

# Geo-Supervised Visual Depth Prediction

Xiaohan Fei, Alex Wong, and Stefano Soatto

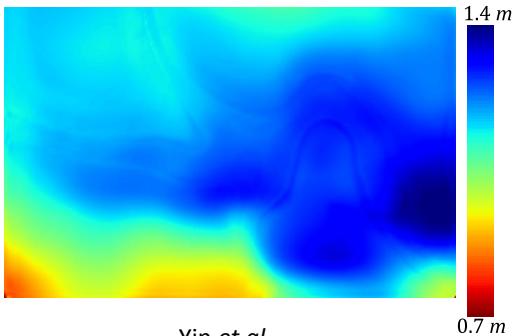
**UCLAVISIONLAB**

# Gravity Biases Shape

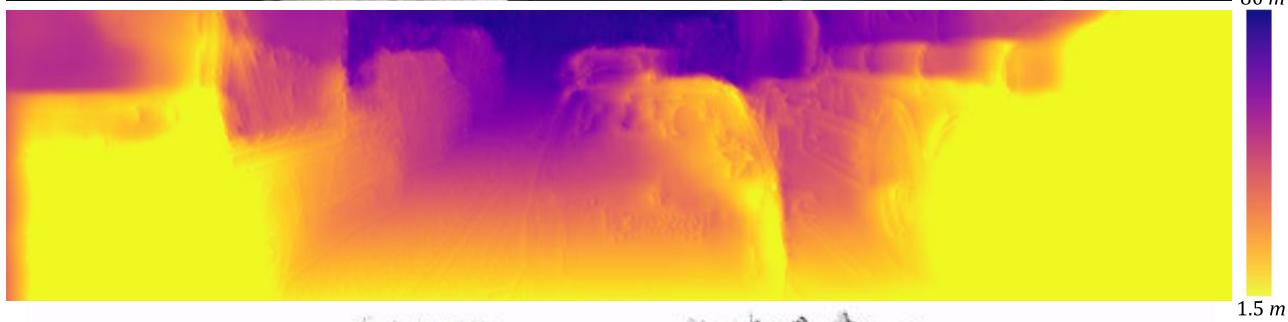


# Generic Regularization

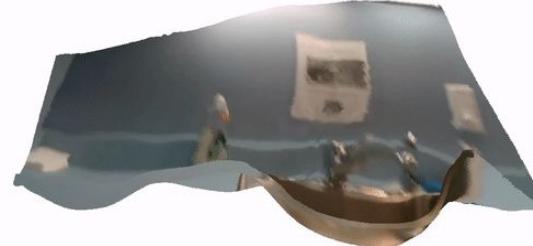




Yin *et al.*



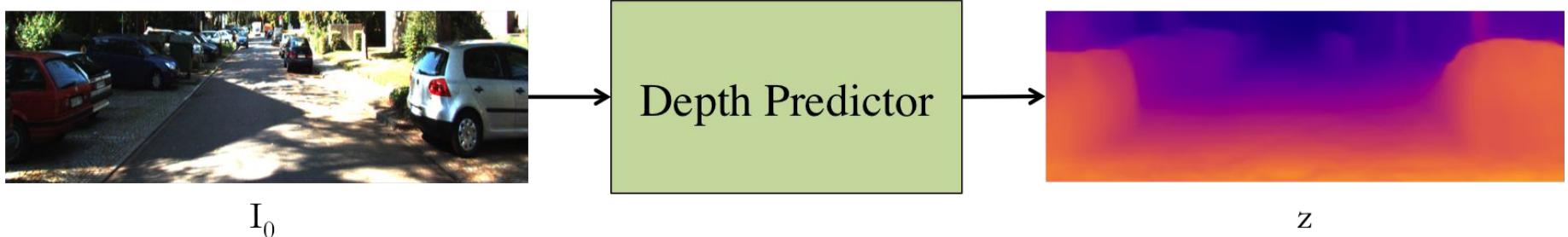
Godard *et al.*



[1] Yin *et al.*, GeoNet: Unsupervised learning of depth, optical flow and camera pose. CVPR, 2018.

