



Geo-Supervised Visual Depth Prediction

Xiaohan Fei
 <feixh@cs.ucla.edu>

Alex Wong
 <alexw@cs.ucla.edu>

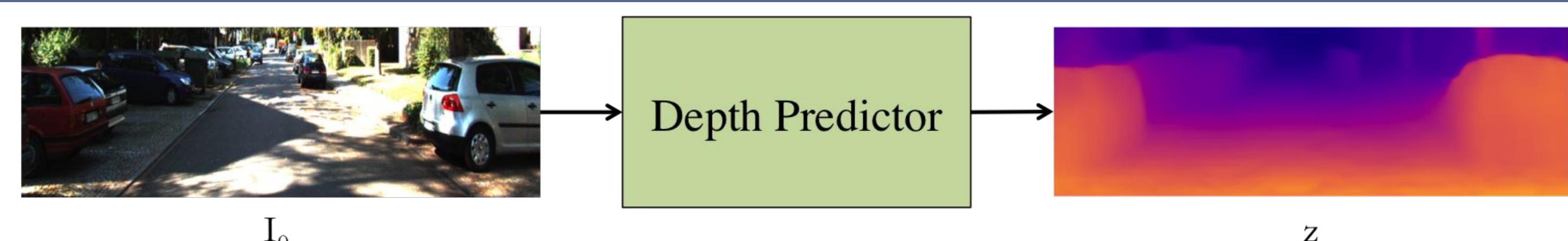
Stefano Soatto
 <soatto@cs.ucla.edu>

UCLAVISIONLAB



Motivation

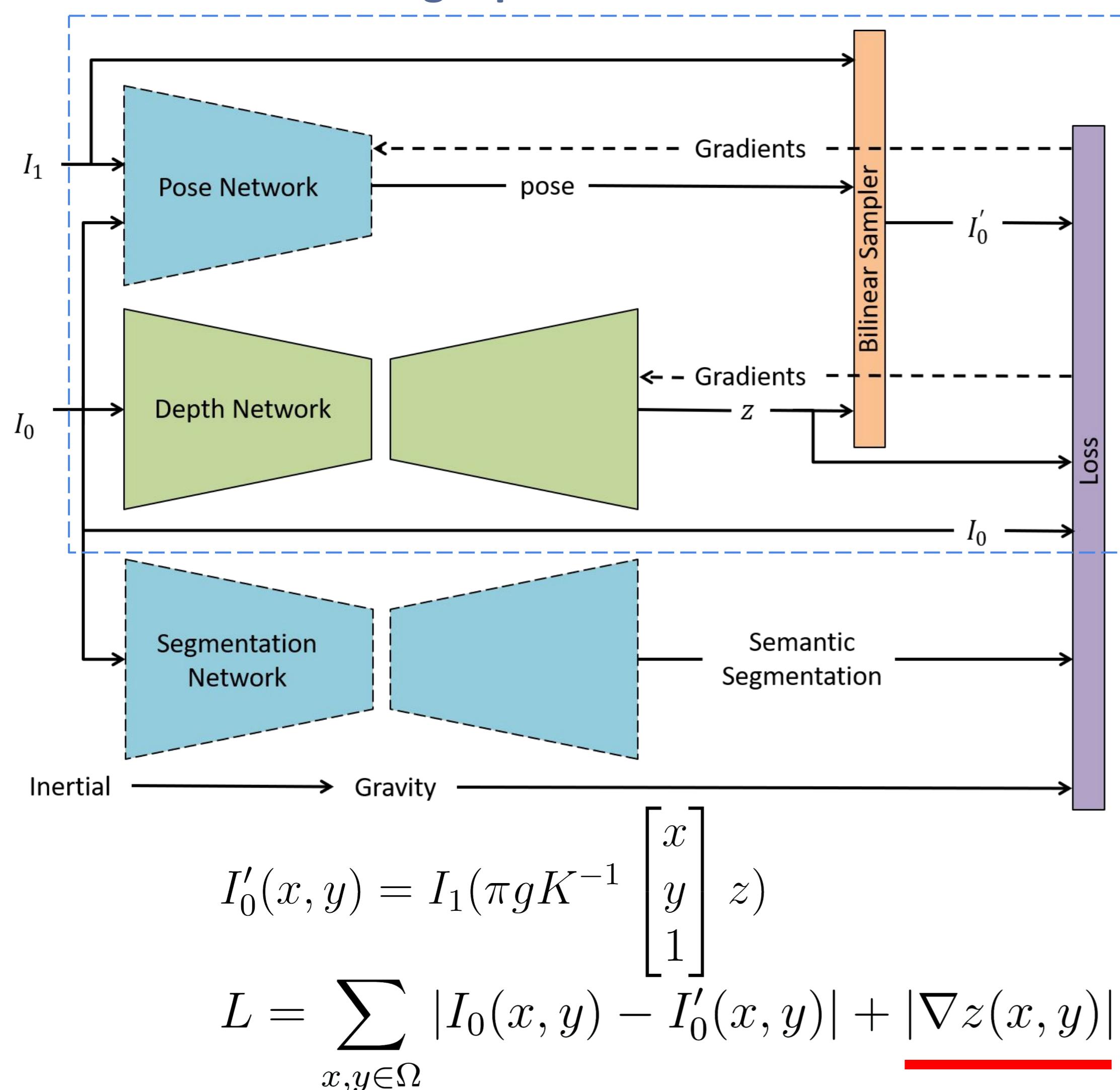
- ❖ Agents can benefit from understanding the shape of objects.
- ❖ The shape of objects surrounding us is biased by gravity.
- ❖ Gravity provides a persistent global orientation reference and can be easily inferred from inertial sensors.
- ❖ Many devices today have both visual and inertial sensors.
- ❖ **Goal:** to leverage gravity to impose priors on shapes to improve visual depth prediction.



