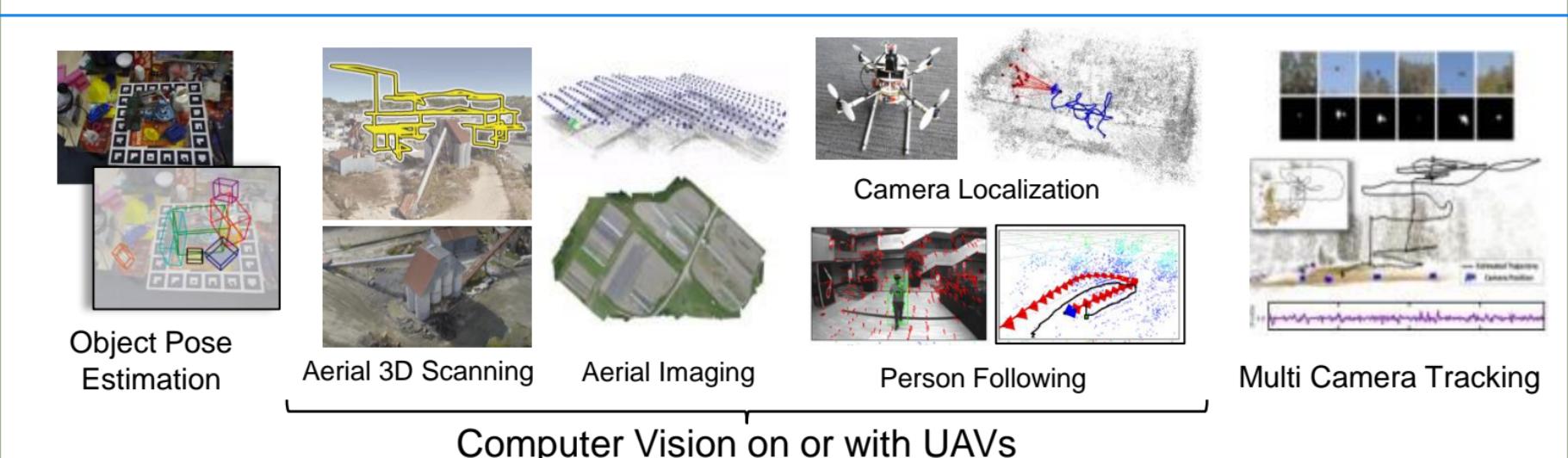
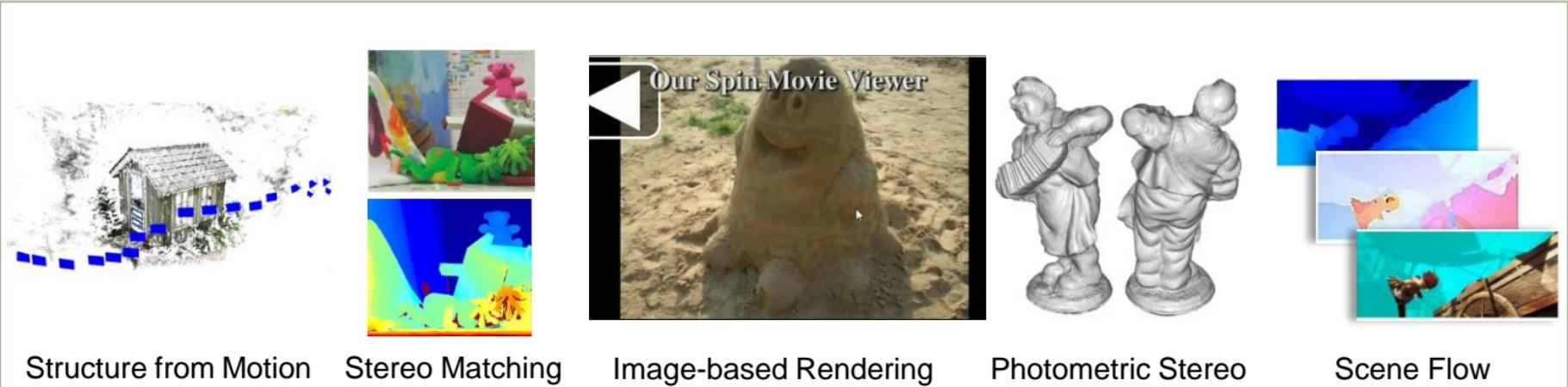


Recent Progress in Stereo Matching, Scene Flow and Object Pose Estimation

Sudipta Sinha

Microsoft Research

Talk at Georgia Tech, August 24, 2018





Structure from Motion



Stereo Matching



Image-based Rendering



Photometric Stereo



Scene Flow



**Object Pose
Estimation**



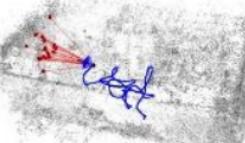
Aerial 3D Scanning



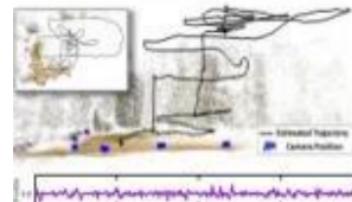
Aerial Imaging



Camera Localization



Person Following



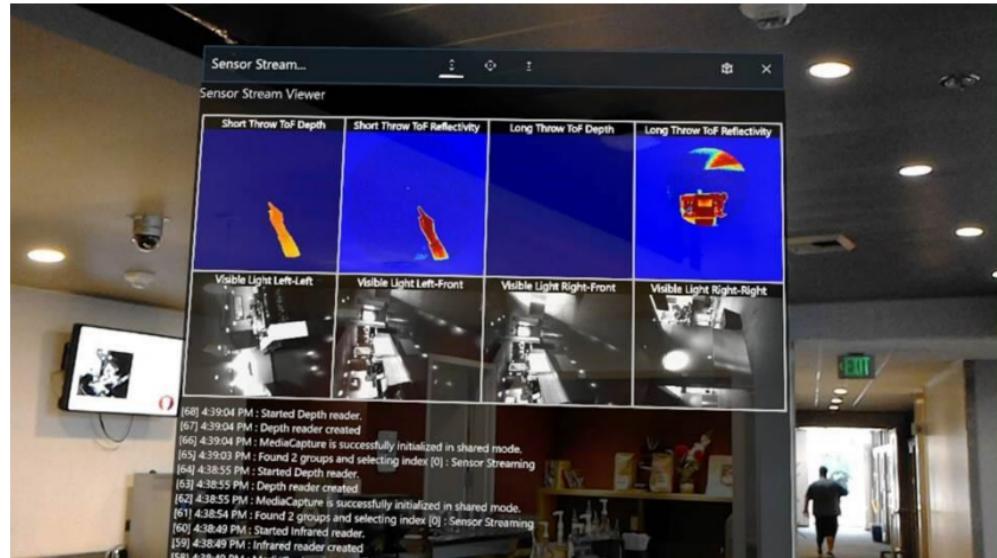
Multi Camera Tracking

Computer Vision on or with UAVs

HoloLens Research Mode

<https://docs.microsoft.com/en-us/windows/mixed-reality/cvpr-2018>

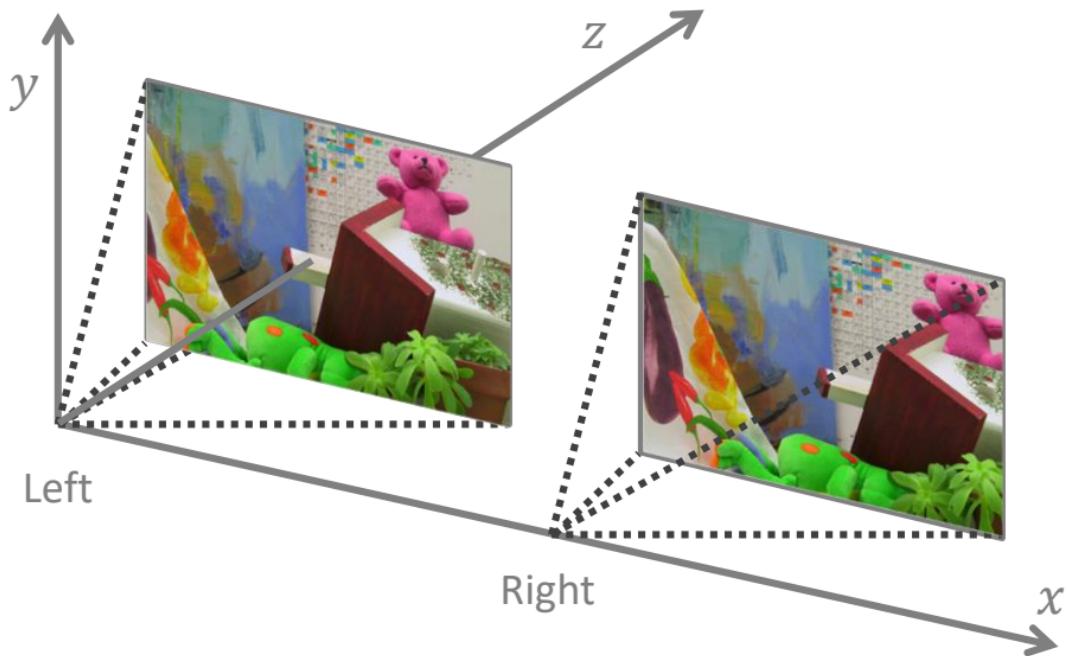
- On-device API
- Gives access to
 - audio
 - 6-dof pose
 - surface mesh
 - raw video streams
 - 4x greyscale cameras
 - Color camera
 - ToF Depth & IR Reflectivity (near, far)



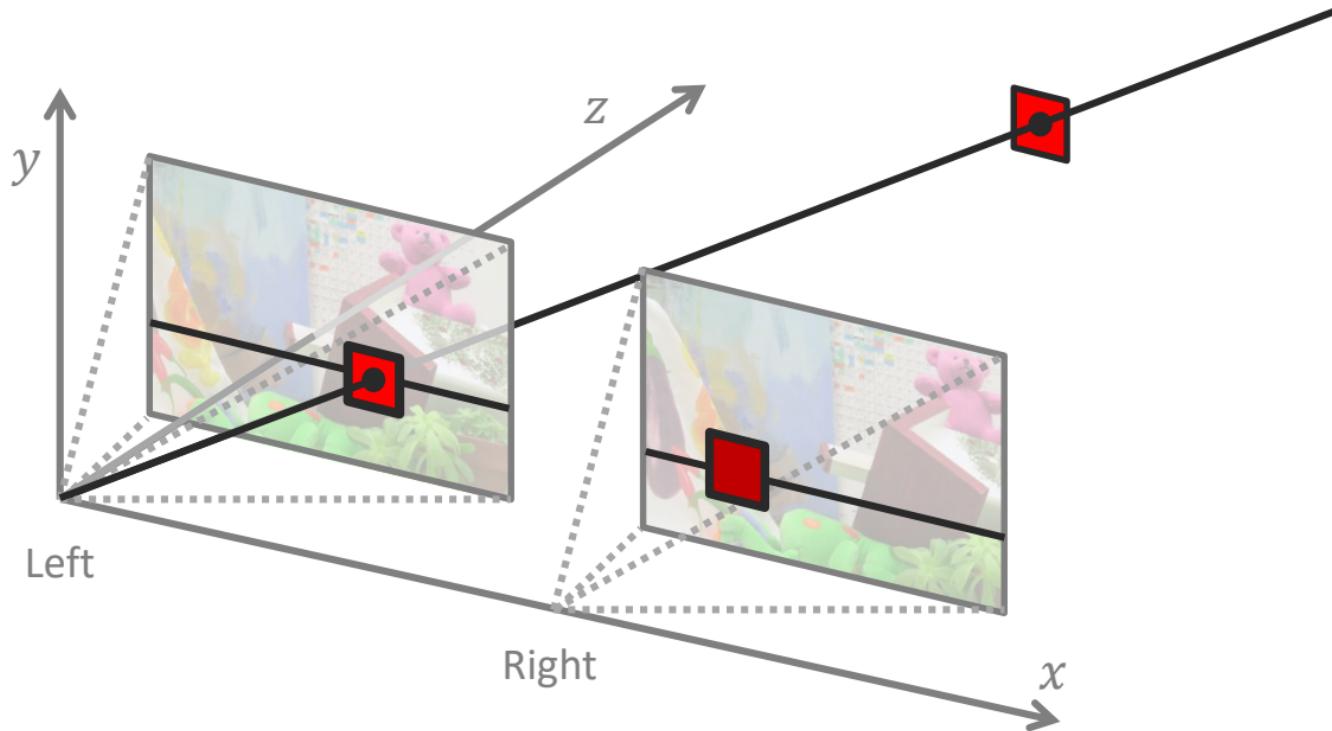
Outline

- Stereo Matching: New Trends
- Semi-Global Stereo Matching (SGM)
 - SGM with Surface Orientation Priors
 - Learning to Fuse Proposals
- Stereo Scene Flow with Motion Segmentation
- Camera Trajectory Planning for Aerial 3D Scanning
- Deep 6-DoF Object Pose Prediction

Stereo Matching



Stereo Matching



Solved Problem?

Still difficult ...



Fore-shortening



Specular surfaces



Transparency, reflections



Different lighting



High Dynamic range

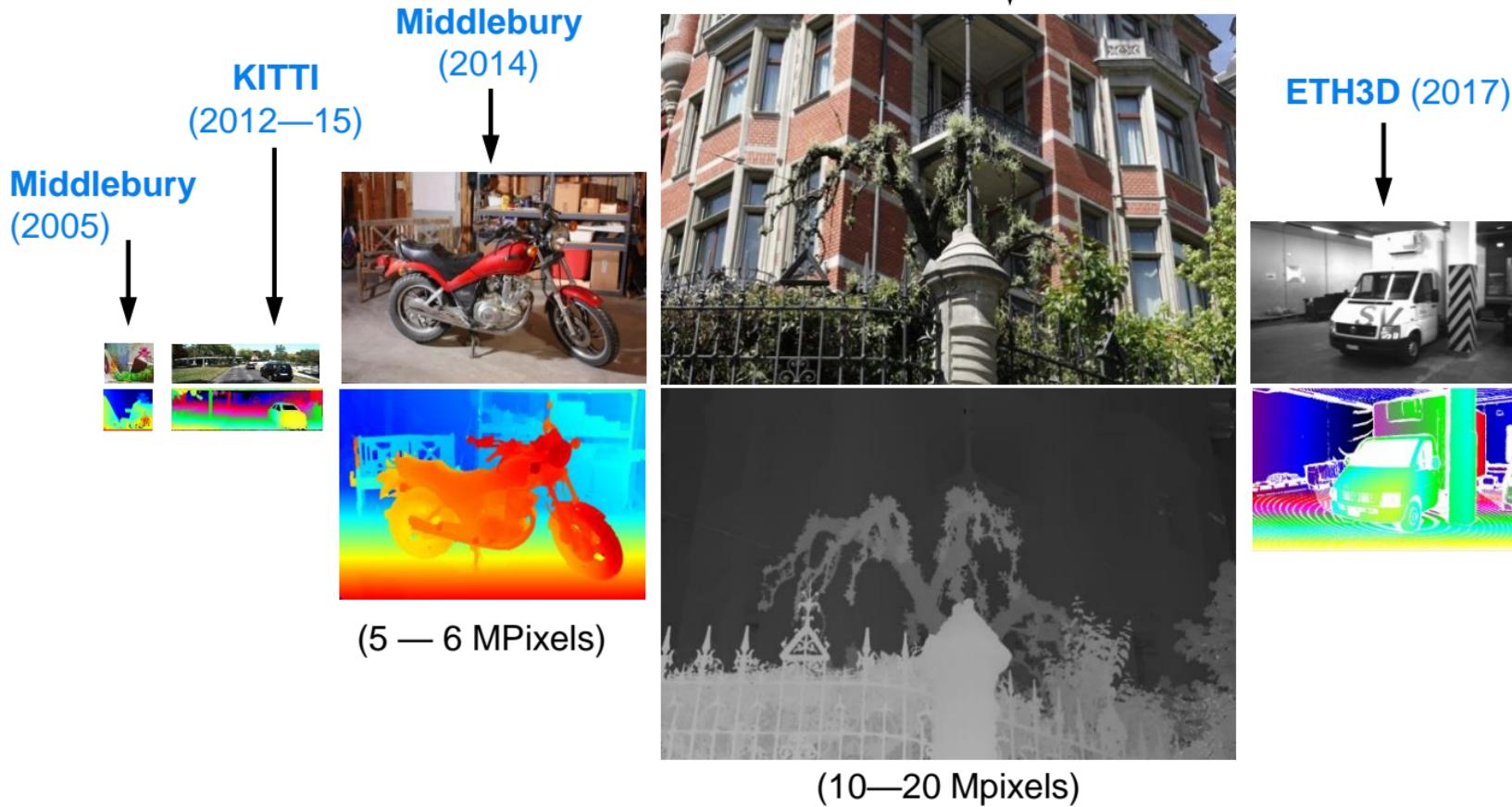


Untextured slanted surfaces

- Top methods still inaccurate in corner cases and infeasible for real-time, low power systems