[2] Godard *et al.*, Unsupervised Monocular Depth Estimation with Left-Right Consistency. CVPR, 2017.

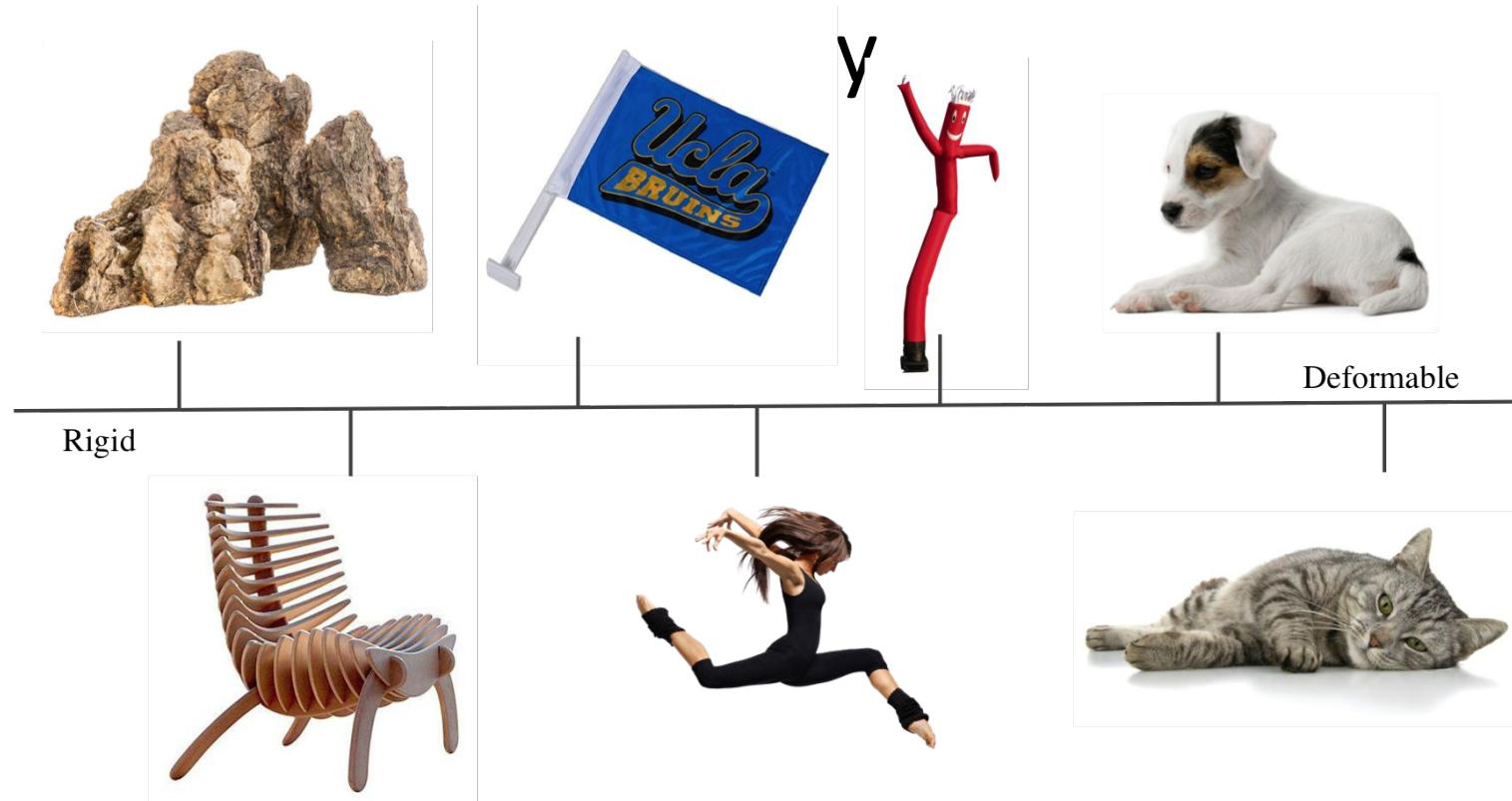
# Single Image Depth Prediction



Single image depth prediction is an ill-posed problem

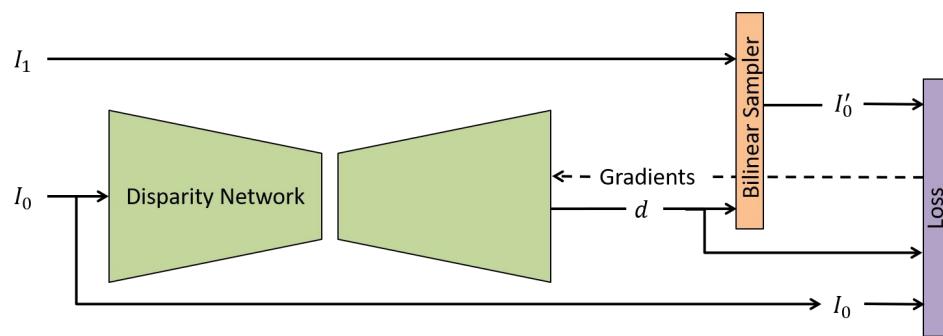
We **must** rely on a prior -- one that exploits the bias induced by gravity