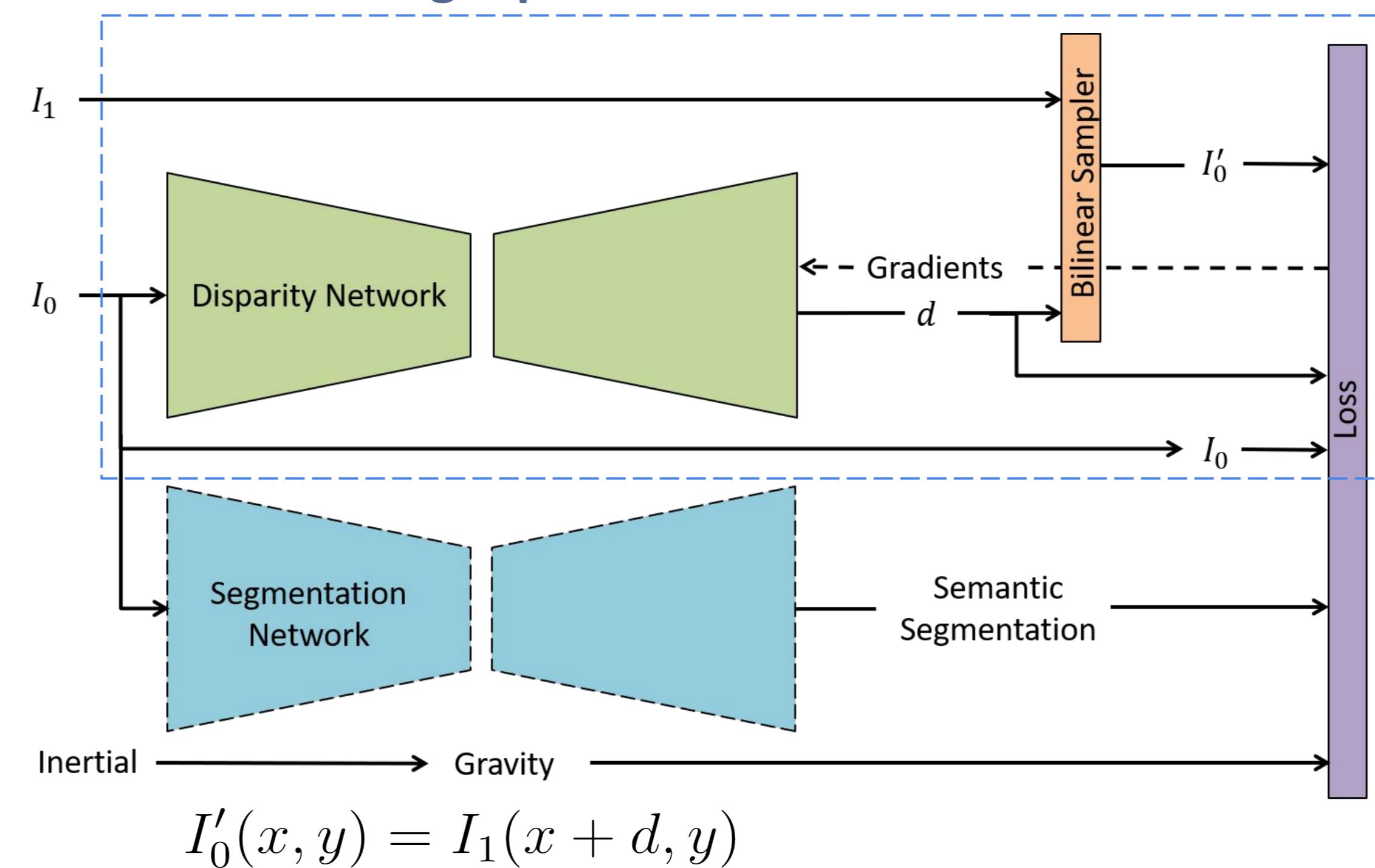
- ❖ Inferring shape from images is ill-posed so prior/regularization is needed.
- ❖ Not every object is biased by gravity so imposing the bias induced by gravity has to be done selectively, in a way that is object dependent.