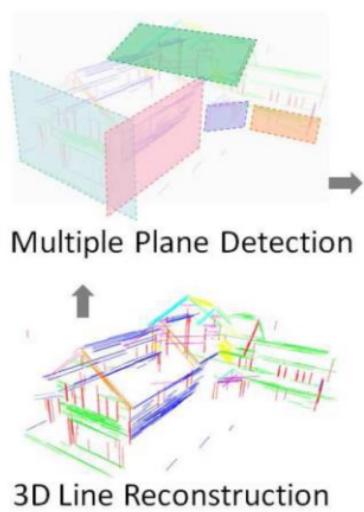
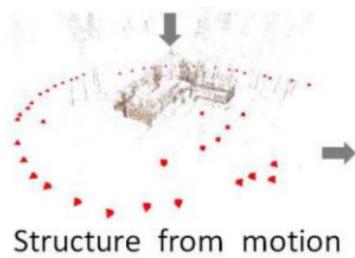
Stereo benchmarks

Disney Research
(2013)

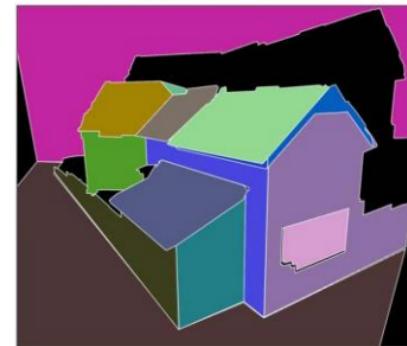


Piecewise Planar Stereo

[Sinha, Steedly, Szeliski 2009]



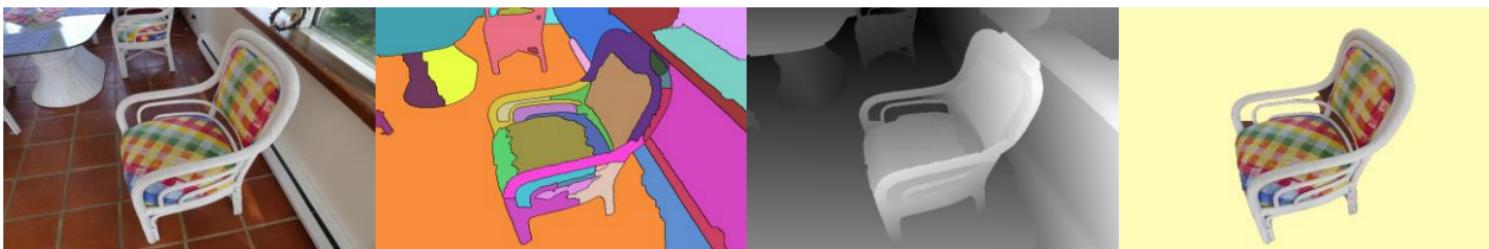
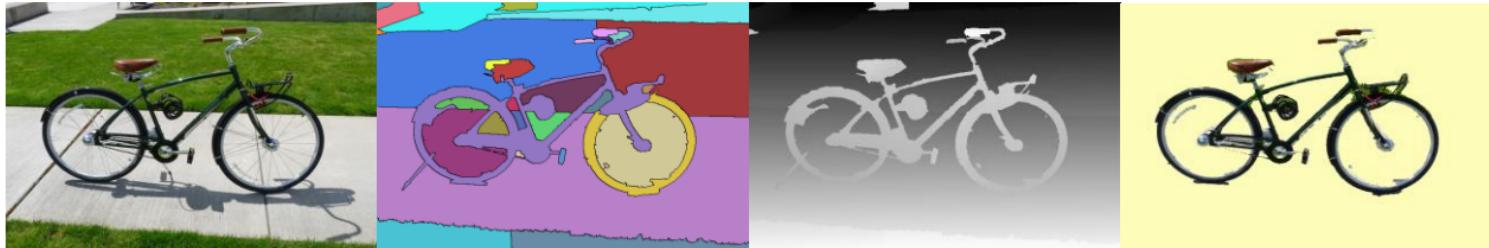
Discrete Optimization (Graph Cuts)



Novel Rendered View

Piecewise Planar Stereo Revisited

[Kowdle, Sinha, Szeliski 2012]



Labeling Problems and Markov Random Fields

- Find a per-pixel label (disparity map) D , by minimizing energy:

$$E(D) = E_{\text{data}}(D) + E_{\text{smooth}}(D)$$

$$= \sum_{p \in I} C_p(d_p) + \sum_{(p,q) \in N} V_{pq}(d_p, d_q)$$

- Data and smoothness terms encode matching costs and prior
- Discrete or continuous labels
- Inference: Graph Cuts, Belief Propagation, PatchMatch, ...

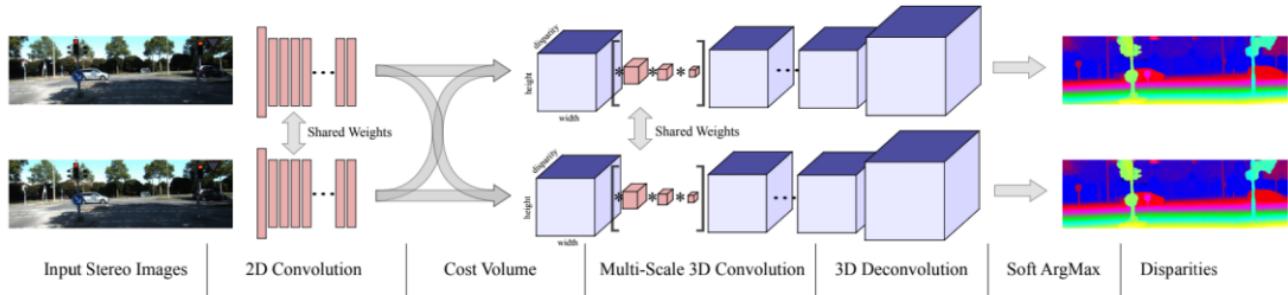
[Boykov+ 1998; Kolmogorov and Zabih 2002; Felzenszwalb and Huttenlocher 2004; Bleyer+ 2011; Besse+ 2012]

New Trends

- Learning the matching cost:
 - MC-CNN [[Zbontar + Lecun 2015](#)], [Chen+ 2015](#), [Luo+ 2016](#)
- Continuous MRFs:
 - Piecewise Planar Stereo [[Taniai+ 2017](#)] (Top rank on Middlebury v3 since 2017)
- Deep stereo regression (end to end training)
 - FlowNet [[Dosovitskiy+ 2015](#)]
 - DispNet [[Mayer+ 2016](#)]
- Unsupervised learning
 - Left-Right Consistency Check → training samples, iterate ... [[Zhou+ ICCV 2017](#)]

New Trends

Kendall+ 2017

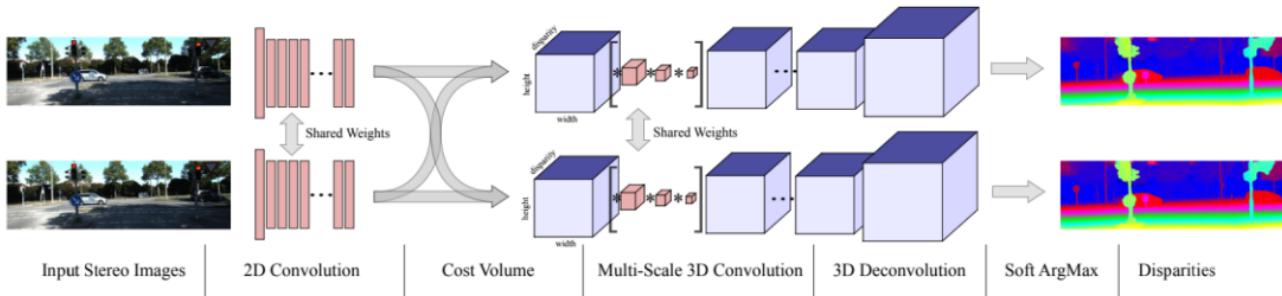


- **Return of “Correlation”**

- DispNetCorr [Mayer+ 2016]
- GC-Net [Kendall+ 2017]

New Trends

Kendall+ 2017

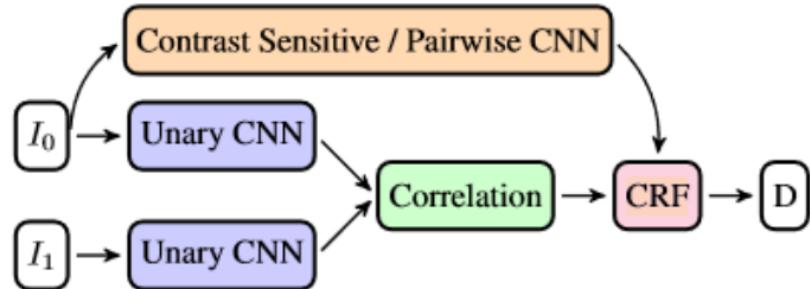


- **Return of “Correlation”**

- DispNetCorr [Mayer+ 2016]
- GC-Net [Kendall+ 2017]

- **Return of “CRFs”**

- SGM-Nets [Seki and Pollefeys 2017]
- Hybrid CNN-CRF [Knobelreiter+ 2017]



Knobelreiter+ 2017

Analysis of Benchmarks (~June 2018)

Middlebury 2014

Mouseover the table cells to see the produced disparity map. Clicking a cell will blink the ground truth for comparison. To change the table type, click the links below. For more information, please see the [description of new features](#).

Submit and evaluate your own results. See [snapshots of previous results](#). See the [evaluation v.2](#) (no longer active).

Set: [test dense](#) [test sparse](#) training dense training sparse

Metric: [bad](#) [bad1.0](#) [bad2.0](#) [bad4.0](#) [avgerr](#) [rms](#) [A50](#) [A90](#) [A95](#) [A99](#) [time](#) [time/MP](#) [time/GD](#)

Mask: [nonecc](#) [all](#)

plot selected show invalid Reset sort Reference list

Date	Name	Res	Avg	Austr		Austri		Bucov		Class		Classi		Compu		Crossa		Crossa		Dambi		Dambi		Dambi		Housa		Liverp		Nubia		Plants		Stairs					
				MP-5.6																																			
09/06/16	NOSS	H	5.04	3.57	2.84	3.99	1.93	5.15	3.34	3.32	3.15	2.32	8.55	2.22	7.49	7.03	12.5	5.20	10.0	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7					
06/22/17	LocalExp	H	5.93	3.65	2.87	2.88	1.99	5.59	2.37	3.43	3.55	2.06	1.03	9.73	8.57	14.4	5.40	9.55	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7					
01/24/17	3DMS	H	5.92	3.71	2.78	4.75	2.72	7.36	4.28	3.44	3.78	2.35	12.6	11.5	8.96	14.0	5.35	8.87	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7					
09/10/17	MC-CNN+TDSR	F	3.95	4.45	4.15	8.80	3.46	10.7	4.05	5.01	5.19	2.82	10.8	9.62	6.59	11.4	6.01	7.04	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7					
09/12/16	PMSC	H	8.71	3.46	2.68	6.19	2.54	6.02	4.54	3.96	4.04	2.37	13.1	12.3	12.7	16.7	5.88	10.8	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7					
10/19/16	LW-CNN	H	7.04	4.65	3.75	5.30	2.63	11.2	5.41	4.32	4.22	2.43	12.2	13.4	13.6	14.8	14.9	4.72	12.0	14	14	14	14	14	14	14	14	14	14	14	14	14	14	14	14	14			
04/12/16	MeshStereEx	H	7.08	4.41	3.98	5.40	3.17	10.7	6.23	4.87	4.77	3.49	12.7	12.4	10.4	14.5	7.89	8.85	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7				
10/12/17	FEN-DOOR	H	7.23	4.59	4.11	5.03	3.03	8.42	4.05	6.05	4.90	5.28	3.20	11.5	14.1	13.4	13.9	5.09	14.3	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7			
05/28/16	APAP-Stereo	H	7.26	6.43	4.91	21.51	5.17	21.39	6.99	4.31	4.23	3.24	14.3	9.78	7.32	13.2	4.34	6.30	8.49	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7	10.7			
03/11/18	SGM-Forest	H	7.37	4.71	3.89	4.93	3.18	11.1	5.73	5.57	5.81	2.65	14.5	13.2	13.1	14.8	5.53	11.2	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10				
01/19/16	NTDE	H	7.44	6.72	4.76	11.92	5.92	2.80	10.4	5.71	5.30	5.54	2.40	13.5	14.1	12.6	13.9	6.39	12.2	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10			
02/28/18	FDR	H	7.89	5.41	4.22	12.4	4.20	2.73	10.2	5.40	6.40	4.57	4.72	11.2	5.5	13.4	16.5	5.22	13.0	17	17	17	17	17	17	17	17	17	17	17	17	17	17	17	17	17	17		
09/25/15	MC-CNN-act	H	8.08	5.59	4.55	5.96	2.83	11.4	5.81	8.32	8.89	2.71	16.3	14.1	13.2	13.0	6.40	14.0	11.1	11	11	11	11	11	11	11	11	11	11	11	11	11	11	11	11	11	11		
11/03/15	MC-CNN+RBS	H	8.42	6.05	5.10	27.2	3.27	11.1	4.36	1.87	9.83	3.21	15.1	15	15	12.8	13.5	7.04	9.99	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10		
09/30/16	SNP-RSM	H	8.75	5.46	4.85	21.60	3.37	10.4	7.31	1.87	9.27	3.58	14.3	14.7	14.9	12.8	13.1	10.1	10.8	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10		
12/11/17	CVOD	H	8.87	4.74	3.64	5.51	4.82	12.8	6.51	1.91	9.96	3.13	16.6	14.9	14.1	15.4	16.9	6.92	14.3	12.0	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	
01/21/16	MCCNN_Layoff	H	8.94	5.53	5.63	5.06	3.59	12.8	7.23	7.53	8.99	5.79	23.0	22	13.6	15.0	14.7	5.85	10.4	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10	10		
04/26/16	MCCNN-FM	H	8.47	7.35	24	5.07	7.18	4.71	16.8	24	8.47	7.37	6.97	2.82	20.7	21	17.4	15.4	15.1	7.90	12.6	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12	12