# Not All Objects are Biased by



# Typical Architecture of Depth Prediction Systems

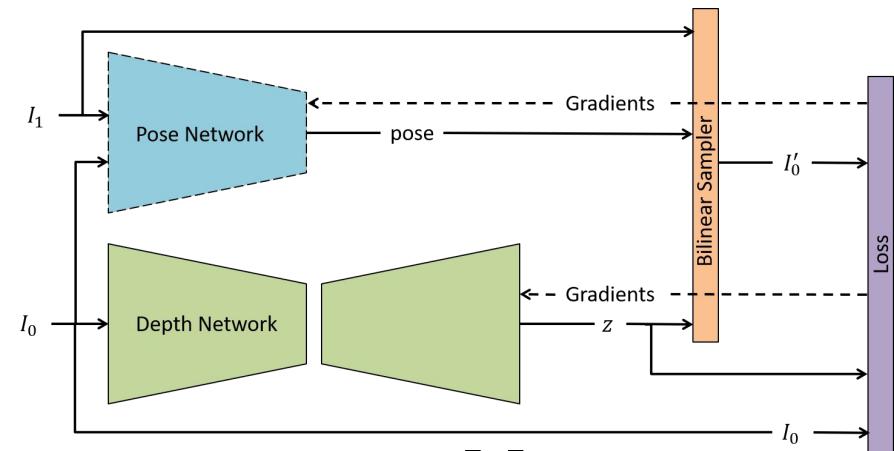
Stereo Training Pipeline [2]



$$I'_0(x, y) = I_1(x + d, y)$$

$$L = \sum_{x,y \in \Omega} |I_0(x, y) - I'_0(x, y)| + |\nabla d(x, y)|$$

Monocular Training Pipeline [3]



