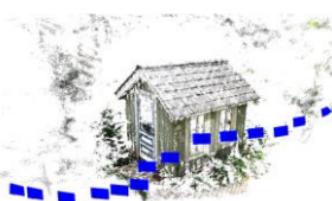


# **Efficient and Accurate 3D Scene Reconstruction and Object Pose Prediction**

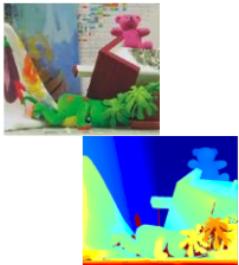
Sudipta Sinha  
Microsoft Research

*University of Illinois at Urbana-Champaign  
May 9, 2018*

# Overview



Structure from Motion



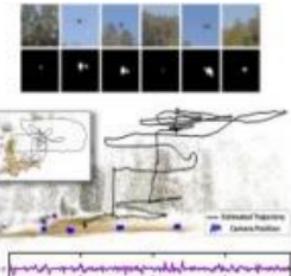
Dense Stereo



Photometric Stereo



Scene Flow



Multiview Tracking



Image-based Rendering



Aerial Mapping with UAVs

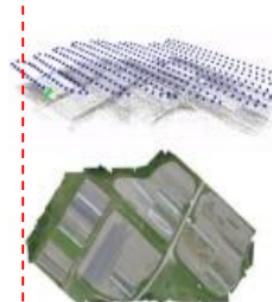


Image-based Camera Localization

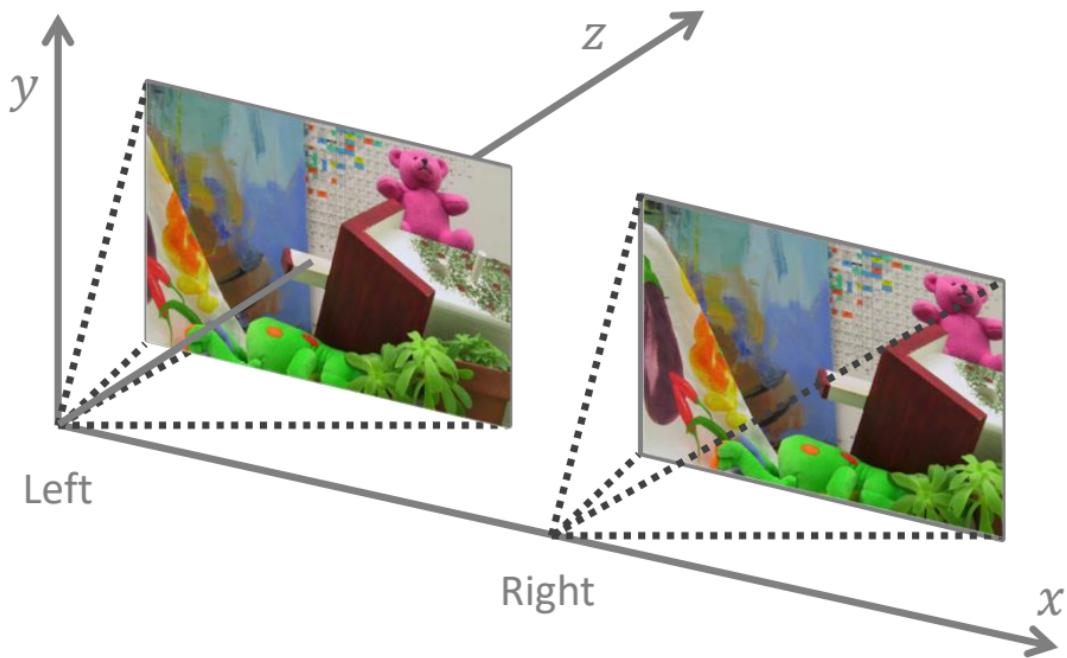


Object 6D Pose Augmented Reality

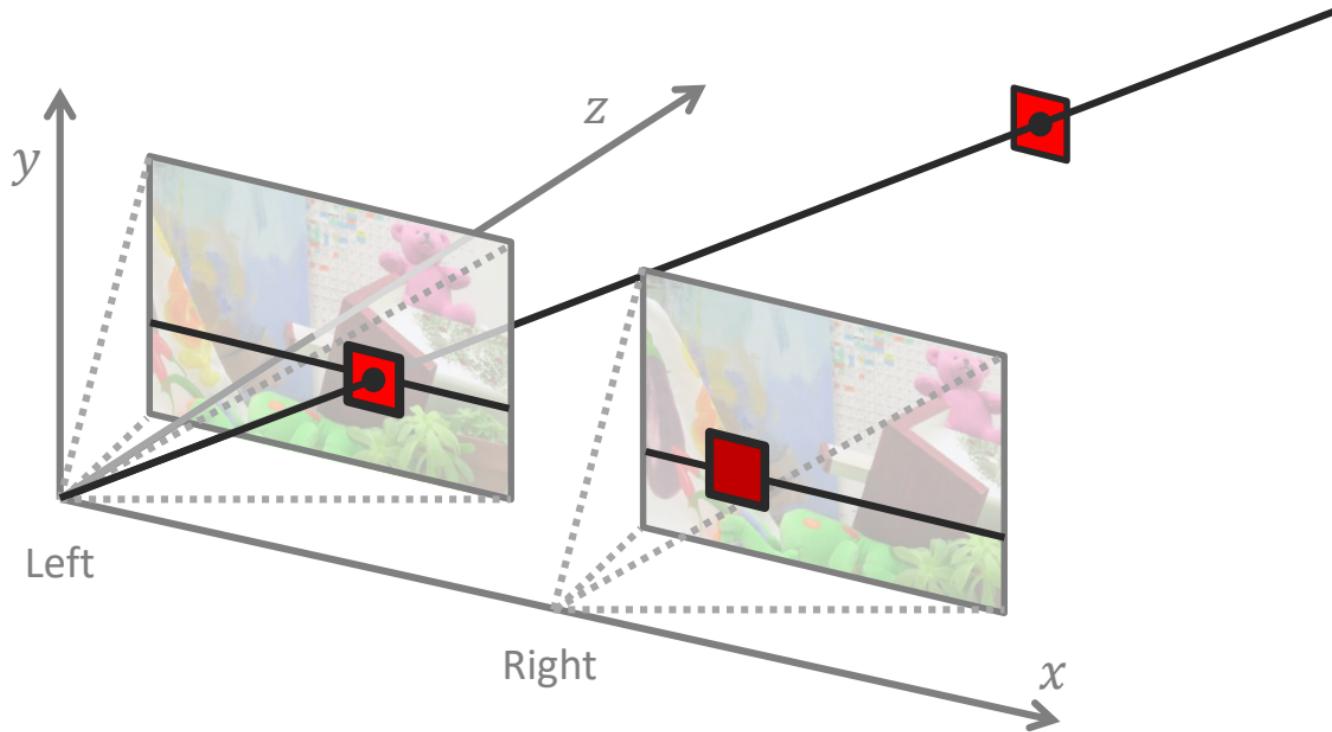