KITTI 2015

KITTI 2015									
Environment									
1 core	5.7 Gb (C/C++)	5.7 Gb (Python)	5.7 Gb (C/C++)	5.7 Gb (Python)	5.7 Gb (C/C++)	5.7 Gb (Python)	5.7 Gb (C/C++)	5.7 Gb (Python)	5.7 Gb (C/C++)
GPU	2.5 Gb (Python)	2.5 Gb (Python)	2.5 Gb (Python)	2.5 Gb (Python)	2.5 Gb (Python)	2.5 Gb (Python)	2.5 Gb (Python)	2.5 Gb (Python)	2.5 Gb (Python)
1 core	1.0 Gb (Python)	1.0 Gb (Python)	1.0 Gb (Python)	1.0 Gb (Python)	1.0 Gb (Python)	1.0 Gb (Python)	1.0 Gb (Python)	1.0 Gb (Python)	1.0 Gb (Python)
GPU	0.4 Gb (Python)	0.4 Gb (Python)	0.4 Gb (Python)	0.4 Gb (Python)	0.4 Gb (Python)	0.4 Gb (Python)	0.4 Gb (Python)	0.4 Gb (Python)	0.4 Gb (Python)
1 core	0.2 Gb (Python)	0.2 Gb (Python)	0.2 Gb (Python)	0.2 Gb (Python)	0.2 Gb (Python)	0.2 Gb (Python)	0.2 Gb (Python)	0.2 Gb (Python)	0.2 Gb (Python)
GPU	0.1 Gb (Python)	0.1 Gb (Python)	0.1 Gb (Python)	0.1 Gb (Python)	0.1 Gb (Python)	0.1 Gb (Python)	0.1 Gb (Python)	0.1 Gb (Python)	0.1 Gb (Python)
1 core	0.05 Gb (Python)	0.05 Gb (Python)	0.05 Gb (Python)	0.05 Gb (Python)	0.05 Gb (Python)	0.05 Gb (Python)	0.05 Gb (Python)	0.05 Gb (Python)	0.05 Gb (Python)
GPU	0.02 Gb (Python)	0.02 Gb (Python)	0.02 Gb (Python)	0.02 Gb (Python)	0.02 Gb (Python)	0.02 Gb (Python)	0.02 Gb (Python)	0.02 Gb (Python)	0.02 Gb (Python)
1 core	0.01 Gb (Python)	0.01 Gb (Python)	0.01 Gb (Python)	0.01 Gb (Python)	0.01 Gb (Python)	0.01 Gb (Python)	0.01 Gb (Python)	0.01 Gb (Python)	0.01 Gb (Python)
GPU	0.005 Gb (Python)	0.005 Gb (Python)	0.005 Gb (Python)	0.005 Gb (Python)	0.005 Gb (Python)	0.005 Gb (Python)	0.005 Gb (Python)	0.005 Gb (Python)	0.005 Gb (Python)
KITTI 2015									
Environment	KITTI 2015	KITTI 2015	KITTI 2015	KITTI 2015	KITTI 2015	KITTI 2015	KITTI 2015	KITTI 2015	KITTI 2015
1 core	0.01 Gb (Python)	0.01 Gb (Python)	0.01 Gb (Python)	0.01 Gb (Python)	0.01 Gb (Python)	0.01 Gb (Python)	0.01 Gb (Python)	0.01 Gb (Python)	0.01 Gb (Python)
GPU	0.005 Gb (Python)	0.005 Gb (Python)	0.005 Gb (Python)	0.005 Gb (Python)	0.005 Gb (Python)	0.005 Gb (Python)	0.005 Gb (Python)	0.005 Gb (Python)	0.005 Gb (Python)
1 core	0.002 Gb (Python)	0.002 Gb (Python)	0.002 Gb (Python)	0.002 Gb (Python)	0.002 Gb (Python)	0.002 Gb (Python)	0.002 Gb (Python)	0.002 Gb (Python)	0.002 Gb (Python)
GPU	0.001 Gb (Python)	0.001 Gb (Python)	0.001 Gb (Python)	0.001 Gb (Python)	0.001 Gb (Python)	0.001 Gb (Python)	0.001 Gb (Python)	0.001 Gb (Python)	0.001 Gb (Python)
1 core	0.001 Gb (Python)	0.001 Gb (Python)	0.001 Gb (Python)	0.001 Gb (Python)	0.001 Gb (Python)	0.001 Gb (Python)	0.001 Gb (Python)	0.001 Gb (Python)	0.001 Gb (Python)
GPU	0.0005 Gb (Python)	0.0005 Gb (Python)	0.0005 Gb (Python)	0.0005 Gb (Python)	0.0005 Gb (Python)	0.0005 Gb (Python)	0.0005 Gb (Python)	0.0005 Gb (Python)	0.0005 Gb (Python)
1 core	0.0002 Gb (Python)	0.0002 Gb (Python)	0.0002 Gb (Python)	0.0002 Gb (Python)	0.0002 Gb (Python)	0.0002 Gb (Python)	0.0002 Gb (Python)	0.0002 Gb (Python)	0.0002 Gb (Python)
GPU	0.0001 Gb (Python)	0.0001 Gb (Python)	0.0001 Gb (Python)	0.0001 Gb (Python)	0.0001 Gb (Python)	0.0001 Gb (Python)	0.0001 Gb (Python)	0.0001 Gb (Python)	0.0001 Gb (Python)
1 core	0.00005 Gb (Python)	0.00005 Gb (Python)	0.00005 Gb (Python)	0.00005 Gb (Python)	0.00005 Gb (Python)	0.00005 Gb (Python)	0.00005 Gb (Python)	0.00005 Gb (Python)	0.00005 Gb (Python)
GPU	0.00002 Gb (Python)	0.00002 Gb (Python)	0.00002 Gb (Python)	0.00002 Gb (Python)	0.00002 Gb (Python)	0.00002 Gb (Python)	0.00002 Gb (Python)	0.00002 Gb (Python)	0.00002 Gb (Python)
1 core	0.00001 Gb (Python)	0.00001 Gb (Python)	0.00001 Gb (Python)	0.00001 Gb (Python)	0.00001 Gb (Python)	0.00001 Gb (Python)	0.00001 Gb (Python)	0.00001 Gb (Python)	0.00001 Gb (Python)
GPU	0.000005 Gb (Python)	0.000005 Gb (Python)	0.000005 Gb (Python)	0.000005 Gb (Python)	0.000005 Gb (Python)	0.000005 Gb (Python)	0.000005 Gb (Python)	0.000005 Gb (Python)	0.000005 Gb (Python)
1 core	0.000002 Gb (Python)	0.000002 Gb (Python)	0.000002 Gb (Python)	0.000002 Gb (Python)	0.000002 Gb (Python)	0.000002 Gb (Python)	0.000002 Gb (Python)	0.000002 Gb (Python)	0.000002 Gb (Python)
GPU	0.000001 Gb (Python)	0.000001 Gb (Python)	0.000001 Gb (Python)	0.000001 Gb (Python)	0.000001 Gb (Python)	0.000001 Gb (Python)	0.000001 Gb (Python)	0.000001 Gb (Python)	0.000001 Gb (Python)
1 core	0.0000005 Gb (Python)	0.0000005 Gb (Python)	0.0000005 Gb (Python)	0.0000005 Gb (Python)	0.0000005 Gb (Python)	0.0000005 Gb (Python)	0.0000005 Gb (Python)	0.0000005 Gb (Python)	0.0000005 Gb (Python)
GPU	0.0000002 Gb (Python)	0.0000002 Gb (Python)	0.0000002 Gb (Python)	0.0000002 Gb (Python)	0.0000002 Gb (Python)	0.0000002 Gb (Python)	0.0000002 Gb (Python)	0.0000002 Gb (Python)	0.0000002 Gb (Python)
1 core	0.0000001 Gb (Python)	0.0000001 Gb (Python)	0.0000001 Gb (Python)	0.0000001 Gb (Python)	0.0000001 Gb (Python)	0.0000001 Gb (Python)	0.0000001 Gb (Python)	0.0000001 Gb (Python)	0.0000001 Gb (Python)
GPU	0.00000005 Gb (Python)	0.00000005 Gb							

Analysis of Benchmarks (~June 2018)

Middlebury 2014

Mouseover the table cells to see the produced disparity map. Clicking a cell will blink the ground truth for comparison. To change the table type, click the links below. For more information, please see the [description of new features](#).

Submit and evaluate your own results. See [snapshots of previous results](#). See the [evaluation v.2](#) (no longer active).

Set: **test dense** **test sparse** **training dense** **training sparse**
 Metric: **bad**, **bad1.0**, **bad2.0**, **bad4.0**, **avgerr**, **rms**, **A50**, **A90**, **A95**, **A99**, **time**, **time/MP**, **time/GO**

Mask: **nomono**, **all**

plot selected show invalid Reset sort Reference list

Date	Test	Res	Avg	Austr	Austr	Bucov2	Class	ClassE	Compu	Crossa	Crossa	Dambi	Dambi	Hevesi	Livgren	Neubu	Plants	Stans	
06/06/16	NOSS	H	8.04	3.57	2.94	3.99	1.93	3.34	3.22	2.18	2.32	8.55	7.49	7.09	12.5	5.20	10.0		
06/22/17	LocalExp	H	5.43	3.65	2.87	2.98	1.99	5.59	2.37	3.49	3.35	2.06	10.3	9.79	8.57	14.4	5.40	9.55	
01/04/17	3DMS	H	5.92	3.71	2.70	4.76	2.72	7.36	4.28	3.42	3.78	2.35	12.6	11.5	8.96	14.0	5.35	8.87	
09/19/17	MC-CNN+TDSR	SR	4.95	4.45	4.45	17.00	2.46	10.7	0.05	5.01	5.19	2.82	10.8	9.42	6.59	11.4	6.01	7.24	
05/12/16	PMSC	F	8.71	3.48	2.55	19	2.54	6.92	4.54	3.99	4.57	2.37	13.1	12.3	12.7	16.2	5.88	10.8	
10/15/16	LW-CNN	H	7.04	4.45	4.45	17.00	2.46	10.7	0.05	5.01	5.19	2.82	10.8	9.42	6.59	11.4	6.01	7.24	
04/12/16	MeshStereosExt	H	7.00	4.41	4.41	17.00	2.46	10.7	0.05	5.01	5.19	2.82	10.8	9.42	6.59	11.4	6.01	7.24	
10/12/17	FEN-D2DR	H	7.23	4.59	4.11	5.03	3.03	8.42	4.05	4.90	5.22	3.20	11.5	14.1	11.4	13.9	5.05	14.3	
05/29/16	APAP-Steer	H	7.26	4.43	4.91	21.51	5.17	21.6	11.9	4.31	4.23	3.24	14.3	9.78	7.32	13.4	6.30	16.49	
03/11/18	SGM-Forest	H	7.37	4.71	3.89	1.11	5.37	5.57	5.81	2.65	14.5	12.2	13.1	11	14.8	5.53	11.2		
01/19/16	NTDE	H	7.44	6.72	17.46	5.92	2.83	10.4	5.71	5.20	5.64	2.40	13.5	11.4	12.8	13.9	6.39	12.2	
05/29/18	CD2D	H	7.69	6.41	4.41	6.76	3.73	10.7	4.40	4.61	6.76	4.73	11.7	16.4	13.4	16.5	6.32	13.0	
06/25/15	MC-CNN-acrt	H	8.08	5.59	4.55	17.56	2.83	11.4	5.81	8.32	8.89	2.71	16.3	14.1	13.2	13.0	6.40	11.1	
11/03/15	MC-CNN+RBS	H	8.42	5.15	5.15	27.2	3.27	11.1	4.4	8.87	9.83	3.21	15.1	15.9	12.8	13.5	7.04	9.99	
09/13/16	SNP-RSM	H	7.75	4.46	14.85	21.65	3.37	10.4	7.17	8.73	9.27	3.58	14.3	14.7	14.9	12.8	10.1	10.8	
12/11/17	CVOD	H	8.87	4.74	3.64	5.51	4.82	12.8	6.51	14.93	9.96	3.13	16.6	14.8	14.1	15.4	6.92	13.2	
07/21/16	MC-CNN_Layout	H	8.94	5.53	5.63	5.06	3.59	12.6	7.53	8.69	5.79	23.0	13.6	15.0	14.7	5.85	10.4		
01/26/16	MC-CNN-Net	H	8.47	7.35	24	5.07	7.18	4.71	16.8	24	8.47	7.37	6.97	2.82	20.7	21	17.4	15.4	15.1

Group A

#13

MC-CNN acrt
[Zbontar Lecun 2015]

#15

KITTI 2015									
1	Baseline	2.14	8	3.46	2.36	100.00	W	0.23	X
2	Baseline	2.10	8	3.49	2.40	100.00	W	0.25	X
3	Baseline	2.10	8	3.49	2.44	100.00	W	0.22	X
4	Baseline	2.10	8	4.01	2.47	100.00	W	0.8	X
5	Baseline	2.10	8	2.35	2.50	100.00	W	0.1	X
6	Baseline	2.10	8	4.75	2.99	100.00	W	0.9	X
7	Baseline	2.10	8	4.59	2.81	100.00	W	0.28	X
8	Baseline	2.10	8	3.49	2.47	100.00	W	0.47	X
9	Baseline	2.10	8	3.49	2.47	100.00	W	0.47	X
10	Baseline	2.10	8	4.59	2.81	100.00	W	0.28	X
11	Baseline	2.10	8	3.49	2.47	100.00	W	0.47	X
12	Baseline	2.10	8	3.49	2.47	100.00	W	0.47	X
13	Baseline	2.10	8	3.49	2.47	100.00	W	0.47	X
14	Baseline	2.10	8	3.49	2.47	100.00	W	0.47	X
15	Baseline	2.10	8	3.49	2.47	100.00	W	0.47	X
16	Baseline	2.10	8	3.49	2.47	100.00	W	0.47	X
17	Baseline	2.10	8	3.49	2.47	100.00	W	0.47	X
18	Baseline	2.10	8	3.49	2.47	100.00	W	0.47	X
19	Baseline	2.10	8	3.49	2.47	100.00	W	0.47	X
20	Baseline	2.10	8	3.49	2.47	100.00	W	0.47	X
21	MC-CNN	2.05	8	4.49	2.31	100.00	W	0.9	X
22	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
23	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
24	MC-CNN	2.07	8	4.49	2.31	100.00	W	0.9	X
25	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
26	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
27	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
28	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
29	MC-CNN	2.07	8	4.49	2.31	100.00	W	0.9	X
30	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
31	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
32	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
33	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
34	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
35	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
36	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
37	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
38	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
39	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
40	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
41	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
42	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
43	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
44	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
45	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
46	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
47	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
48	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
49	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
50	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
51	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
52	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
53	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
54	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
55	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
56	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
57	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
58	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
59	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
60	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
61	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
62	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
63	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
64	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
65	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
66	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
67	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
68	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
69	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
70	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
71	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
72	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
73	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
74	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
75	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
76	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
77	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
78	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
79	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
80	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
81	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
82	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
83	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
84	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
85	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
86	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
87	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
88	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
89	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
90	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
91	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
92	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
93	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
94	MC-CNN	2.07	8	3.49	2.32	100.00	W	0.8	X
95	MC-CNN								

Summary

- Dataset bias exists.
- Middlebury stereo pairs from different scenes; makes learning difficult.
- Need a better way to evaluate “deep stereo regression”.

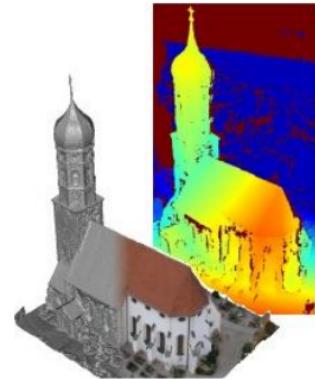
Must train one model on combined training set and submit to all benchmarks!

Outline

- Stereo Matching: New Trends
- **Semi-Global Stereo Matching (SGM)**
 - SGM with Surface Orientation Priors
 - Learning to Fuse Proposals
- Stereo Scene Flow with Motion Segmentation
- Camera Trajectory Planning for Aerial 3D Scanning
- Deep 6-DoF Object Pose Prediction

Semi Global Matching [Hirschmüller 2005]

- MRF inference (Graph Cuts, BP, ..) too slow
- SGM: Approximate even more; use heuristics
 - Widely used: assisted driving, robotics, aerial mapping ...
 - Runs on GPUs, FPGAs ...



Scansline Optimization (1D)

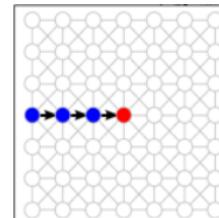
Minimize:

$$E(D) = \sum_{p \in I} C_p(d_p) + \sum_{(p,q) \in N} V_{pq}(d_p, d_q)$$

- Consider the above problem on a 1D scanline.
- Compute an aggregated matching cost

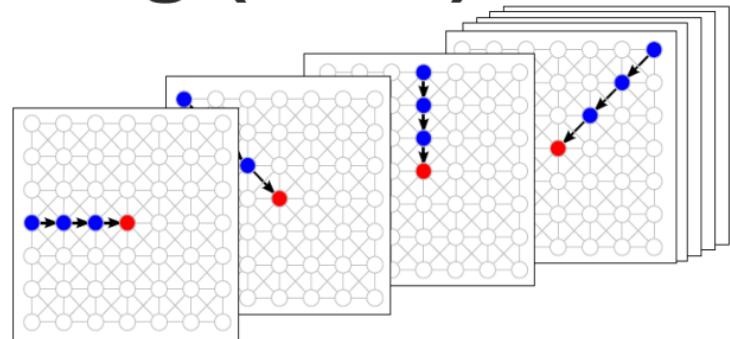
$$L_{\mathbf{r}}(\mathbf{p}, d) = C_{\mathbf{p}}(d) + \min_{d' \in \mathcal{D}} (L_{\mathbf{r}}(\mathbf{p} - \mathbf{r}, d') + V(d, d')).$$

- $\mathbf{r} = (1, 0)$: start at leftmost pixel, scan left



Semi Global Matching (SGM)

- For 8 directions
 - calculate aggregated costs



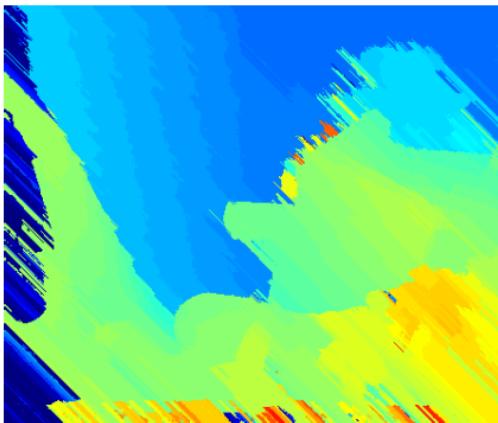
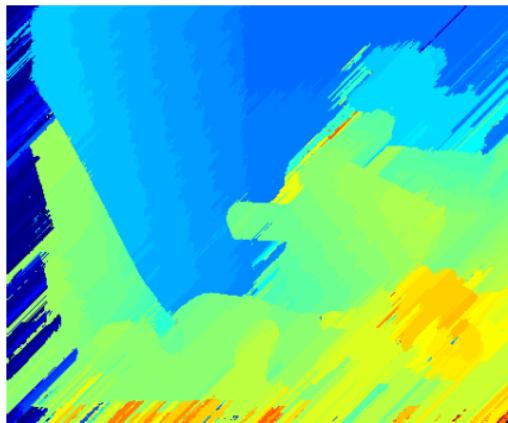
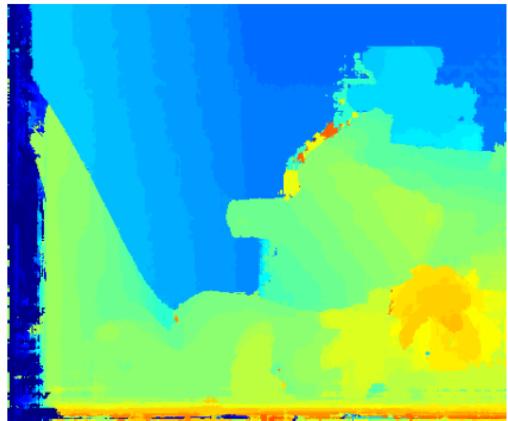
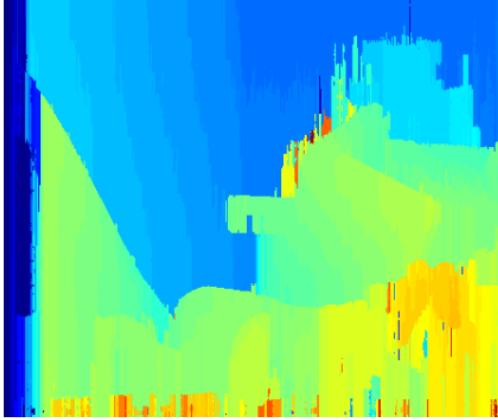
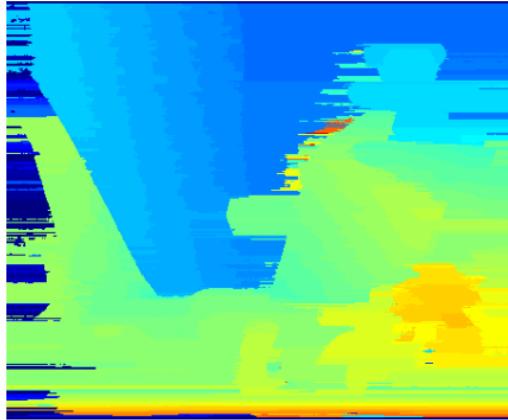
$$L_{\mathbf{r}}(\mathbf{p}, d) = C_{\mathbf{p}}(d) + \min_{d' \in \mathcal{D}} (L_{\mathbf{r}}(\mathbf{p} - \mathbf{r}, d') + V(d, d')).$$