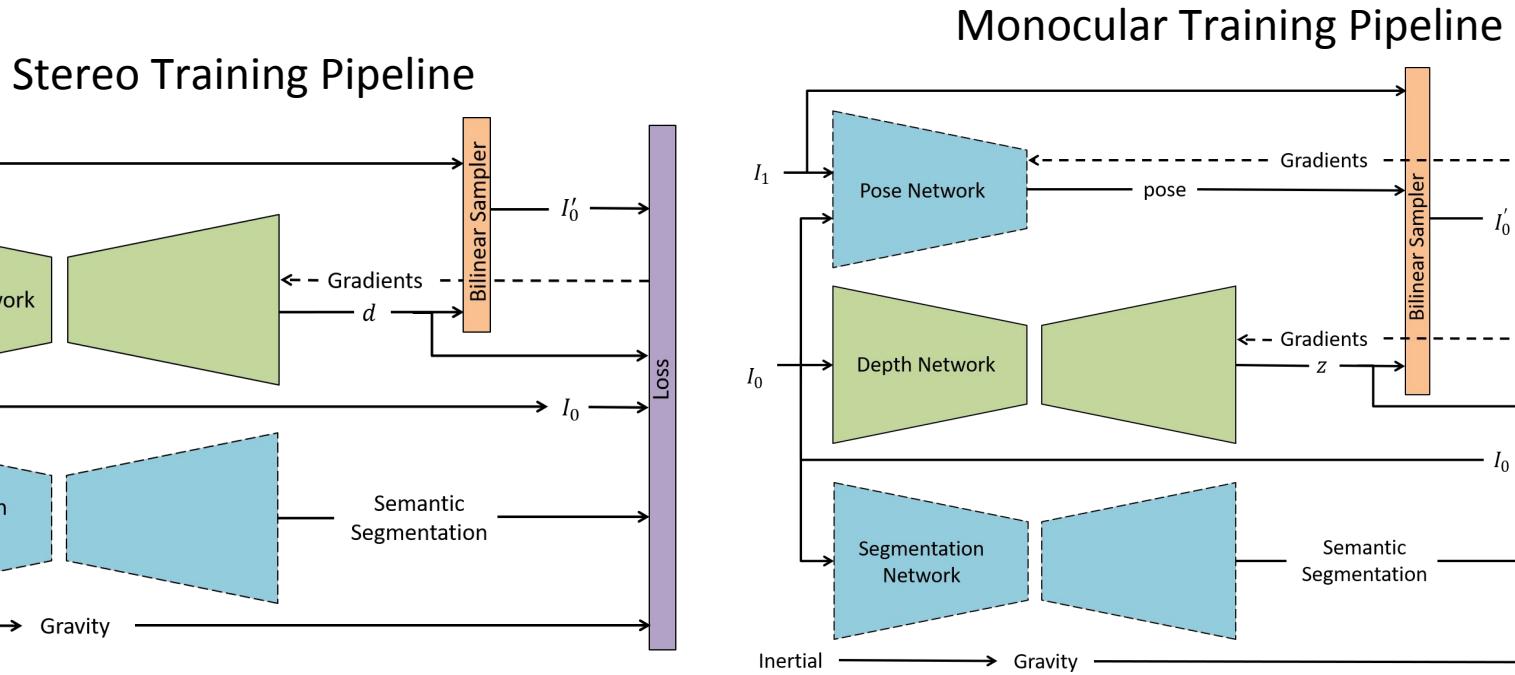
$$I'_0(x, y) = I_1(\pi g K^{-1} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} z)$$

$$L = \sum_{x,y \in \Omega} |I_0(x, y) - I'_0(x, y)| + |\nabla z(x, y)|$$

[2] Godard *et al.*, Unsupervised Monocular Depth Estimation with Left-Right Consistency. CVPR, 2017.

[3] Zhou *et al.*, Unsupervised learning of depth and ego-motion from vision. CVPR, 2017.

# Our Proposed System



Incorporating a Semantically Induced Geometric Loss (SIGL):

$$L^* = L + \underbrace{L_{HP} + L_{VP}}_{\text{SIGL}}$$

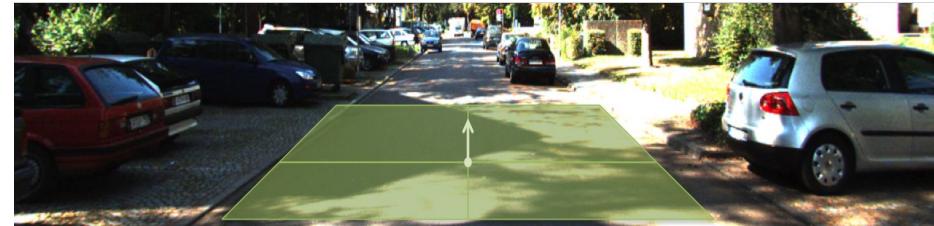
# SIGL: Horizontal Plane Loss

$$\begin{aligned} L_{HP}(\Omega_{HP}) &= \frac{1}{|\Omega_{HP}|} \|(\mathbf{I} - \frac{1}{|\Omega_{HP}|} \mathbf{1}\mathbf{1}^\top) \bar{\mathbf{X}}_\gamma\|_2^2 \\ &= \frac{1}{|\Omega_{HP}|} \sum_{i=1}^{|\Omega_{HP}|} ((\mathbf{X}_i - \mu)^\top \gamma)^2 \end{aligned}$$

