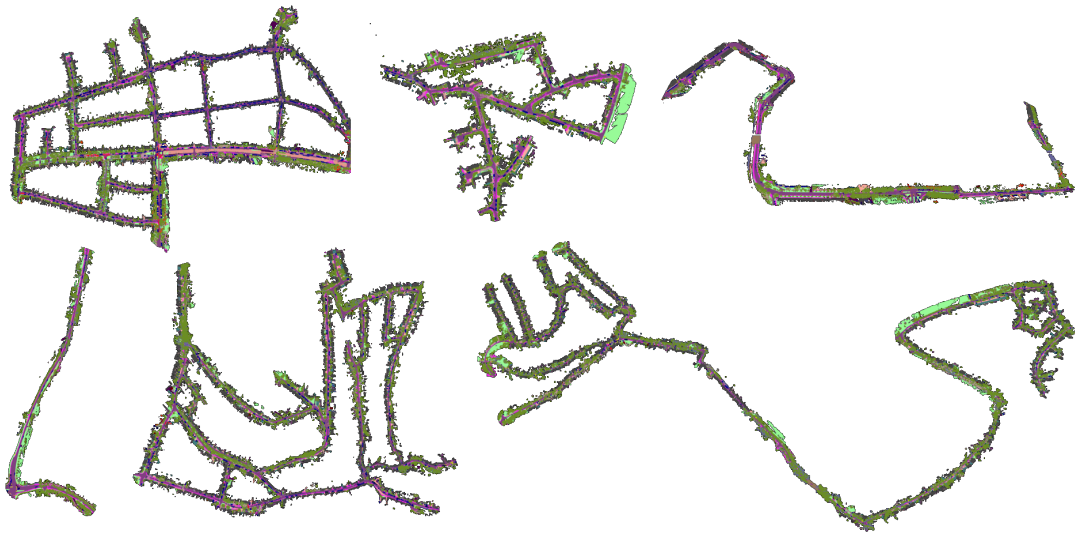**Conclusion:**

► Joint estimation helps

► This is only a first step

► Object-level reconstruction remains challenging with limited observations

► Scaling to larger scenes

► Scaling to scenes with external illumination

Schmitt, Donné, Riegler, Koltun and Geiger: On Joint Estimation of Pose, Geometry and svBRDF From a Handheld Scanner. CVPR, 2020.                26

How to obtain training data with semantic labels?

**KITTI-360**

# KITTI-360



Xie, Kiefel, Sun, Geiger: Semantic Instance Annotation of Street Scenes by 3D to 2D Label Transfer. CVPR, 2016
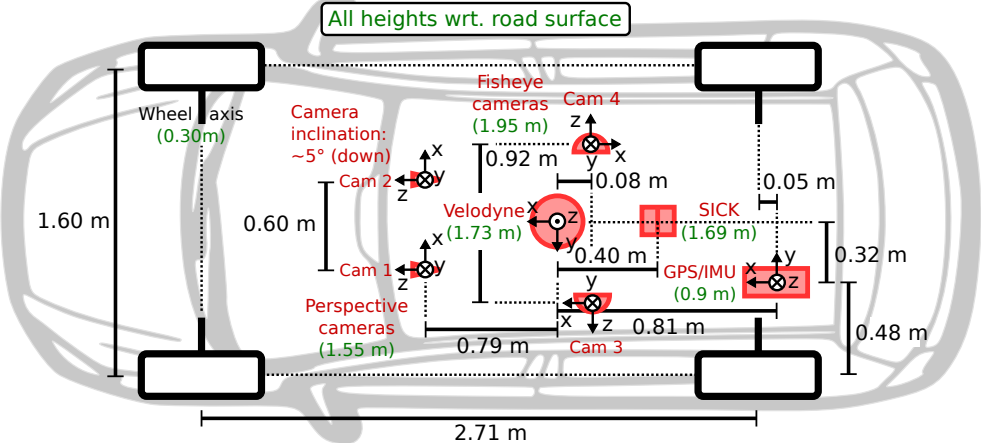
# KITTI-360