- Finally, sum the costs and select per-pixel minima.

$$S(\mathbf{p}, d) = \sum_{\mathbf{r}} L_{\mathbf{r}}(\mathbf{p}, d)$$

$$D_{\mathbf{p}} = \arg \min_d S(\mathbf{p}, d).$$

Semi Global Matching (SGM)

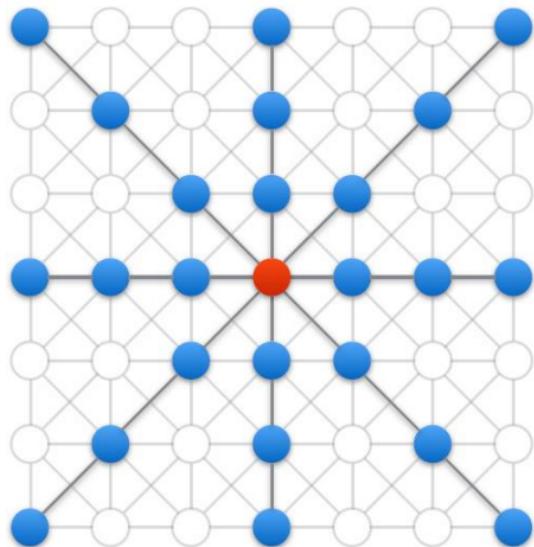


Semi Global Matching (SGM)

Approximates 2D MRF using 1D optimization
along 8 cardinal directions

$$E(D) = \sum_{\mathbf{p}} C_{\mathbf{p}}(d_{\mathbf{p}}) + \sum_{\mathbf{p}, \mathbf{q} \in \mathcal{N}} V(d_{\mathbf{p}}, d_{\mathbf{q}})$$

- Related to Belief Propagation, TRW-S
[Drory et al. 2014]



Semi Global Stereo Matching with Surface Orientation Priors



Daniel Scharstein

Middlebury College



Tatsunori Taniai

RIKEN, Tokyo



Sudipta Sinha

Microsoft Research

3DV 2017

Semi Global Matching [Hirschmüller 2005]

Approximates 2D MRF using 1D optimization along 8 cardinal directions

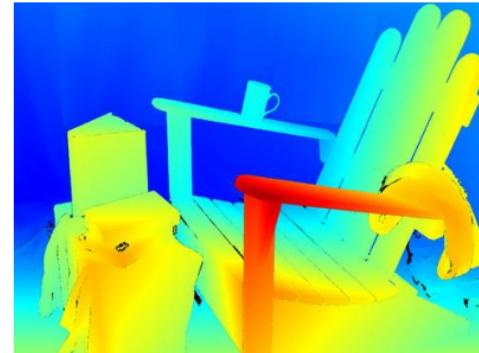
$$E(D) = \sum_{\mathbf{p}} C_{\mathbf{p}}(d_{\mathbf{p}}) + \sum_{\mathbf{p}, \mathbf{q} \in \mathcal{N}} V(d_{\mathbf{p}}, d_{\mathbf{q}})$$

$$\begin{cases} 0 & \text{if } d = d' \\ \infty & \text{otherwise} \end{cases}$$

- Fronto parallel bias
- Inaccurate on slanted untextured surfaces

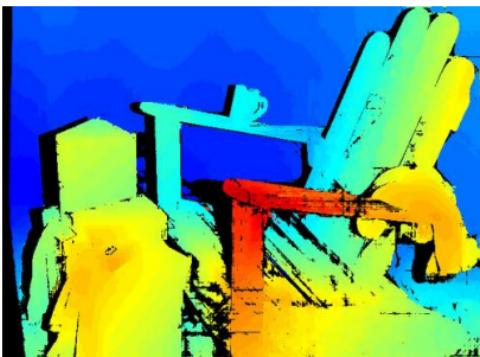


Left Image

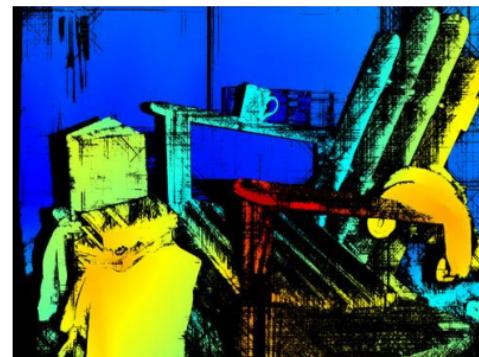


Ground Truth

SGM* (quarter resolution)



SGM* (full resolution (6 MP))



* Confidence measure used to prune uncertain pixels (black holes)

SGM-P: SGM with orientation priors

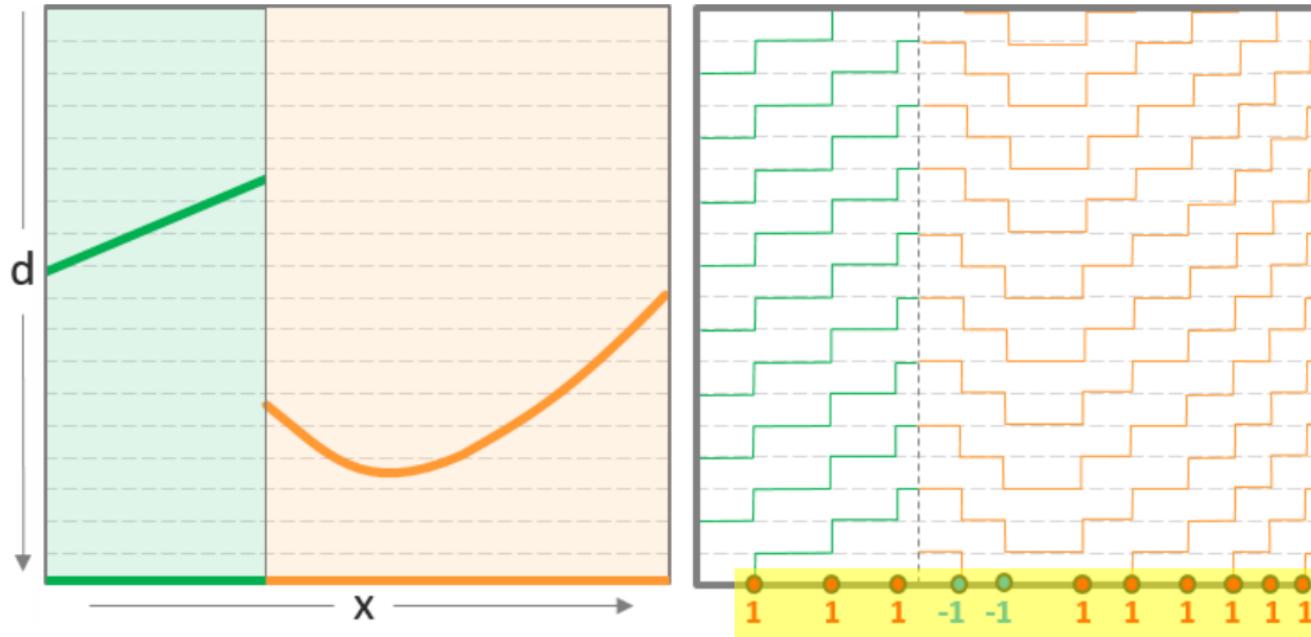
[Scharstein, Taniai, Sinha, 3DV 2017]

- What if we knew the surface slant?

Idea:

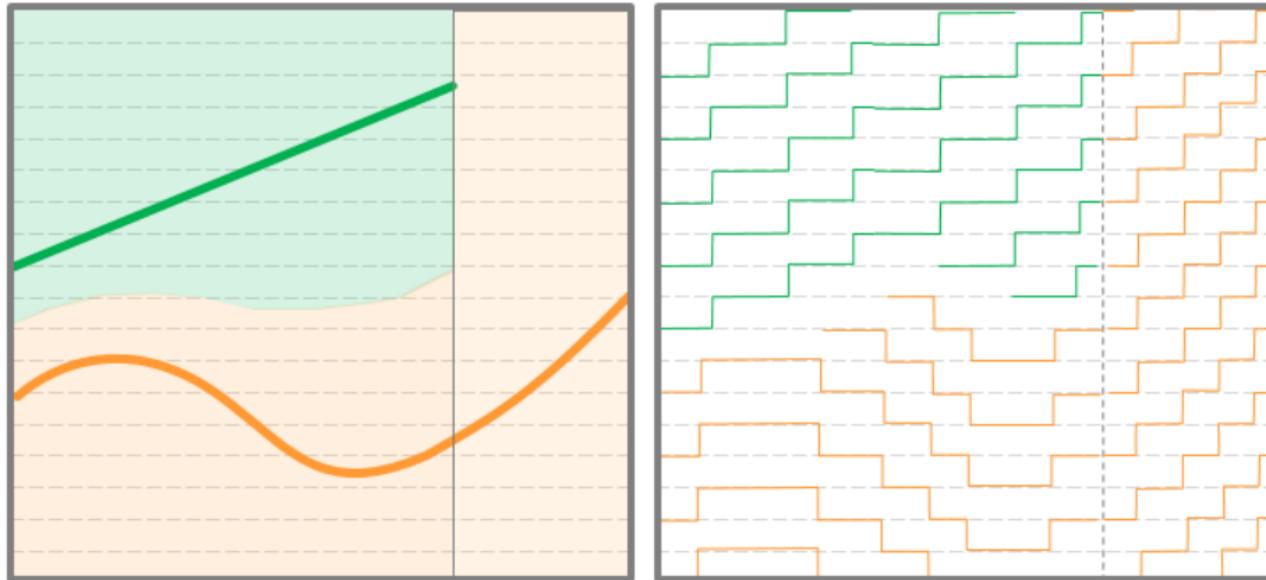
- *Rasterize* disparity surface prior (at arbitrary depth)
- Adjust $V(d, d')$ to follow discrete disparity “steps”

SGM-P: 2D orientation priors



$$V_S(d_p, d'_p) = V(d_p + j_p, d'_p)$$

SGM-P: 3D orientation priors



$$V_S(d_p, d'_p) = V(d_p + j_p(d_p), d'_p)$$

Jump locations
vary with disparity

SGM-P: Where do we get priors?

- Matched features + triangulation
- Matched features + plane fitting
- Low-res matching + plane fitting
- Ground truth oracle
- Semantic analysis
- Manhattan-world assumptions