$\Omega_{HP} \subset \mathbb{R}^2$  is a subset of the image plane whose associated semantic classes have **horizontal surfaces**

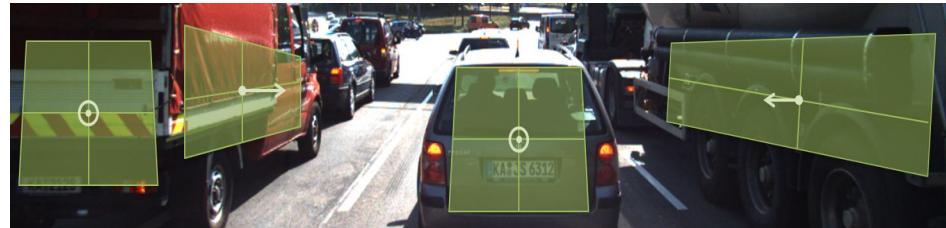
$\bar{\mathbf{X}}$  is the collection of 3D points which project to  $\Omega_{HP}$

$\gamma$  denotes gravity



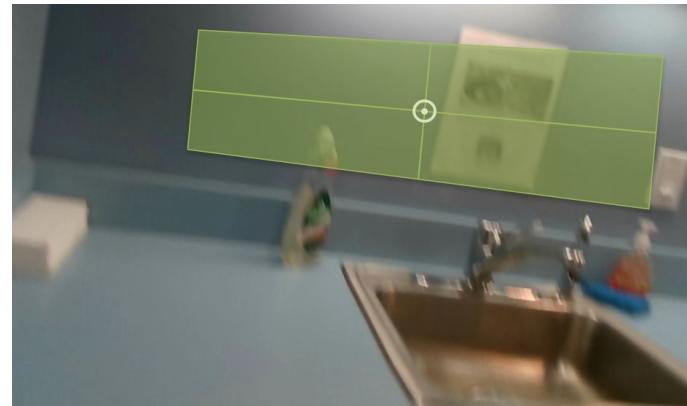
# SIGL: Vertical Plane Loss

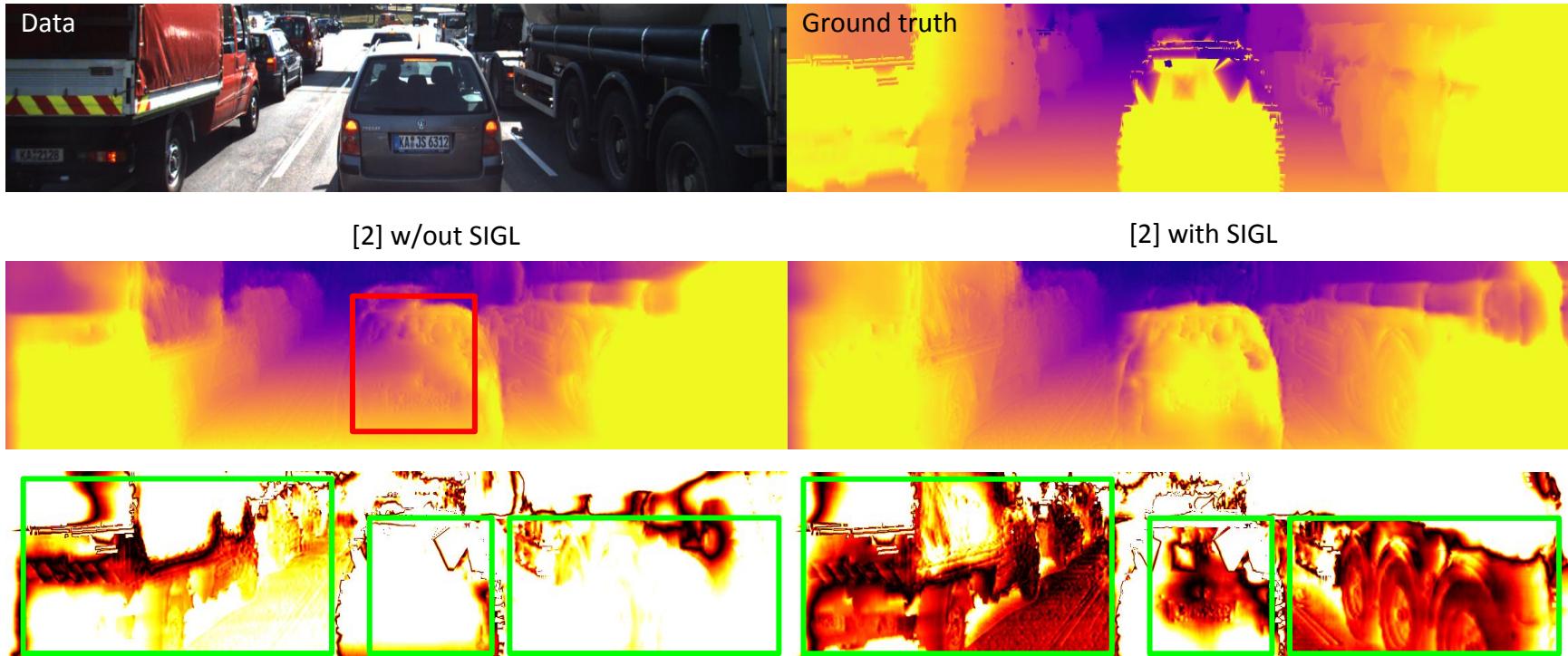
$$L_{VP}(\Omega_{VP}) = \min_{\substack{N \in \mathcal{N}(\gamma) \\ \|N\|=1}} \frac{1}{|\Omega_{VP}|} \|(\mathbf{I} - \frac{1}{|\Omega_{VP}|} \mathbf{1}\mathbf{1}^\top) \bar{\mathbf{X}} N\|_2^2$$