**Sensors:**

- ► Front-facing stereo camera
- ► 360° fisheye cameras
- ► Velodyne HDL 64 laser scanner
- ► SICK pushbroom laser scanner
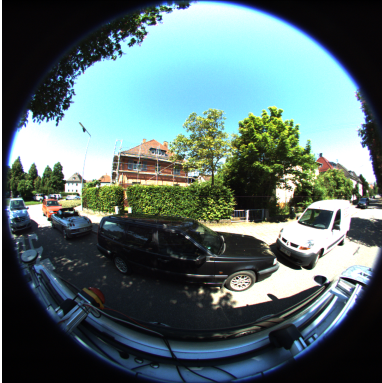- ► IMU/GPS localization system

**Features:**

- ► Driving distance: **73.7 km**    Frames: **4 × 83,000**
- ► All frames accurately **geolocalized** ($\Rightarrow$ OpenStreetMap)
- ► Semantic label definition consistent with Cityscapes, **19 classes** for evaluation
- ► Each instance assigned with a **consistent instance ID** across all frames

# Sensors



All heights wrt. road surface

Wheel axis (0.30 m)

Camera inclination: ~5° (down)

Fisheye cameras (1.95 m)
Cam 4

Cam 2

Velodyne (1.73 m)

SICK (1.69 m)

Cam 1

GPS/IMU (0.9 m)

Perspective cameras (1.55 m)

Cam 3

1.60 m
0.60 m
0.92 m
0.08 m
0.05 m
0.40 m
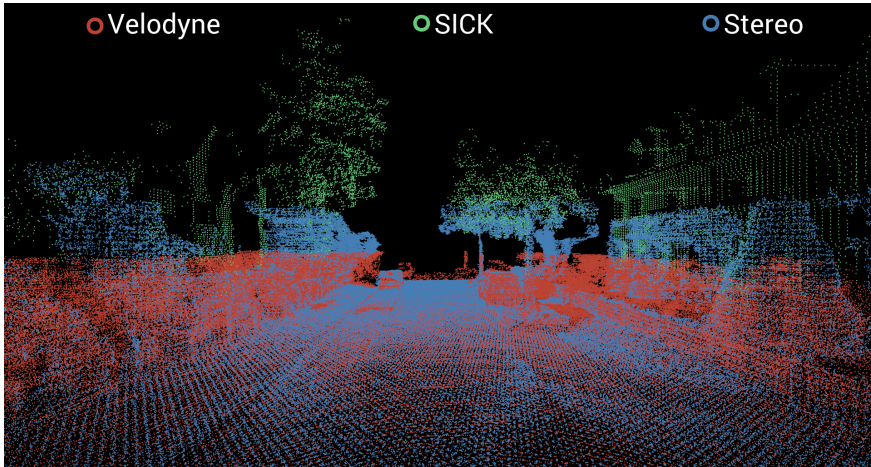0.32 m
0.81 m
0.48 m
0.79 m
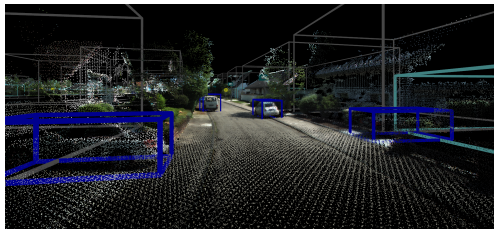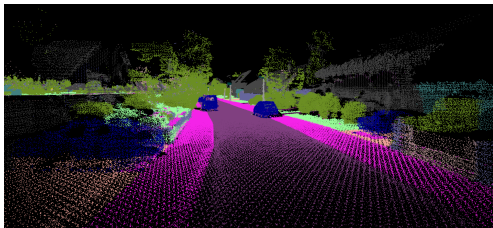2.71 m

# 360° 2D Sensors

# 360° 3D Sensors

# 3D Annotations


RGB


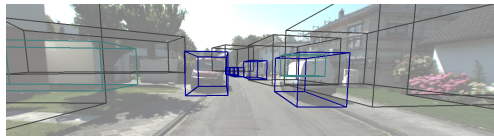Bounding Box


Semantic


Instance

33

# 2D Annotations



Semantic

Instance

Confidence

Bounding Box

Thank you!

http://autonomousvision.github.io