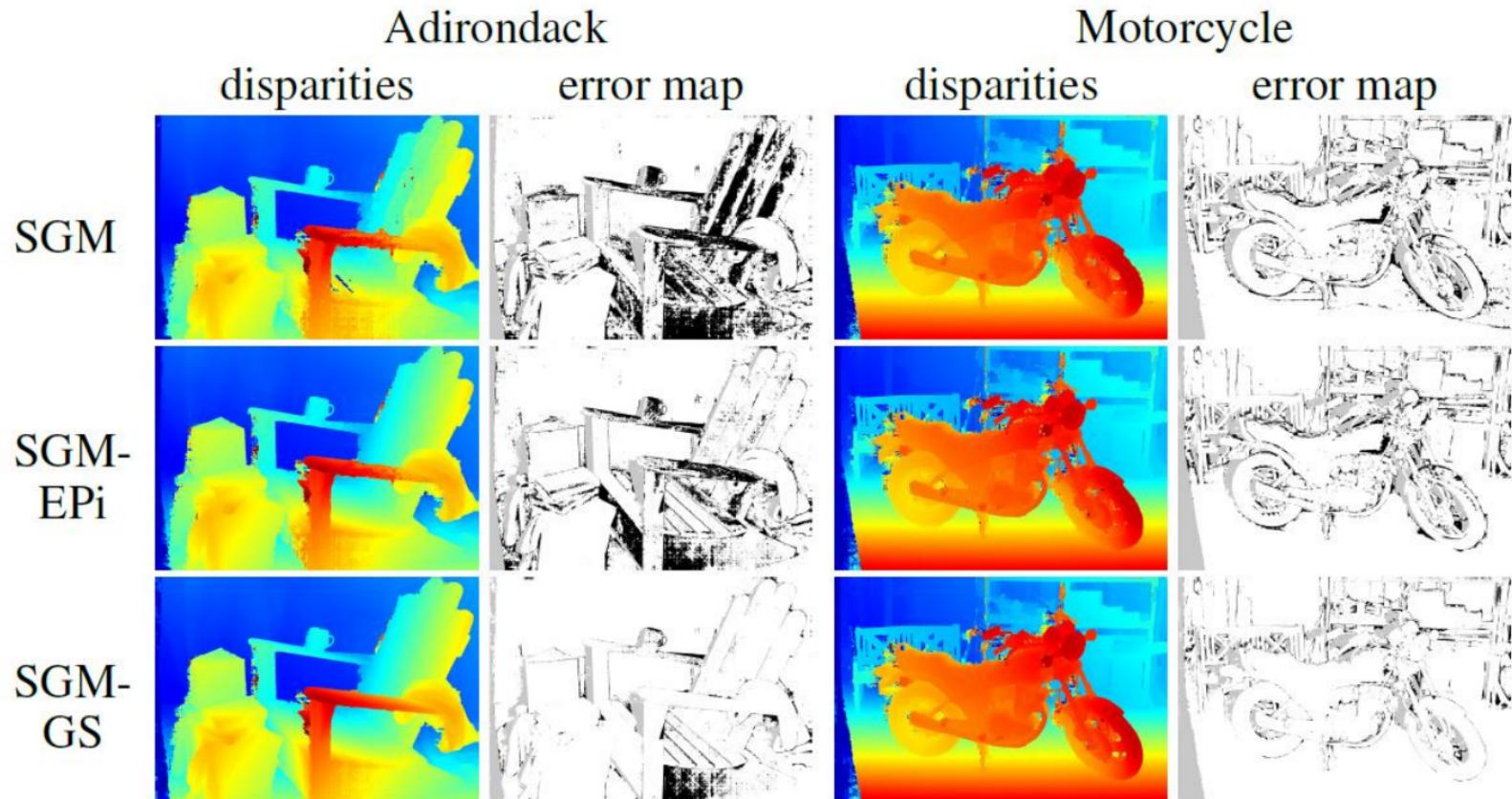


Input

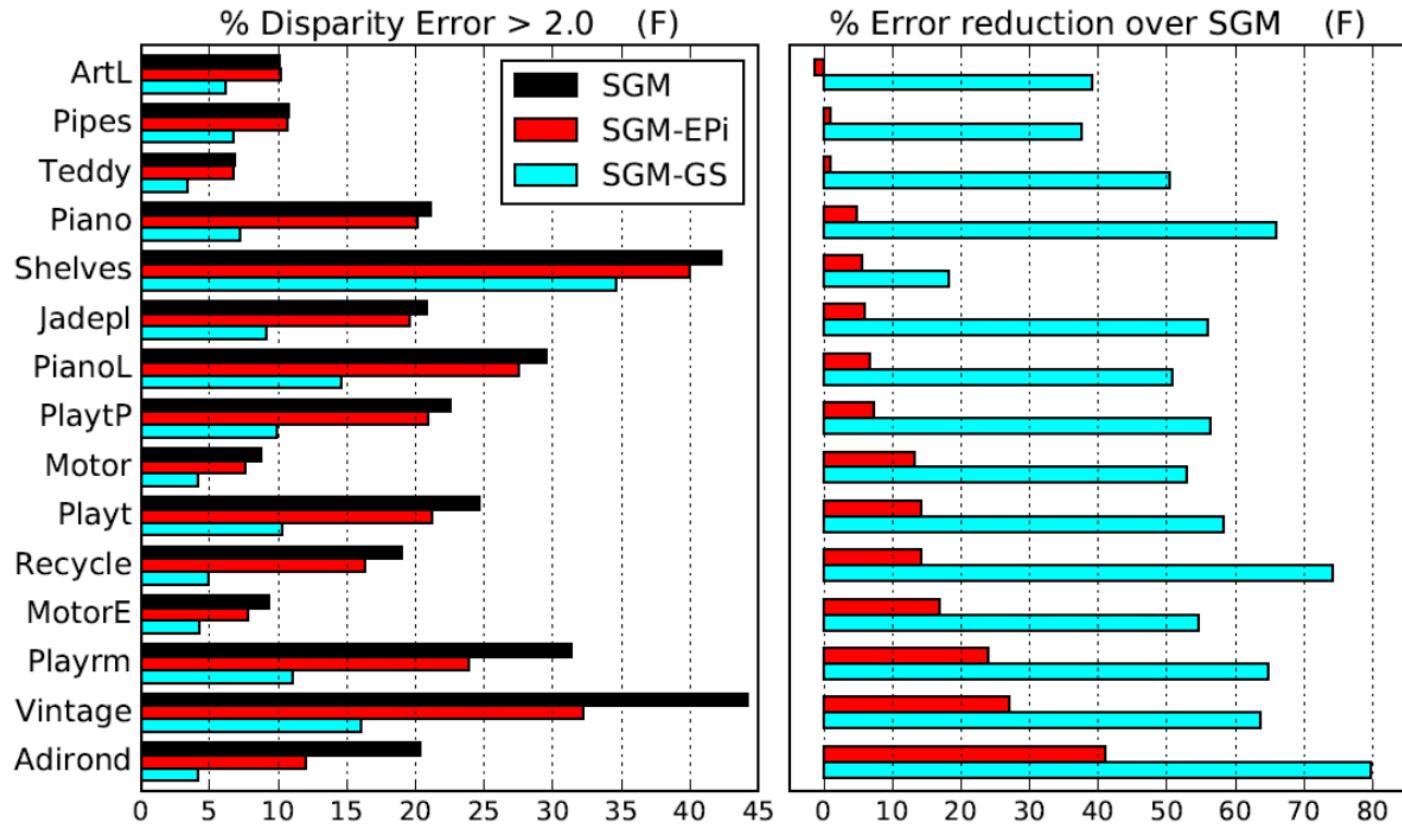
SGM-EPi

SGM-GS

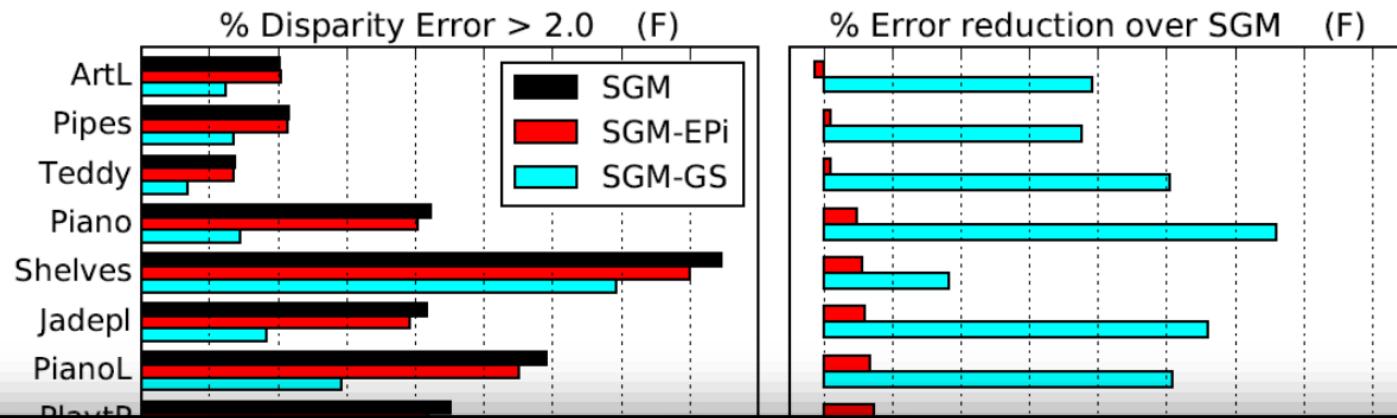
SGM-P: Results



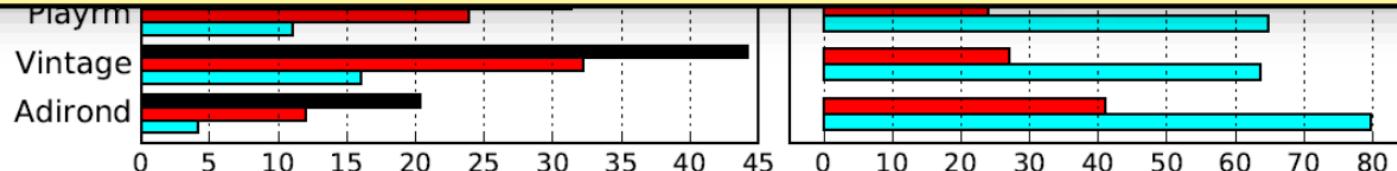
SGM-P: Results



SGM-P: Results



- Huge performance gains for slanted untextured scenes
- Soft constraint; inaccurate normals don't hurt accuracy



Learning to Fuse Proposals from Multiple Scanline Optimizations in Semi-Global Matching



Johannes Schönberger

ETH Zurich

Sudipta Sinha

Microsoft Research

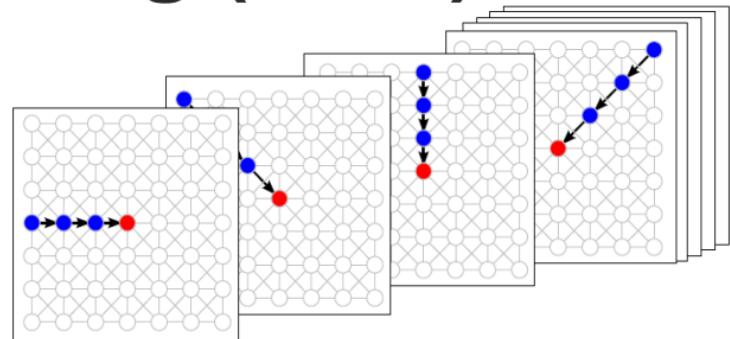
Marc Pollefeys

Microsoft / ETH Zurich

ECCV 2018 (to appear)

Semi Global Matching (SGM)

- For 8 directions
 - calculate aggregated costs



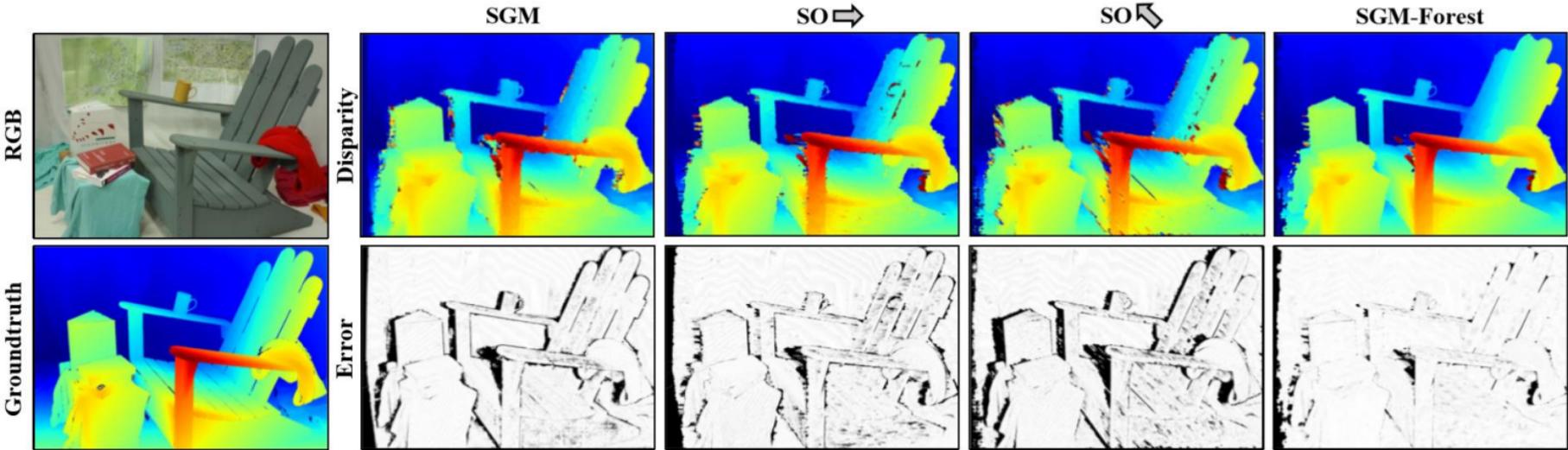
$$L_{\mathbf{r}}(\mathbf{p}, d) = C_{\mathbf{p}}(d) + \min_{d' \in \mathcal{D}} (L_{\mathbf{r}}(\mathbf{p} - \mathbf{r}, d') + V(d, d')).$$

- Finally, sum the costs and select per-pixel minima. → Ad-hoc step

$$S(\mathbf{p}, d) = \sum_{\mathbf{r}} L_{\mathbf{r}}(\mathbf{p}, d)$$

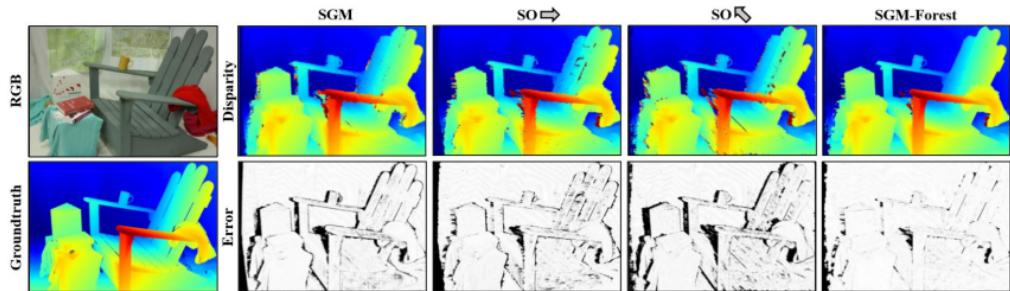
$$D_{\mathbf{p}} = \arg \min_d S(\mathbf{p}, d).$$

Main Idea

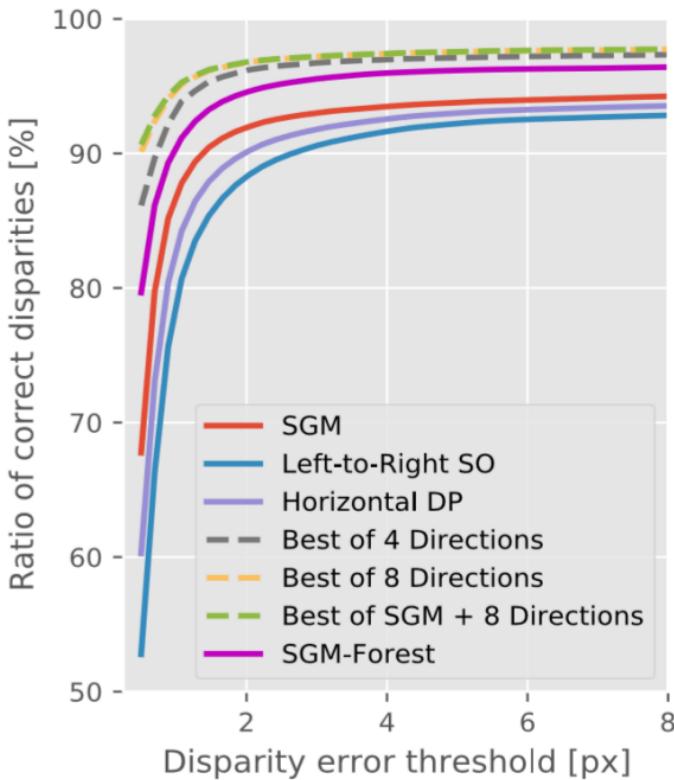


- Treat the 1D SO results as candidate solutions
- We propose to learn a function to select the best proposal
 - **SGM-Forest** is based on random forest classifiers

Main Idea

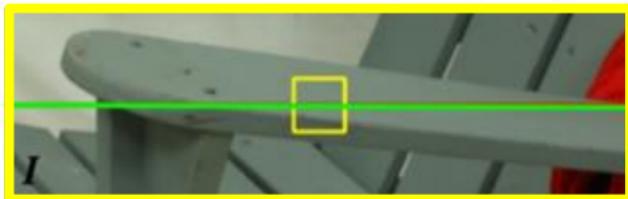
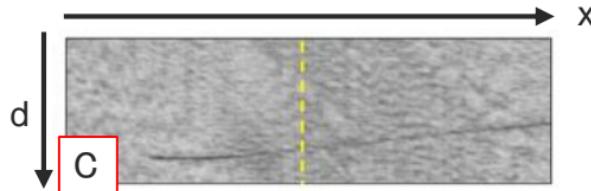
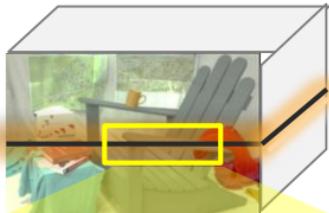


- “Best of k directions” oracle gives an upper bound
- Large gap between SGM and oracle
- SGM-Forest closes the gap a fair bit



Analyzing the Scanline Optimization Costs

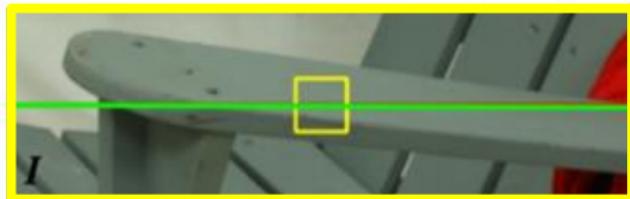
DSI



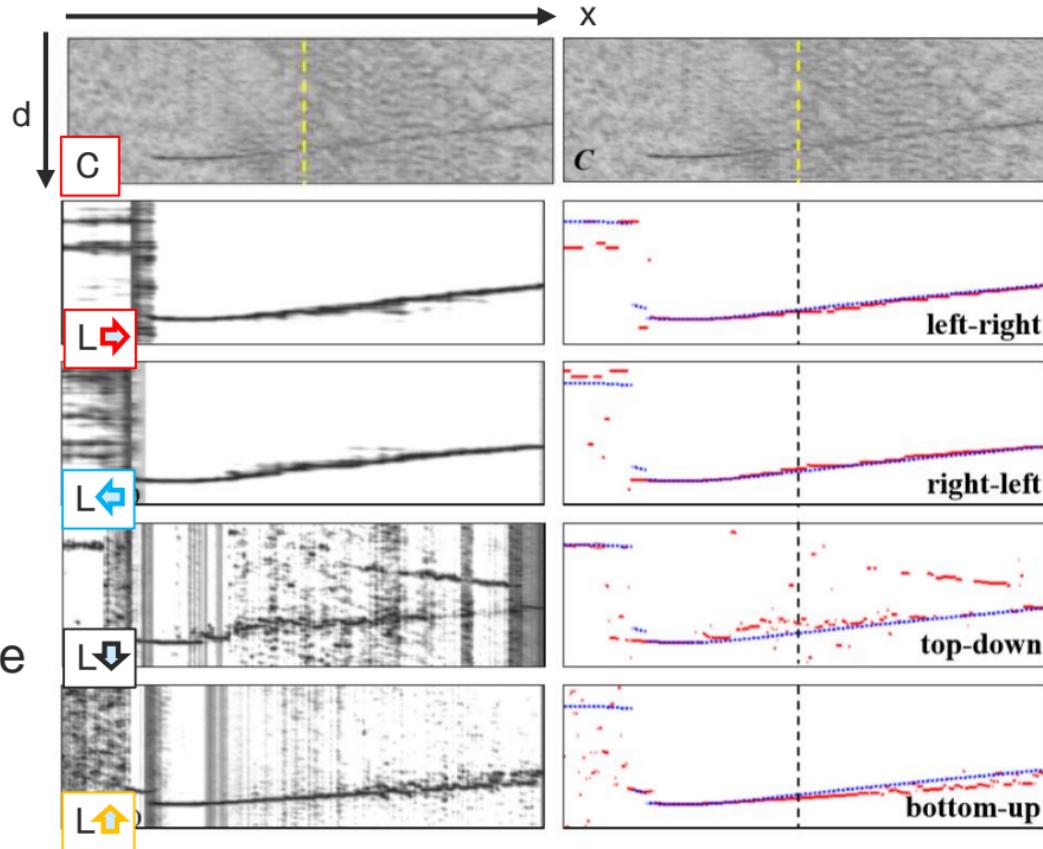
- C: unary cost volume
- L: scanline opt. cost volume
- Horizontal slices for the green scanline