$\Omega_{VP} \subset \mathbb{R}^2$  is a subset of the image plane whose associated semantic classes have **vertical surfaces**

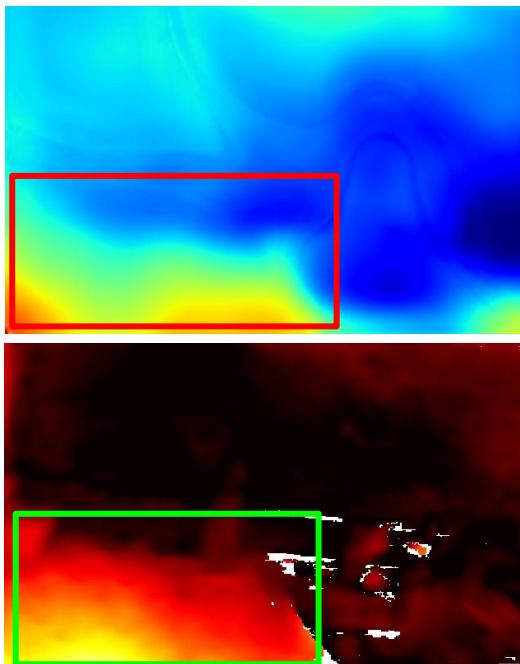
$\mathcal{N}(\gamma)$  is the null space of  $\gamma$



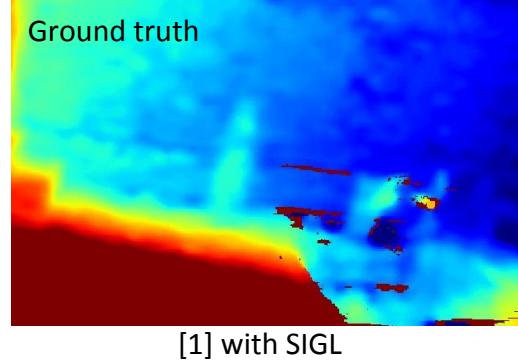
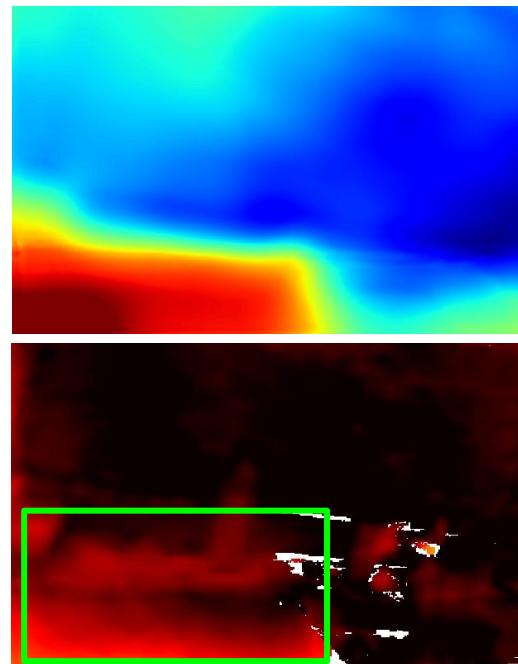


