

Connecting the Dots: Learning Representations for Active Monocular Depth Estimation

Introduction

Passive vs. Active Stereo? Active Stereo can reconstruct textureless regions by projecting into the scene.

Learning representations?

We find that a conv. network is estimating disparities effectively, despite no absolute location information is encoded.

Where does the supervision come from? Self-supervision. We use robust photometric and geometric losses.

That's all?

We further utilize ambient edge information that can be disentangled from the input.

Okay, you got me. Where can I get it?

https://github.com/autonomousvision/connecting_the_dots



Gernot Riegler^{1,*}, Yiyi Liao^{2,*}, Simon Donne², Vladlen Koltun¹, Andreas Geiger² ¹Intel Intelligent Systems Lab, ²Autonomous Vision Group, MPI-IS / University of Tübingen



		Qua	mutati	ve Ke	:Suits			
		0	(0.5)	o(1)	0(2)	o(5)	
	$\begin{array}{cccc} \mathbf{I} & \rightarrow & \mathcal{I} \\ \mathbf{I} & \rightarrow & \mathcal{I} \end{array}$	P_{row}	6.22 8.19 11.83 58.12	3.00 4.35 5.08 28.81	1. 2.4 2.4 8.0	53 40 46 02	0.85 1.07 1.11 2.11	
	$\mathbf{A} ightarrow \mathbf{A}$	$\mathbf{A} \rightarrow \mathcal{D}$ 90.71 81.41 63.53 32.00 Architectural Choices						
	`		<i>o</i> (0.5) o	(1)	o(2)	<i>o</i> (5)	
	Supervis \mathcal{L}_P $\mathcal{L}_P + \mathcal{L}$ $\mathcal{L}_P + \mathcal{L}$	Supervised \mathcal{L}_P $\mathcal{L}_P + \mathcal{L}_D$ $\mathcal{L}_P + \mathcal{L}_D + \mathcal{L}_G$		2 3 2 6 7 4 7 3	.00 .00 .23 <i>:</i> .88	1.63 4.10 2.56 2.57	0.85 2.72 1.52 1.63	
		Influcer	nce of l	_OSS	Func	tion		
		o(0.5)	o(1)	0	2) <i>c</i>	p(5)	$o_u(1)$	$o_u(5)$
Block FastM Hyperl	Matching RF [1] Depth [2]	7.84 12.07 15.01	7.20 8.36 12.63	7.0 6.7 11.8)6 6 71 5 83 11	5.83 5.14 1.49	4.44 5.25 7.39	4.23 3.57 6.73
Ours		6.77 Comp	3.88 arison	2.5 to ba	aselin	. 63 es	1.75	0.70
		•						
		acc comp h. mean						
	Block N FastMI Hyper[Matching RF [1] Depth [2]	551.0 12.6 8.7	82 90 59	3.883 6.971 5.263		7.712 8.999 6.575	

Evaluation on dataset of [1]











Method Details

Photometric Loss - Minimize smooth census loss [3] to maximize photometric consistency:

$$\mathcal{L}_{P}^{i}(\hat{\mathbf{P}}_{i},\mathbf{P},\mathbf{D}_{i}) = \sum_{x,y} \|\hat{\mathbf{p}}_{i}(x,y) - \mathbf{p}(x - D_{i}(x,y),y)\|_{C}.$$
 (1)

Disparity Loss - Encourage disparity discontinuities at edges of the ambient image:

$$\mathcal{L}_D^i(\mathbf{D}'_i, \mathbf{E}_i) = \sum_{x, y} -\log\left(p(D'_i(x, y), E_i(x, y))\right) \,. \tag{2}$$

Edge Loss - Avoid trivial solution to disparity loss:

$$\mathcal{C}_{E}^{k}(\mathbf{E}_{k},\mathbf{A}_{k}) = -\sum_{x,y} A_{k}'(x,y) \log E_{k}(x,y) + w \left(1 - A_{k}'(x,y)\right) \log(1 - E_{k}(x,y)).$$
(3)

Geometric Loss - Utilize multiple views of the same static scene:

$$\mathcal{L}_{G}^{ij}(\mathbf{D}_{i}, \mathbf{X}_{j}') = \sum_{\mathbf{x} \in \mathbf{X}_{j}'} \min\left(\left|\mathbf{x}_{z} - bfD_{i}^{-1}(\mathbf{K}\,\mathbf{x})\right|, \tau\right) \,. \tag{4}$$

- I... Input image **D**... Predicted disparity map $\hat{\mathbf{P}} \dots$ LCN input image **A** . . . Ambient image
 - ${f P}\ldots {f LCN}$ reference pattern
 - ${f E}\ldots$ Predicted edge map
 - $\mathbf{D'}\ldots$ Disparity magnitude $=|
 abla \mathbf{D}|$
 - ${f X}\ldots$ Point cloud from predicted disparity map



Conclusion

- It is possible to train a network for disparity estimation in a self-supervised fashion
- Obtains similar results to training on ground-truth disparities
- Modeling joint distribution of depth and image edges improves depth boundaries
- Constraints on discontinuities from multiple views improves performance

	References	
[1]	Q. Chen and V. Koltun. "Fast MRF Optimization with Application to Depth Reconstruction". In: CVPR. 2014.	
[2]	S. R. Fanello et al. "HyperDepth: Learning Depth from Structured Light without Matching". In: CVPR. 2016.	
[3]	D. Hafner et al. "Why Is the Census Transform Good for Robust Optic Flow Computation?" In: SSVM. 2013.	