Analyzing the Scanline Optimization Costs

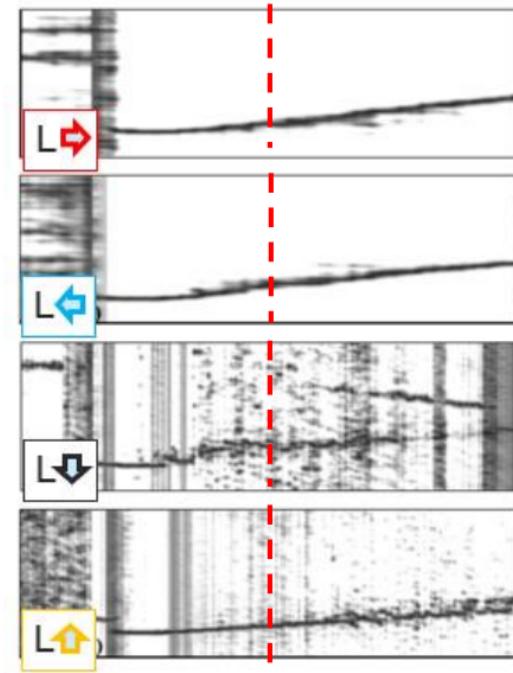
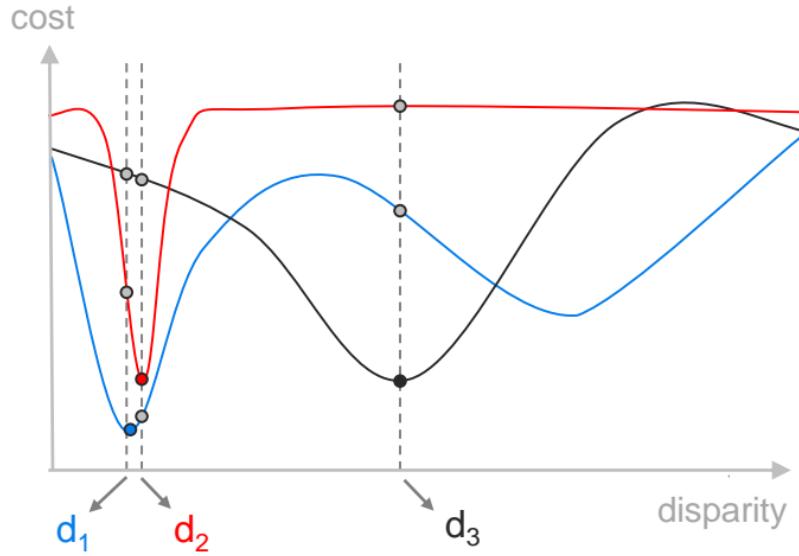
DSI



- C: unary cost volume
- L: scanline opt. cost volume
- Horizontal slices for the green scanline



Proposed Method



Feature vector:

$$f = [d_1, d_2, d_3, c_{11}, c_{12}, c_{13}, c_{21}, c_{22}, c_{23}, c_{31}, c_{32}, c_{33}]$$

Proposed Method

1. Run 1D Scanline Opt. to obtain k disparity map proposals
2. At each pixel
 - Construct $k^2 + k$ dimensional feature vector
 - Do Random Forest Inference
 - outputs probabilities for selecting k proposals
 - Select best disparity
 - Refine disparity
 - Compute probability-weighted average of inlier disparities.
3. Post-processing (using probability image)
 - At each pixel p , find neighboring pixels with high probability and color similar to that of p ; then compute the median filter.

Results

Ablation study: Middlebury 2014 training set (15 pairs) used for evaluation

Method	Left View Scanlines	Right View Scanlines	Filtering	Training Dataset	bad 0.5px [%]	bad 1px [%]	bad 2px [%]	bad 4px [%]	Time [s]
all									
SGM	all		–	65.58	36.08	20.66	16.24	3.0	
SGM – $\min_d L_r(\mathbf{p}, d)$	all		–	66.79	38.35	23.32	18.36	3.1	
SGM – $\min_d \text{median}_r L_r(\mathbf{p}, d)$	all		–	67.53	39.75	23.34	18.12	3.2	
SGM-SVM	all		M	60.89	32.59	20.31	16.16	323.7	
SGM-MLP	all		M	60.49	32.61	20.25	16.14	21.0	
SGM-Forest	horiz+vert		M	61.09	32.69	18.02	13.19	5.7	
	top-down		M	61.31	32.85	18.31	13.37	5.8	
	bottom-up		M	61.38	32.91	18.42	13.43	5.8	
	all		M	60.28	32.15	17.90	13.14	6.1	
	all	•	M	60.18	32.08	17.69	12.91	6.3	
	all	•	•	E	59.89	30.69	16.78	11.67	8.2
	all	•	•	K	59.70	30.61	16.72	11.67	8.2
	all	•	•	M	59.20	30.58	16.57	11.62	8.2

M: Midd 2005, K: KITTI, E: ETH3D

Results

Ablation study: Middlebury 2014 training set (15 pairs) used for evaluation

- Good cross-domain generalization; low dataset bias
- Random forest inference: fast (on CPU); parallelizable
- SGM Forest retains computational benefits of SGM

SGM – $\min_d L_r(\mathbf{p}, a)$	all	–	66.79	38.55	23.32	18.36	5.1
SGM – $\min_d \text{median}_r L_r(\mathbf{p}, d)$	all	–	67.53	39.75	23.34	18.12	3.2
SGM-SVM	all	M	60.89	32.59	20.31	16.16	323.7
SGM-MLP	all	M	60.49	32.61	20.25	16.14	21.0
SGM-Forest	horiz+vert	M	61.09	32.69	18.02	13.19	5.7
	top-down	M	61.31	32.85	18.31	13.37	5.8
	bottom-up	M	61.38	32.91	18.42	13.43	5.8
	all	M	60.28	32.15	17.90	13.14	6.1
	all	M	60.18	32.08	17.69	12.91	6.3
	all	•	59.89	30.69	16.78	11.67	8.2
	all	•	59.70	30.61	16.72	11.67	8.2
	all	•	59.20	30.58	16.57	11.62	8.2
	all	•	59.20	30.58	16.57	11.62	8.2

M: Midd 2005, K: KITTI, E: ETH3D

Results

Datacost	Method	Middlebury 2014				KITTI 2015				ETH3D 2017			
		0.5px	1px	2px	4px	0.5px	1px	2px	4px	0.5px	1px	2px	4px
all													
NCC	SGM	69.23	42.36	27.96	22.25	60.59	33.79	15.06	8.34	32.52	16.71	10.66	7.69
	SGM-F.	64.00	37.22	22.85	17.09	52.39	25.80	10.11	4.69	22.48	11.26	6.36	4.35
MC-CNN-fast	SGM	65.82	36.22	21.98	17.47	58.48	31.39	13.30	7.02	26.34	10.50	6.13	4.52
	SGM-F.	62.04	32.96	18.22	13.16	51.03	24.05	8.73	3.78	17.62	7.17	3.66	2.51
MC-CNN-acrt	SGM	65.58	36.08	20.66	16.24	57.24	28.55	9.54	5.26	39.03	16.34	9.14	6.67
	SGM-F.	59.20	30.58	16.57	11.62	46.88	19.77	6.51	2.97	27.40	11.89	7.30	5.52

Results

Datacost	Method	Middlebury 2014				KITTI 2015				ETH3D 2017			
		0.5px	1px	2px	4px	0.5px	1px	2px	4px	0.5px	1px	2px	4px
all													
NCC	SGM	69.23	42.36	27.96	22.25	60.59	33.79	15.06	8.34	32.52	16.71	10.66	7.69
	SGM-F.	64.00	37.22	22.85	17.09	52.39	25.80	10.11	4.69	22.48	11.26	6.36	4.35
MC-CNN-fast	SGM	65.82	36.22	21.98	17.47	58.48	31.39	13.30	7.02	26.34	10.50	6.13	4.52
	SGM-F.	62.04	32.96	18.22	13.16	51.03	24.05	8.73	3.78	17.62	7.17	3.66	2.51
MC-CNN-acrt	SGM	65.58	36.08	20.66	16.24	57.24	28.55	9.54	5.26	39.03	16.34	9.14	6.67
	SGM-F.	59.20	30.58	16.57	11.62	46.88	19.77	6.51	2.97	27.40	11.89	7.30	5.52

[Zbontar and Lecun 2015]

- Matching cost (unary): NCC, MC-CNN (-fast, -acrt)
- SGM-Forest outperforms SGM in all cases

Benchmark Results

Middlebury 2014 (MC-CNN-acrt)				
Method	non-occl.	all	Time	
LocalExp	5.43%	#1	11.7%	#1
3DMST	5.92%	#2	12.5%	#3
MC-CNN+TDSR	6.35%	#2	12.1%	#3
PMSC	6.71%	#4	13.6%	#4
LW-CNN	7.04%	#5	17.8%	#15
MeshStereoExt	7.08%	#6	15.7%	#9
FEN-D2DRR	7.23%	#7	16.0%	#11
APAP-Stereo	7.26%	#8	13.7%	#5
SGM-Forest	7.37%	#9	15.5%	#8
NTDE	7.44%	#10	15.3%	#7
				88s

Middlebury 2014 (MC-CNN-fast)				
Method	non-occl.	all	Time	
LocalExp	6.52 %	#1	12.1%	#1
3DMST	7.08 %	#2	12.9%	#2
APAP-Stereo	7.53%	#3	14.3%	#6
FEN-D2DRR	7.89%	#4	14.1%	#4
...				73s
MC-CNN-acrt	10.1%	#12	19.7%	#20
...				106s
SGM-Forest	11.1%	#19	17.8%	#14
...				9s
MC-CNN-fast	11.7%	#21	21.5%	#27
				1s

KITTI 2015		
Method	Error	Time
CNNF+SGM	3.60% (#9)	71.0s
SGM-Net	3.66% (#11)	67.0s
MC-CNN-acrt	3.89% (#12)	67.0s
➡ SGM-Forest	4.38% (#14)	6.0s
MC-CNN-WS	4.97% (#18)	1.4s
➡ SGM_ROB [2]	6.38% (#27)	0.1s
SGM+C+NL	6.84% (#31)	270.0s
SGM+LDOF	6.84% (#32)	86.0s
SGM+SF	6.84% (#33)	2700.0s
➡ CSCT+SGM+MF	8.24% (#35)	6.4ms

ETH3D 2017			
Method	non-occl.	all	Time
➡ SGM-Forest	5.40%	4.96%	5.21s
➡ SGM_ROB [2]	10.08%	10.77%	0.15s
MeshStereo	11.94%	11.52%	159.24s
SPS-Stereo	15.83%	15.04%	1.59s
ELAS	17.99%	16.72%	0.13s

Benchmark Results

Middlebury 2014 (MC-CNN-acrt)				
Method	non-occl.	all	Time	
LocalExp	5.43%	#1	11.7%	#1
3DMST	5.92%	#2	12.5%	#3
MC-CNN+TDSR	6.35%	#2	12.1%	#3

Middlebury 2014 (MC-CNN-fast)				
Method	non-occl.	all	Time	
LocalExp	6.52 %	#1	12.1%	#1
3DMST	7.08 %	#2	12.9%	#2
APAP-Stereo	7.53%	#3	14.3%	#6

- #1 on ETH3D
- #8 on Midd 2014, #14 on KITTI
- Significantly outperform all previous SGM variants
- SGM-ROB: [Hirschmuller 2008] impl.

↳ SGM-Forest	4.38% (#14)	6.0s
MC-CNN-WS	4.97% (#18)	1.4s
↳ SGM_ROB [2]	6.38% (#27)	0.1s
SGM+C+NL	6.84% (#31)	270.0s
SGM+LDOF	6.84% (#32)	86.0s
SGM+SF	6.84% (#33)	2700.0s
↳ CSCT+SGM+MF	8.24% (#35)	6.4ms

↳ SGM-FOREST	5.40%	4.96%	5.21s
↳ SGM_ROB [2]	10.08%	10.77%	0.15s
MeshStereo	11.94%	11.52%	159.24s
SPS-Stereo	15.83%	15.04%	1.59s
ELAS	17.99%	16.72%	0.13s

Outline

- Stereo Matching: New Trends
- Semi-Global Stereo Matching (SGM)
 - SGM with Surface Orientation Priors
 - Learning to Fuse Proposals
- **Stereo Scene Flow with Motion Segmentation**
- Camera Trajectory Planning for Aerial 3D Scanning
- Deep 6-DoF Object Pose Prediction

Fast Multi-frame Stereo Scene Flow with Motion Segmentation



Tatsunori Taniai
RIKEN, Tokyo



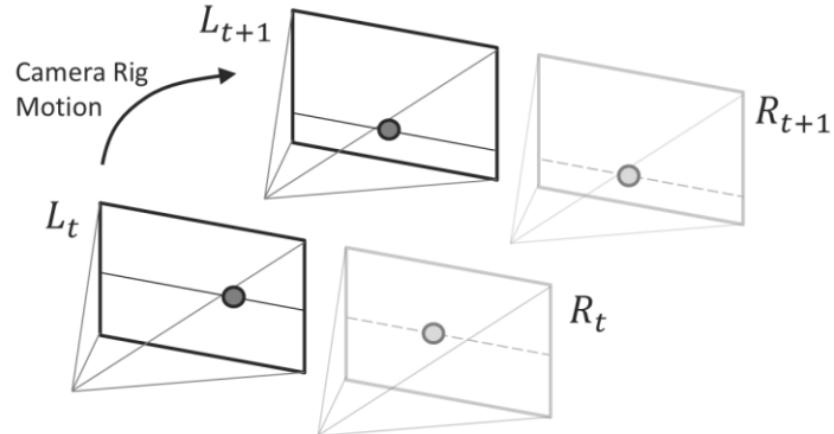
Sudipta Sinha
Microsoft Research



Yoicho Sato
Univ. of Tokyo

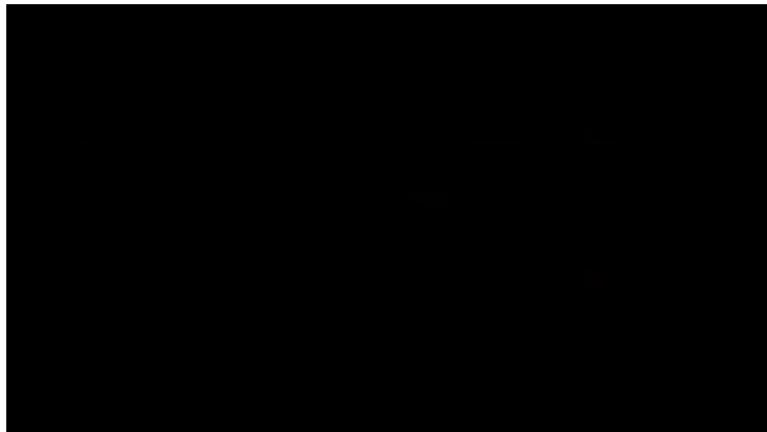
CVPR 2017

Multi-frame Stereo Scene Flow



- Stereo video from moving stereo camera rig (calibrated)
- Scene Flow equivalent to stereo matching and optical flow estimation

Application



Object scene flow for autonomous vehicles

Menze and Geiger 2015

Action recognition by dense trajectories

Wang+ 2011

- Depth and flow sequences are useful in many applications

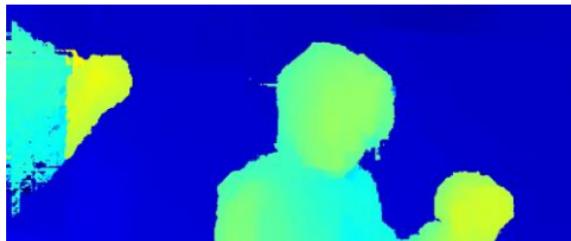
Motivation

Efficient, unified method for

- Stereo
- Optical Flow
- Moving object segmentation
- Visual Odometry (Camera ego-motion)



Disparity Map



Optical Flow



Moving Object Segmentation



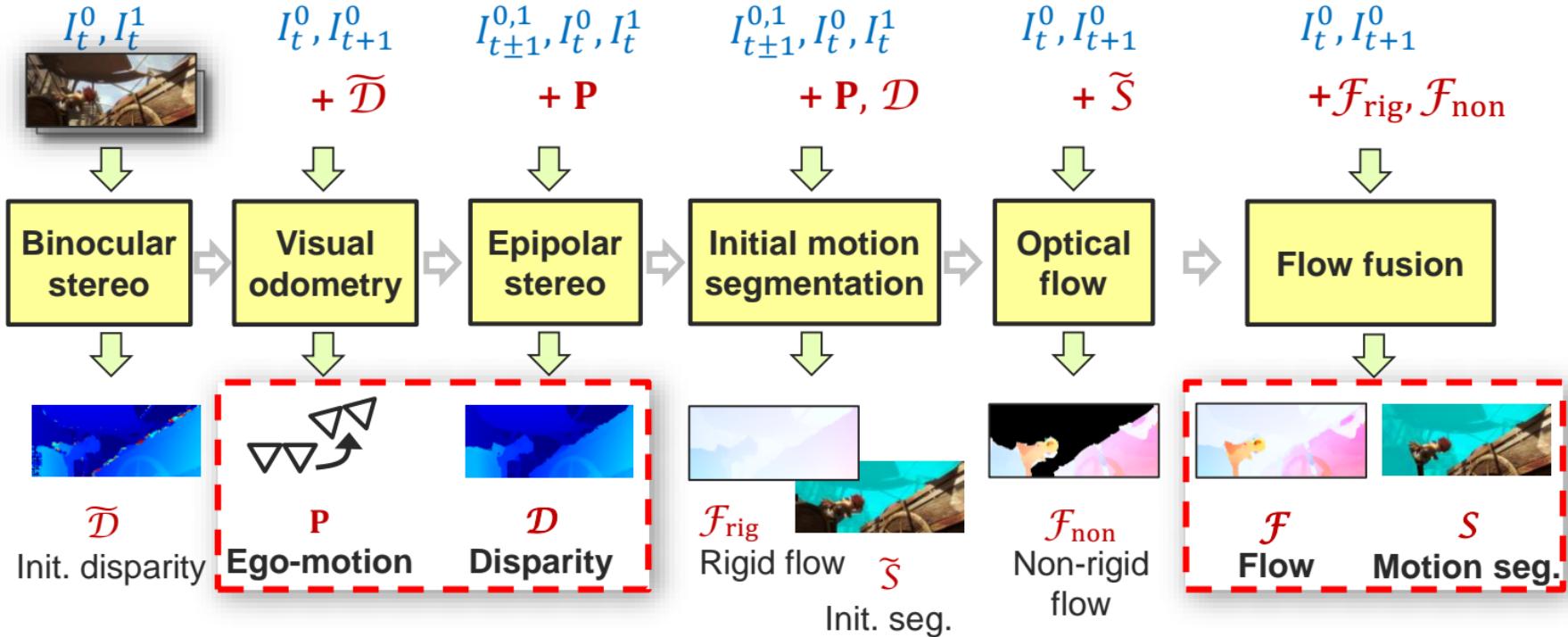
Main Idea: Dominant Rigid Scene Assumption



- Most of the scene is rigid; hence, camera motion determines the *rigid optical flow*.
- Given *rigid flow map*, only find regions with moving objects and recompute their flow.