[Geiger *et al.*, Vision meets robotics: The KITTI dataset. IJRR, 2013.]

Depth



Error



[1] with SIGL

Our Visual-Inertial  
and Depth Dataset

Method	Data	Error metric				Accuracy ( $\bar{\delta} <$ )		
		AbsRel	SqRel	RMSE	RMSElog	1.25	1.25 <sup>2</sup>	1.25 <sup>3</sup>
<b>Stereo Training</b>								
Monodepth <i>VGG</i> [2] +SIGL	K	0.148	1.344	5.937	0.247	0.803	0.922	0.964
	K	<b>0.139</b>	<b>1.211</b>	<b>5.702</b>	<b>0.239</b>	<b>0.816</b>	<b>0.928</b>	<b>0.966</b>
Stereo-Temporal [4] +SIGL	K	0.144	1.391	5.869	0.241	0.803	0.928	0.969
	K	<b>0.137</b>	<b>1.061</b>	<b>5.692</b>	<b>0.239</b>	<b>0.805</b>	<b>0.928</b>	<b>0.969</b>
Monodepth <i>VGG</i> [2] +SIGL	CS+K	0.124	1.076	5.311	0.219	0.847	0.942	0.973
	CS+K	<b>0.114</b>	<b>0.885</b>	<b>4.877</b>	<b>0.203</b>	<b>0.858</b>	<b>0.950</b>	<b>0.978</b>
Monodepth <i>ResNet</i> [2] +SIGL	CS+K	0.114	0.898	4.935	0.206	0.861	0.949	0.976
	CS+K	<b>0.112</b>	<b>0.836</b>	<b>4.892</b>	<b>0.204</b>	<b>0.862</b>	<b>0.950</b>	<b>0.977</b>
<b>Monocular Training</b>								
GeoNet <i>ResNet</i> [1] +SIGL	K	0.155	1.296	5.857	0.233	0.793	0.931	0.973
	K	<b>0.142</b>	<b>1.124</b>	<b>5.611</b>	<b>0.223</b>	<b>0.813</b>	<b>0.938</b>	<b>0.975</b>
DDVO [5] +SIGL	K	0.151	1.257	5.583	0.228	<b>0.810</b>	0.936	0.974
	K	<b>0.146</b>	<b>1.068</b>	<b>5.538</b>	<b>0.224</b>	<b>0.809</b>	<b>0.938</b>	<b>0.975</b>
GeoNet <i>ResNet</i> [1] +SIGL	CS+K	0.153	1.328	5.737	0.232	0.802	0.934	0.972
	CS+K	<b>0.147</b>	<b>1.076</b>	<b>5.468</b>	<b>0.222</b>	<b>0.806</b>	<b>0.938</b>	<b>0.976</b>
DDVO [5] +SIGL	CS+K	0.148	1.187	5.496	0.226	0.812	0.938	0.975
	CS+K	<b>0.142</b>	<b>1.094</b>	<b>5.409</b>	<b>0.219</b>	<b>0.821</b>	<b>0.941</b>	<b>0.976</b>

CS = Cityscapes  
K = KITTI

[1] Yin *et al.*, GeoNet: Unsupervised learning of depth, optical flow and camera pose. CVPR, 2018.