# Outline

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

# Stereo Matching



# Stereo Matching



# Still Challenging



Fore-shortening



Specular



Transparency, reflections



Different lighting



Dynamic range



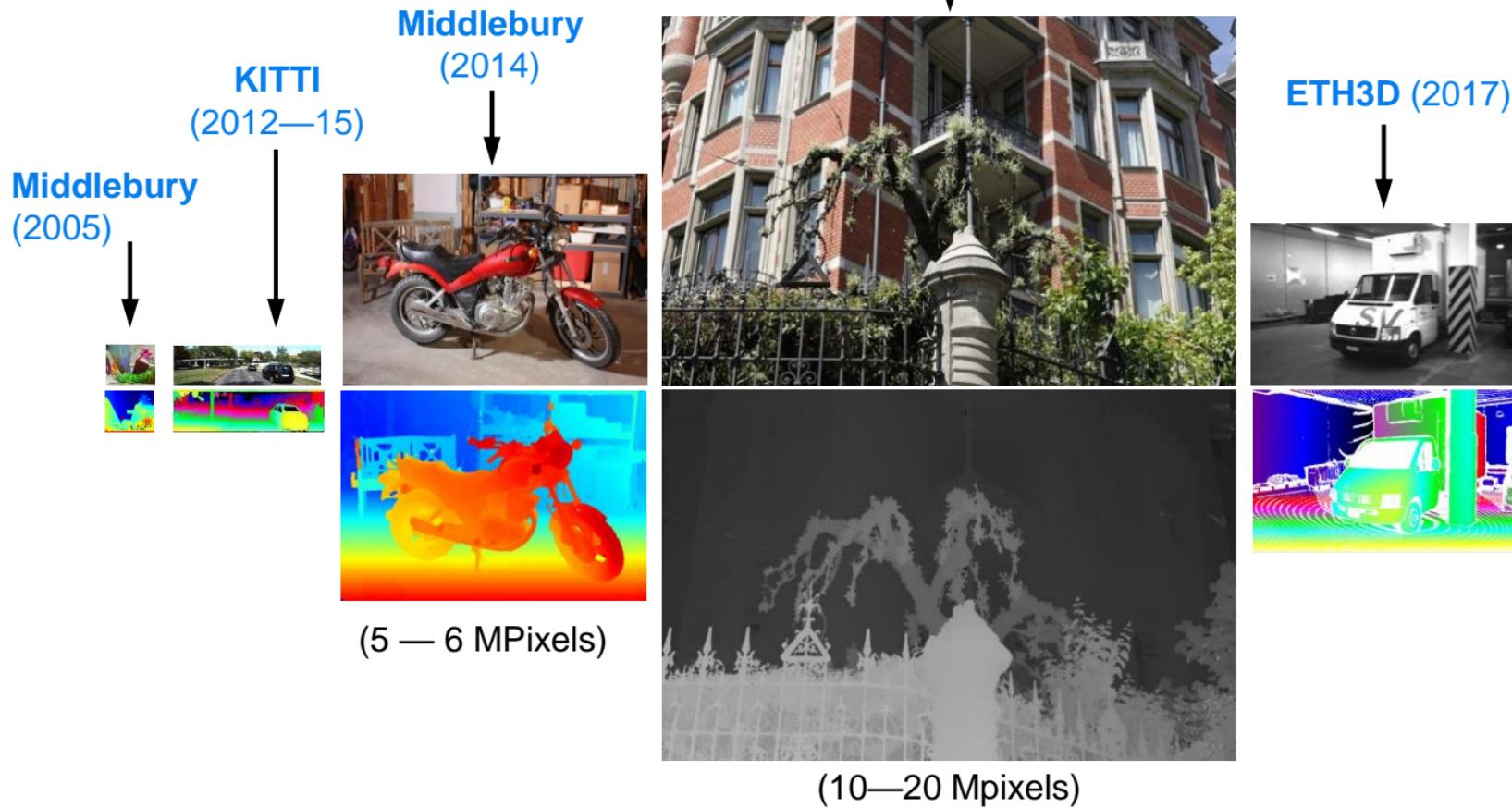
Untextured slanted surfaces

- State of the art methods are still ...

- Inaccurate in many corner cases
- Too slow for real-time, resource constrained systems

# Stereo benchmarks

Kim+ 2013



# Classical Methods (MRF inference)

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

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

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

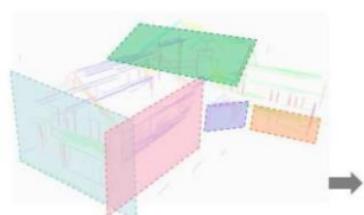