Proposed Approach

Input



Results – KITTI 2015 Scene Flow Benchmark (Nov 2016)

Rank	Method	D1-bg	D1-fg	D1-all	D2-bg	D2-fg	D2-all	Fl-bg	Fl-fg	Fl-all	SF-bg	SF-fg	SF-all	Time
1	PRSM [43]	3.02	10.52	4.27	5.13	15.11	6.79	5.33	17.02	7.28	6.61	23.60	9.44	300 s
2	OSF [30]	4.54	12.03	5.79	5.45	19.41	7.77	5.62	22.17	8.37	7.01	28.76	10.63	50 min
3	FSF+MS (ours)	5.72	11.84	6.74	7.57	21.28	9.85	8.48	29.62	12.00	11.17	37.40	15.54	2.7 s
4	CSF [28]	4.57	13.04	5.98	7.92	20.76	10.06	10.40	30.33	13.71	12.21	36.97	16.33	80 s
5	PR-Sceneflow [42]	4.74	13.74	6.24	11.14	20.47	12.69	11.73	27.73	14.39	13.49	33.72	16.85	150 s
8	PCOF + ACTF [10]	6.31	19.24	8.46	19.15	36.27	22.00	14.89	62.42	22.80	25.77	69.35	33.02	0.08 s (GPU)
12	GCSF [8]	11.64	27.11	14.21	32.94	35.77	33.41	47.38	45.08	47.00	52.92	59.11	53.95	2.4 s

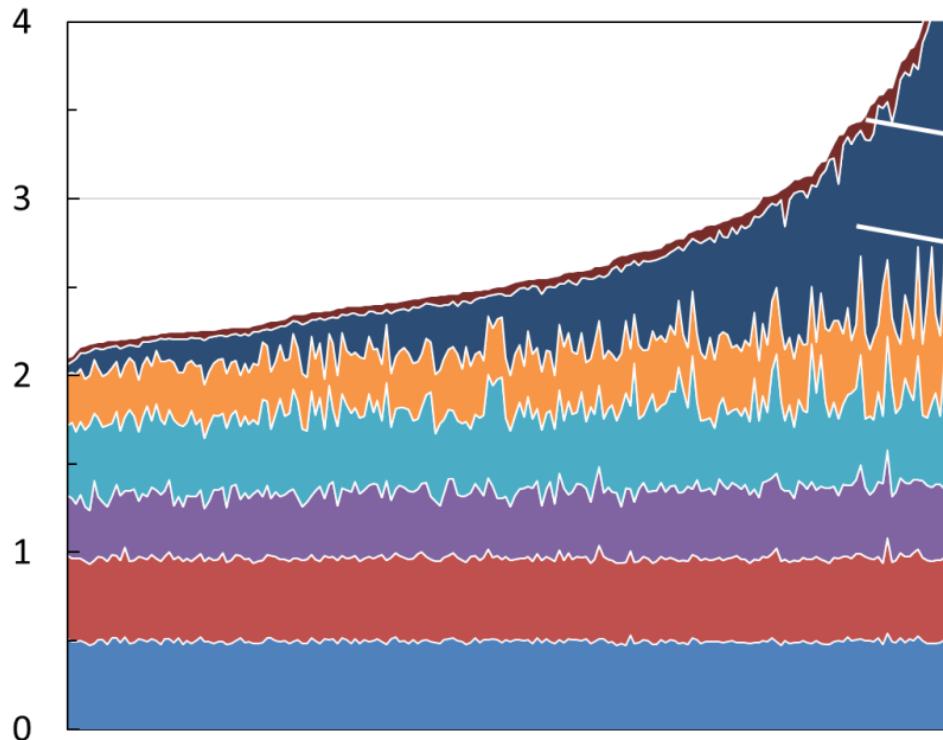


200 road scenes with multiple moving objects

Rank	SF-all	Time
1	9.44	300 s
2	10.63	50 min
3	15.54	2.7 s
4	16.33	80 s
5	16.85	150 s
8	33.02	0.08 s (GPU)
12	53.95	2.4 s

Results – KITTI 2015 Scene Flow Benchmark (Nov 2016)

Running time / frame (sec)



200 scenes from KITTI benchmark

CPU: 3.5 GHz × 4 Cores
Image: (1242 × 375) × 0.65 scale

- 0.07** sec Flow fusion
- 0.48** sec Optical flow
- 0.36** sec Initial segmentation
- 0.47** sec Epipolar stereo
- 0.38** sec Visual odometry
- 0.47** sec Binocular stereo
- 0.72** sec Initialization

2.72 sec per frame

Outline

- Stereo Matching: New Trends
- Semi-Global Stereo Matching (SGM)
 - SGM with Surface Orientation Priors
- Stereo Scene Flow with Motion Segmentation
- Camera Trajectory Planning for Aerial 3D Scanning
- Deep 6-DoF Object Pose Prediction

Submodular Trajectory Optimization for Aerial 3D Scanning

ICCV 2017

Mike Roberts^{1,2} Debadatta Dey² Anh Truong³ Sudipta Sinha²
Shital Shah² Ashish Kapoor² Pat Hanrahan¹ Neel Joshi²

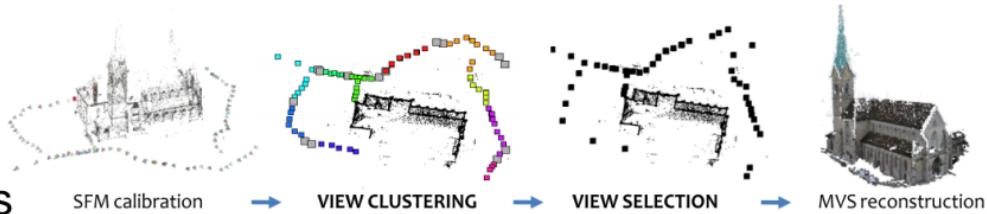
¹Stanford University

²Microsoft Research

³Adobe Research

Related Work

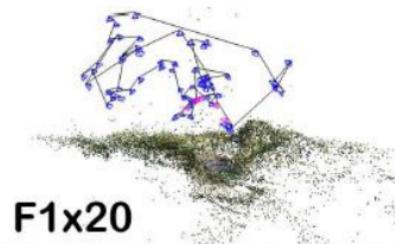
- View selection [Hornung+ 2011]
 - First, acquire dense imagery
 - Later, select subset & process



- Next-best-view planning
 - Information-gain maximization [Isler+ 2016]
 - Robotic RGB-D 3D scanning [Wu+ 2014]
 - No travel budget constraints



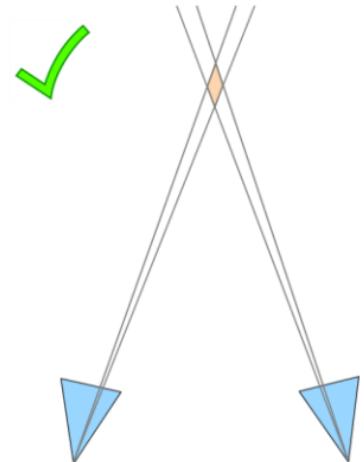
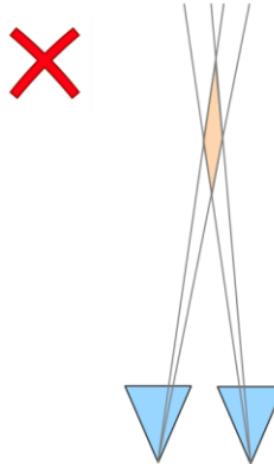
- Real-time Drone view planning
[Mostegel+ 2016, Hepp+ 2018]
 - Greedy technique; heuristic-based



Planning the camera path

1. Fly an easy-to-generate trajectory
2. Compute coarse 3D model (SFM → MVS → mesh)
- 3. Plan *optimized* trajectory**
4. Fly the computed trajectory
5. Run SFM + MVS on all images

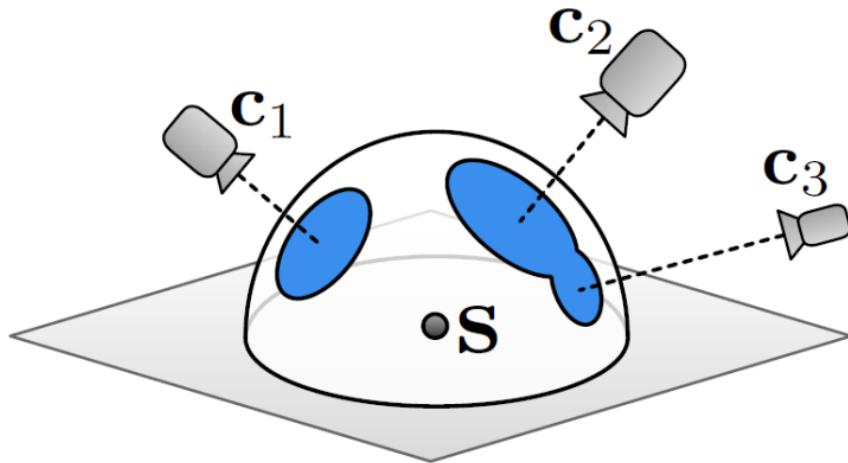
Heuristics Known to Work Well



We prefer images with

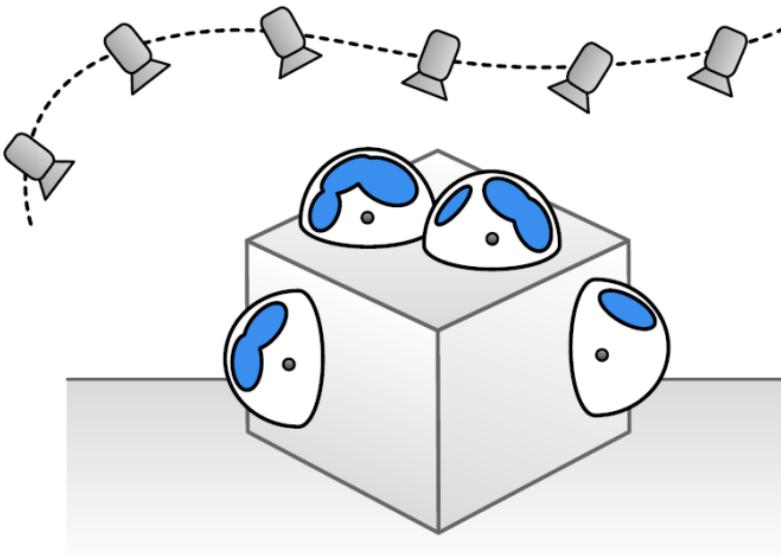
- diverging viewing angles
- close-up views
- fronto parallel views of surfaces

Coverage Measure



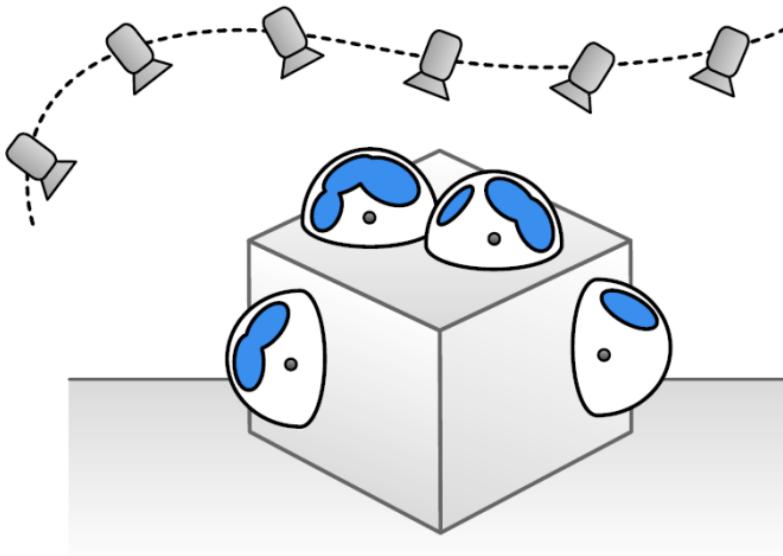
For a surface point S observed from multiple cameras, we define coverage as the area of the union of all the blue disks on a hemisphere.

Coverage Measure



Similarly, we define coverage for multiple surface points observed from multiple camera viewpoints.

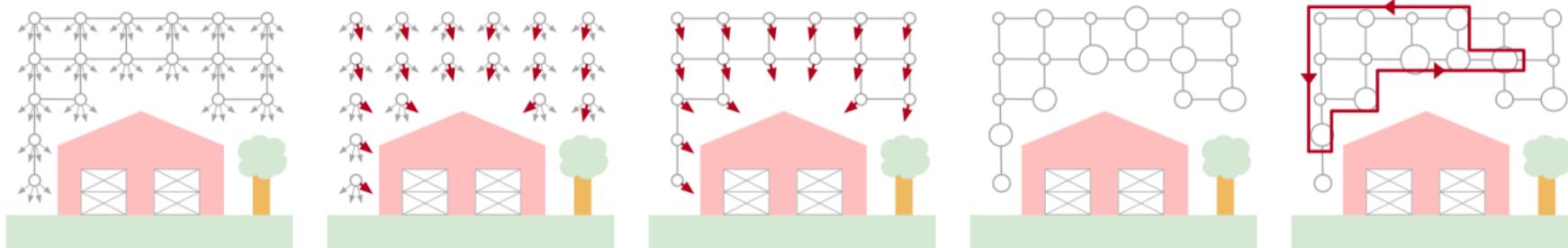
Coverage Measure



Submodular
Function!

Similarly, we define coverage for multiple surface points observed from multiple camera viewpoints.

Planning Optimized Trajectories



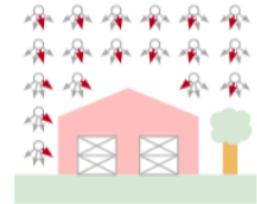
Nodes = possible camera locations; edge weights = pairwise Euclidean distances.