Proposed System

Monocular Training Pipeline



Stereo Training Pipeline



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

SIGL = Semantically Informed Geometric Loss

Typical architecture

- ❖ Encoder-Decoder structure with skip-connections.
- ❖ Learn to predict depth by minimizing a view synthesis loss (data term) and some generic regularizers, such as piecewise smoothness.
- ❖ Two popular training diagrams:
 - Monocular videos
 - Stereo pairs
- ❖ At inference time, only one RGB image is needed.

Proposed system

- ❖ At training time, apply *category-specific* shape priors *selectively*, which requires:
 - 1) a semantic segmentation module to provide per pixel class labels
 - 2) a Visual-Inertial Odometry (VIO) system to provide reliable gravity estimation
- ❖ At inference time, still only one RGB image is needed.

Semantically Informed Geometric Loss (SIGL)

Vertical Plane Loss

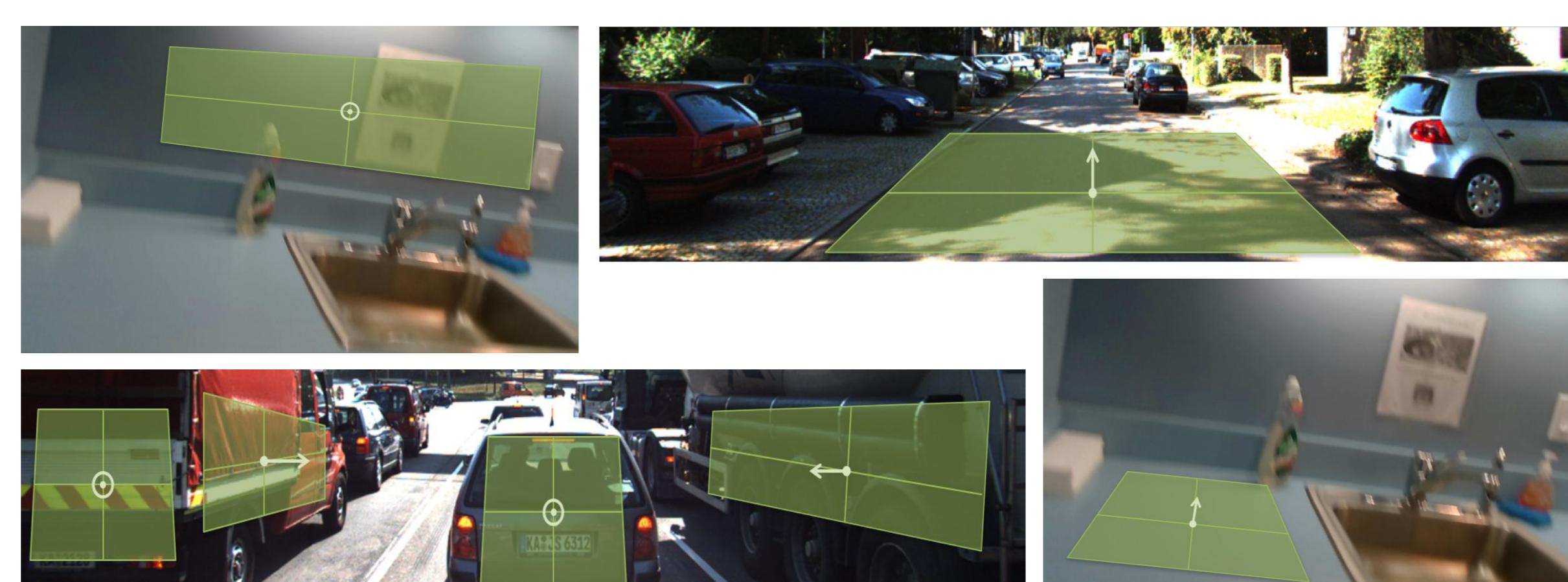
$$L_{VP}(\Omega_{VP}) = \min_{N \in \mathcal{N}(\gamma) \atop \|N\|=1} \frac{1}{|\Omega_{VP}|} \|(\mathbf{I} - \frac{1}{|\Omega_{VP}|} \mathbf{1}\mathbf{1}^\top) \bar{\mathbf{X}} N\|^2$$

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

$\bar{\mathbf{X}} \in \mathbb{R}^{|\Omega_{VP}| \times 3}$: collection of 3D points which project to Ω_{VP}