- Data (cost) term encodes matching costs
- Smoothness (cost) term encodes prior
- Discrete vs. Continuous labels
- Inference: Graph Cuts, Belief Propagation, PatchMatch-style optim.

# Piecewise Planar Stereo

[Sinha, Steedly, Szeliski 2009]



Structure from motion



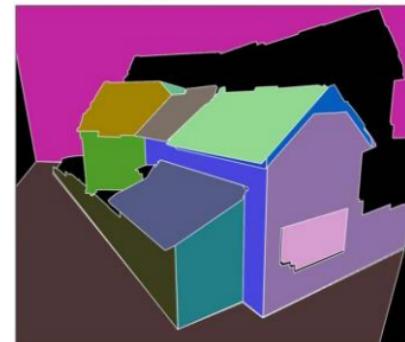
Multiple Plane Detection



3D Line Reconstruction



MRF energy minimization  
via graph cuts

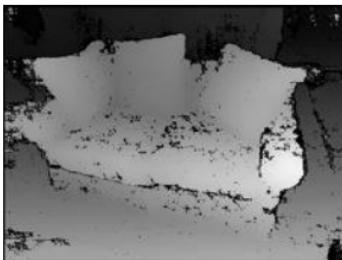


Novel View

# Piecewise Planar Stereo Revisited

- Local plane fitting (more flexible)
- CRF models photo-consistency (stereo cue) and color segmentation (monocular cue)
- Learn color models per-surface
- Alternate between graph cuts and learning