Sub-problems:

1. Solve optimal set of orientations; ignoring path constraints
2. Then, find path by solving a *graph orienteering problem*

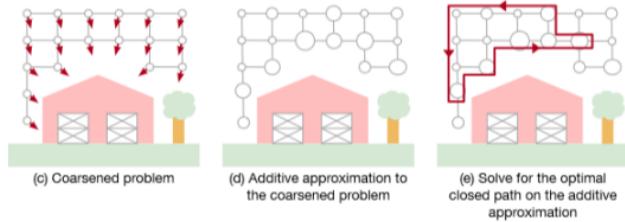
Solving for Camera Orientations

- Coverage set function is submodular
 - Adding new elements to an existing set gives diminishing returns
- Cardinality and Mutual Exclusion Constraint
 - Select exactly one look-at vector at each position
- Constrained submodular maximization
 - Always, pick the next best element with the most marginal reward
 - Greedy algorithm; good theoretical approximation guarantee



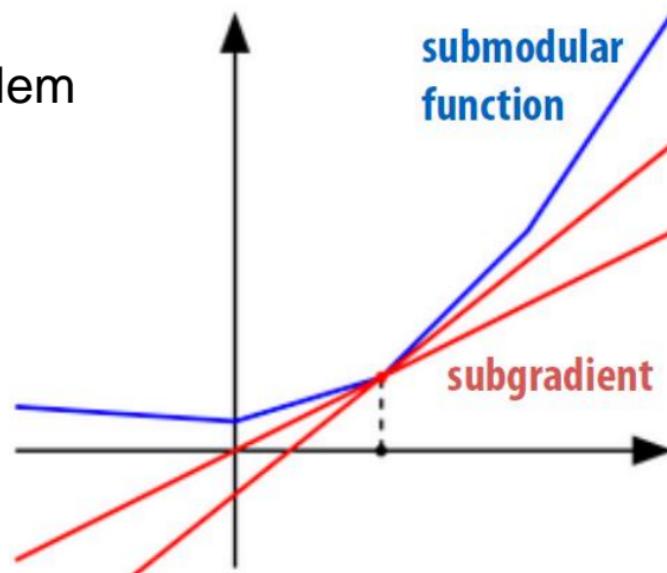
Solving for Camera Positions

- Graph Orienteering Problem
 - NP-Hard; related to TSP and Knapsack
 - Find short paths that let you collect most rewards (at nodes).
- In standard orienteering, rewards are additive.
- But, our reward function is submodular, not additive!
- Hence, we must solve a *submodular orienteering problem*.



Solving for Camera Positions

- Choose a good sub-gradient (additive approximation) for our submodular function
- This gives us a regular orienteering problem
- Solve a integer linear program (ILP)



Outline

- Stereo Matching: New Trends
- Semi-Global Stereo Matching (SGM)
 - SGM with Surface Orientation Priors
- Stereo Scene Flow with Motion Segmentation
- Camera Trajectory Planning for Aerial 3D Scanning
- Deep 6-DoF Object Pose Prediction

Real-Time Seamless Single Shot 6D Object Pose Prediction



Bugra Tekin

EPFL



Sudipta Sinha

Microsoft Research



Pascal Fua

EPFL

CVPR 2018

3D Recognition, 2D-3D Model Alignment

Given a RGB image (with known intrinsics), recognize the object instances and predict their 3D position and orientation within the scene.

Feature-based (RGB, RGB-D)

- Lowe 2001
- Rothganger+ 2005
- Lepetit+ 2005
- Lai+ 2010
- Hinterstoisser+ 2012

Coordinate Regression

- Brachmann+ 2014, 2016

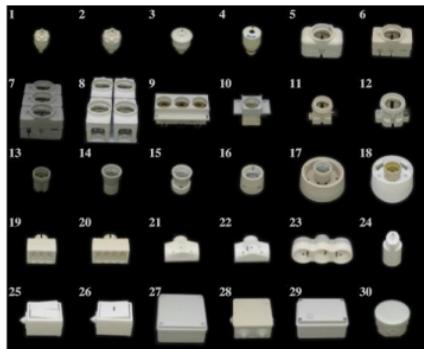
CNNs

- Rad + Lepetit 2017
- Kehl+ 2017
- Xiang+ 2017
- Tekin+ 2018
- Oberweger+ 2018

Texture-less Object 6D Pose Datasets



LINEKIT [2012]
15 objects



T-LESS [2017]
30 objects



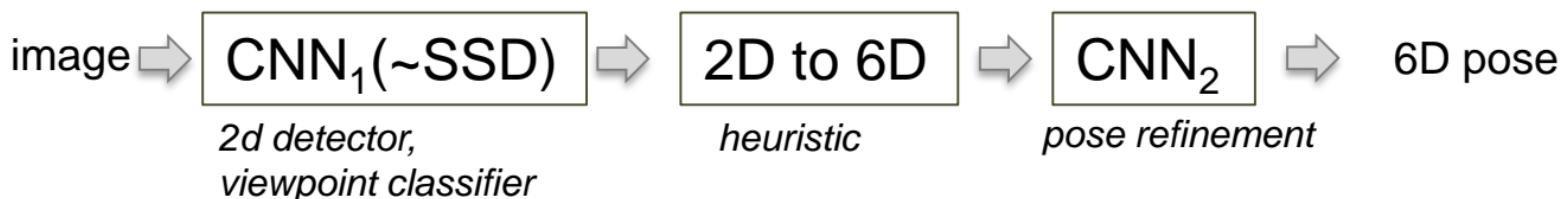
YCB-VIDEO [2018]
21 objects

Deep 6-dof object pose estimation

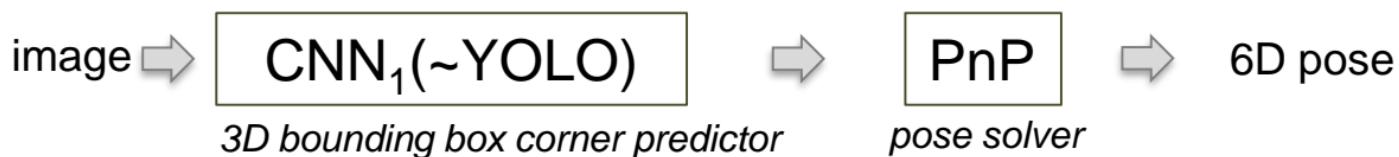
- BB8 [Rad and Lepetit 2017]



- SSD-6D [Kehl+ 2017]

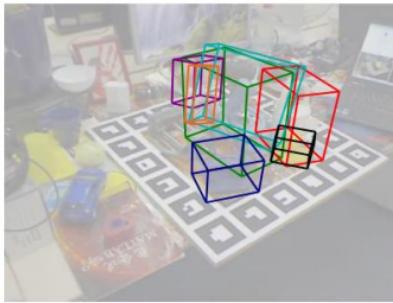


- Ours



Our Method

- Single-shot 2D object detection (YOLO, SSD)
- Our CNN predicts 2D projections of 3D bounding box vertices (and the centroid). We run PnP solver on 9 2D-3D correspondences.
- Accurate, fast (50-90 fps); detects multiple objects in one pass.



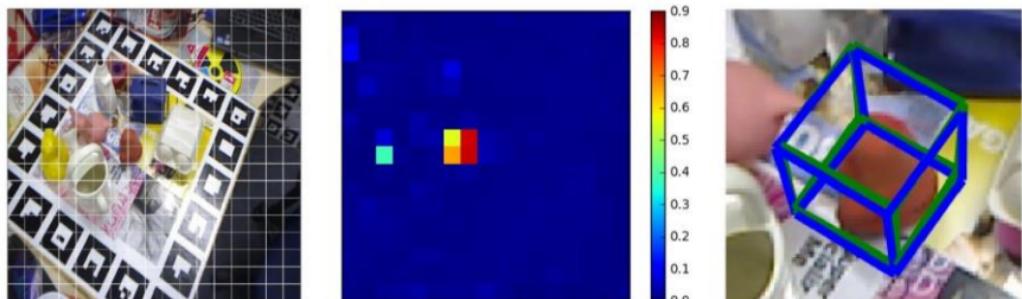
Our Method

Training:

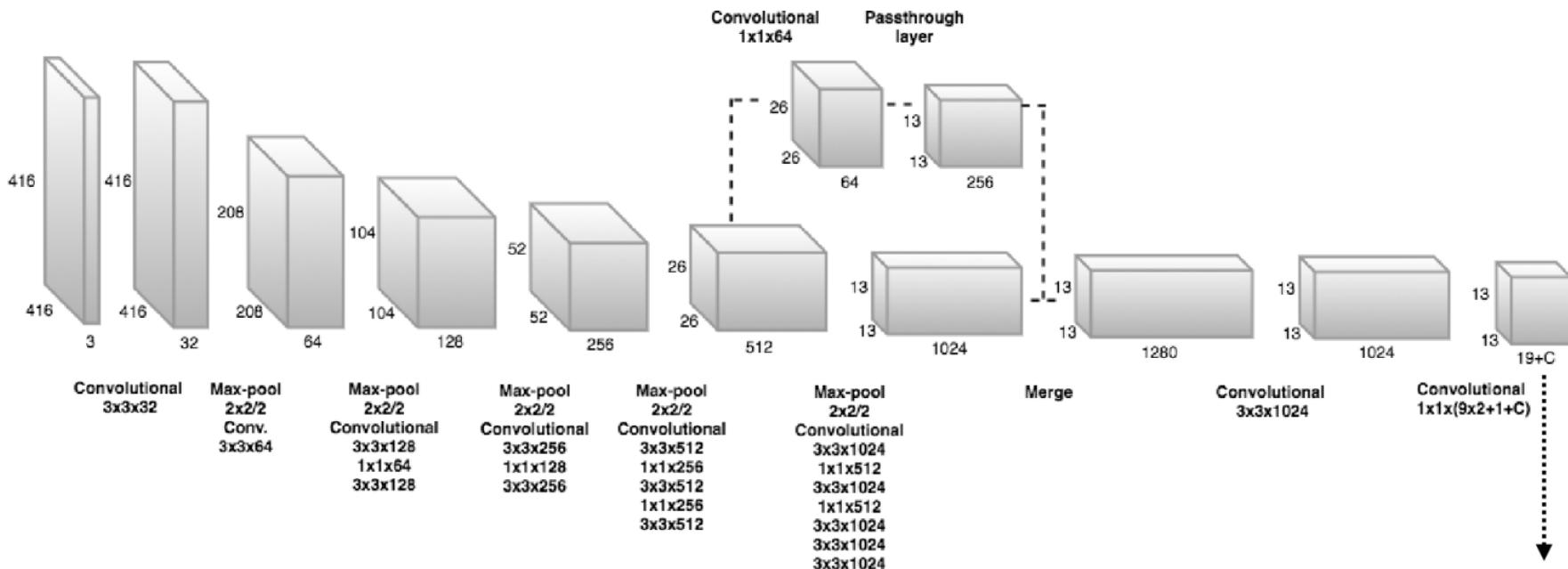
- ground truth 2D coordinates of the 9 control points are the targets
- modify YOLO loss function (for confidence estimation)
- data augmentation

Testing:

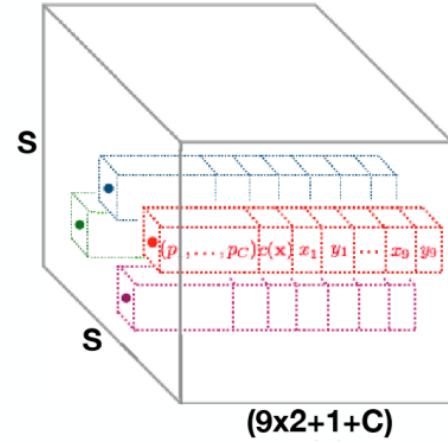
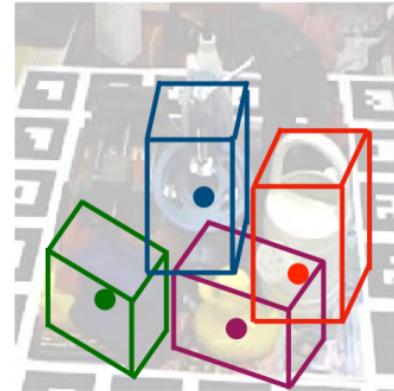
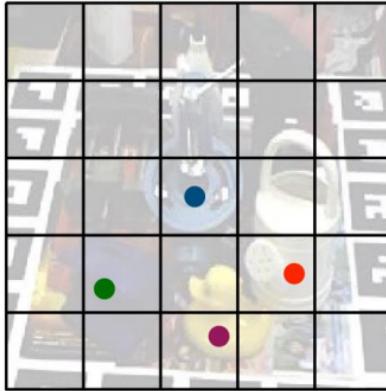
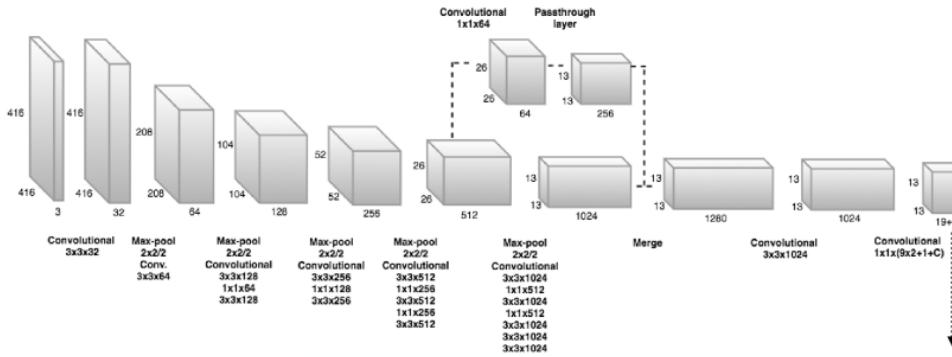
- Subpixel refinement
- PnP (RANSAC, least squares)



CNN Architecture



CNN Architecture



Results on LineMOD dataset

- Two accuracy metrics (2D image projection, 3D model overlap).
- Percentage of test images where the error was lower than specified thresholds.

2D metric

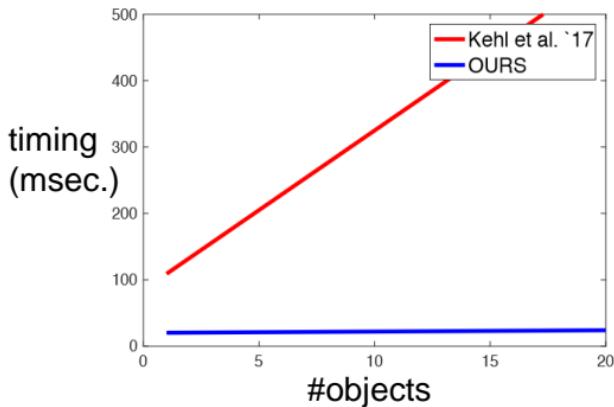
Method	w/o Refinement			w/ Refinement	
	Brachmann [2]	BB8 [25]	OURS	Brachmann [2]	BB8 [25]
Object					
Ape	-	95.3	92.10	85.2	96.6
Benchvise	-	80.0	95.06	67.9	90.1
Cam	-	80.9	93.24	58.7	86.0
Can	-	84.1	97.44	70.8	91.2
Cat	-	97.0	97.41	84.2	98.8
Driller	-	74.1	79.41	73.9	80.9
Duck	-	81.2	94.65	73.1	92.2
Eggbox	-	87.9	90.33	83.1	91.0
Glue	-	89.0	96.53	74.2	92.3
Holepuncher	-	90.5	92.86	78.9	95.3
Iron	-	78.9	82.94	83.6	84.8
Lamp	-	74.4	76.87	64.0	75.8
Phone	-	77.6	86.07	60.6	85.3
Average	69.5	83.9	90.37	73.7	89.3

3D metric

Method	w/o Refinement				w/ Refinement		
	Brachmann [2]	BB8 [25]	SSD-6D [10]	OURS	Brachmann [2]	BB8 [25]	SSD-6D [10]
Object							
Ape	-	27.9	0	21.62	33.2	40.4	65
Benchvise	-	62.0	0.18	81.80	64.8	91.8	80
Cam	-	40.1	0.41	36.57	38.4	55.7	78
Can	-	48.1	1.35	68.80	62.9	64.1	86
Cat	-	45.2	0.51	41.82	42.7	62.6	70
Driller	-	58.6	2.58	63.51	61.9	74.4	73
Duck	-	32.8	0	27.23	30.2	44.3	66
Eggbox	-	40.0	8.9	69.58	49.9	57.8	100
Glue	-	27.0	0	80.02	31.2	41.2	100
Holepuncher	-	42.4	0.30	42.63	52.8	67.2	49
Iron	-	67.0	8.86	74.97	80.0	84.7	78
Lamp	-	39.9	8.20	71.11	67.0	76.5	73
Phone	-	35.2	0.18	47.74	38.1	54.0	79
Average	32.3	43.6	2.42	55.95	50.2	62.7	79

Results on LineMOD dataset

- Running Times:
 - On TitanX or similar GPU.
 - using cuDNN



Method	Overall speed for 1 object	Refinement runtime
Brachmann et al. [2]	2 fps	100 ms/object
Rad & Lepetit [25]	3 fps	21 ms/object
Kehl et al. [10]	10 fps	24 ms/object
OURS	50 fps	-

Method	2D projection metric	Speed
416×416	89.71	94 fps
480×480	90.00	67 fps
544×544	90.37	50 fps
688×688	90.65	43 fps



Up to 90 fps on smaller images

Conclusions

- New trends and challenges in stereo matching
- Improving Semi-Global Stereo (SGM) Matching
 - Incorporating orientation priors
 - Learned fusion step
- Fast scene flow with motion segmentation
- Camera path planning for improved multi-view stereo
- Deep single shot 6D object pose estimation

Collaborators



Rick Szeliski
Facebook



Adarsh Kowdle
Google



Daniel Scharstein
Middlebury College



Tatsunori Taniai
RIKEN, Tokyo



Yoichi Sato
Univ. of Tokyo



Bugra Tekin
EPFL



Pascal Fua
EPFL



Johannes
Schoenberger
ETH Zurich



Marc Pollefeys
Microsoft / ETH
Zurich