$\gamma \in \mathbb{R}^3$: gravity

$\mathcal{N}(\gamma)$: Null space of γ



Horizontal Plane Loss

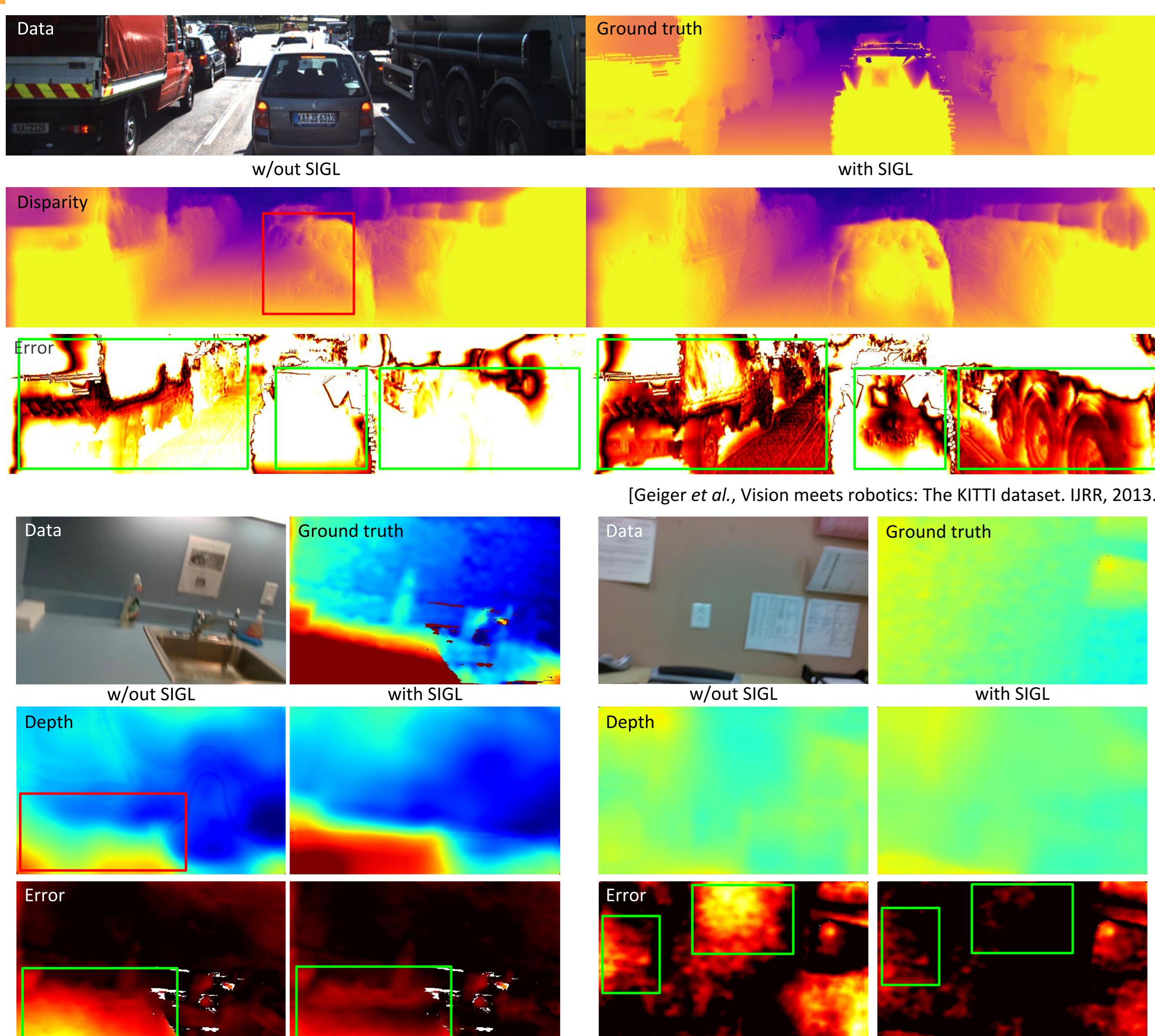
$$L_{HP}(\Omega_{HP}) = \frac{1}{|\Omega_{HP}|} \|(\mathbf{I} - \frac{1}{|\Omega_{HP}|} \mathbf{1}\mathbf{1}^\top) \bar{\mathbf{X}} \gamma\|^2$$

$$= \frac{1}{|\Omega_{HP}|} \sum_{i=1}^{|\Omega_{HP}|} ((\bar{\mathbf{X}}_i - \mu)^\top \gamma)^2$$

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

$\bar{\mathbf{X}} \in \mathbb{R}^{|\Omega_{HP}| \times 3}$: collection of 3D points which project to Ω_{HP}

Results



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

K=KITTI, CS=Cityscapes

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

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

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

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

Our Visual Inertial and Depth Dataset

- ❖ Extended version of the VISMA dataset used in paper
- ❖ RGB-D of VGA size @ 30 Hz
- ❖ Accel & Gyro @ 400 Hz
- ❖ Non-trivial 6 DoF motion
- ❖ > 40,000 frames, suitable for learning-based visual-inertial sensor fusion

Acknowledgement

This work was supported by ONR N00014-17-1-2072 and ARO W911NF-17-0-304.