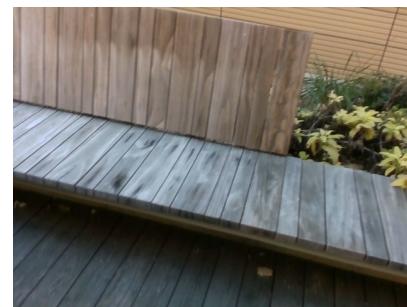
[2] Godard *et al.*, Unsupervised Monocular Depth Estimation with Left-Right Consistency. CVPR, 2017.

[4] Zhan *et al.*, Unsupervised learning of monocular depth estimation and visual odometry with deep feature reconstruction. CVPR, 2017.

[5] Wang *et al.*, Learning depth from monocular video using direct methods. CVPR, 2018.

Method	Data	Error metric				Accuracy ( $\bar{\delta} <$ )		
		AbsRel	SqRel	RMSE	RMSElog	1.25	1.25 <sup>2</sup>	1.25 <sup>3</sup>
<b>Stereo Training</b>								
Monodepth <i>VGG</i> [2] +SIGL	K	0.148	1.344	5.937	0.247	0.803	0.922	0.964
	K	<b>0.139</b>	<b>1.211</b>	<b>5.702</b>	<b>0.239</b>	<b>0.816</b>	<b>0.928</b>	<b>0.966</b>
Stereo-Temporal [4] +SIGL	K	0.144	1.391	5.869	0.241	0.803	0.928	0.969
	K	<b>0.137</b>	<b>1.061</b>	<b>5.692</b>	<b>0.239</b>	<b>0.805</b>	<b>0.928</b>	<b>0.969</b>
Monodepth <i>VGG</i> [2] +SIGL	CS+K	0.124	1.076	5.311	0.219	0.847	0.942	0.973
	CS+K	<b>0.114</b>	<b>0.885</b>	<b>4.877</b>	<b>0.203</b>	<b>0.858</b>	<b>0.950</b>	<b>0.978</b>
Monodepth <i>ResNet</i> [2] +SIGL	CS+K	0.114	0.898	4.935	0.206	0.861	0.949	0.976
	CS+K	<b>0.112</b>	<b>0.836</b>	<b>4.892</b>	<b>0.204</b>	<b>0.862</b>	<b>0.950</b>	<b>0.977</b>
<b>Monocular Training</b>								
GeoNet <i>ResNet</i> [1] +SIGL	K	0.155	1.296	5.857	0.233	0.793	0.931	0.973
	K	<b>0.142</b>	<b>1.124</b>	<b>5.611</b>	<b>0.223</b>	<b>0.813</b>	<b>0.938</b>	<b>0.975</b>
DDVO [5] +SIGL	K	0.151	1.257	5.583	0.228	<b>0.810</b>	0.936	0.974
	K	<b>0.146</b>	<b>1.068</b>	<b>5.538</b>	<b>0.224</b>	<b>0.809</b>	<b>0.938</b>	<b>0.975</b>
GeoNet <i>ResNet</i> [1] +SIGL	CS+K	0.153	1.328	5.737	0.232	0.802	0.934	0.972
	CS+K	<b>0.147</b>	<b>1.076</b>	<b>5.468</b>	<b>0.222</b>	<b>0.806</b>	<b>0.938</b>	<b>0.976</b>
DDVO [5] +SIGL	CS+K	0.148	1.187	5.496	0.226	0.812	0.938	0.975
	CS+K	<b>0.142</b>	<b>1.094</b>	<b>5.409</b>	<b>0.219</b>	<b>0.821</b>	<b>0.941</b>	<b>0.976</b>

CS = Cityscapes  
K = KITTI



# Limitations and Further Extensions

Exploit inertials at test time

Automatic selection of classes that are gravity-biased

Dynamic objects

# Limitations and Further Extensions

Exploit inertials at test time

Automatic selection of classes that are gravity-biased

Dynamic objects

# Limitations and Further Extensions

Exploit inertials at test time

Automatic selection of classes that are gravity-biased

Dynamic objects

Project site: [shorturl.at/cgiAS](http://shorturl.at/cgiAS)