[Kowdle, Sinha, Szeliski 2012]



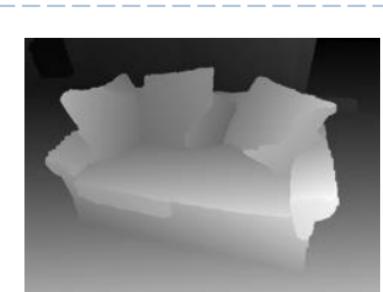
Semi-global stereo (SGM)



Find planes



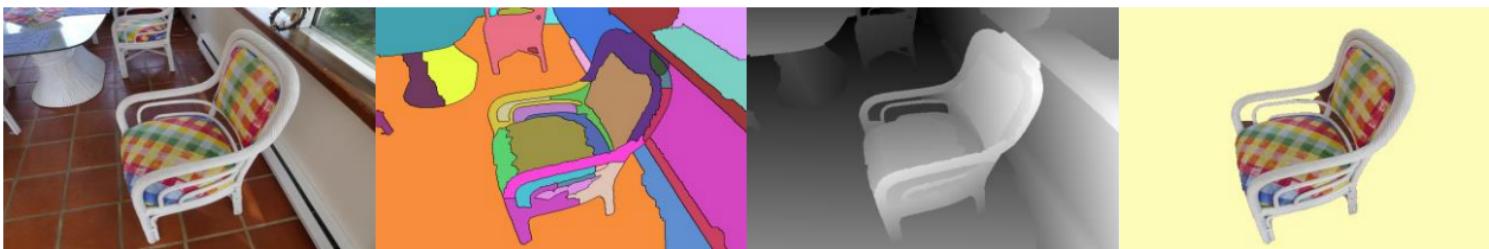
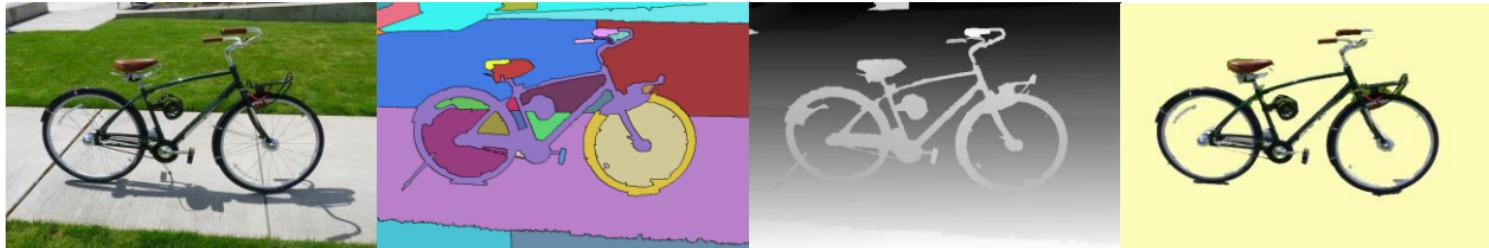
Label map



Depth map

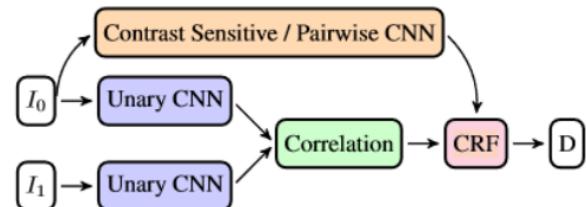
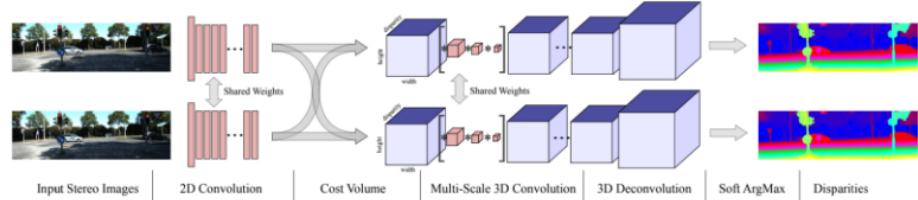
# Piecewise Planar Stereo Revisited

[Kowdle, Sinha, Szeliski 2012]



# New Trends

- Learning the matching cost:
  - MC-CNN [Zbontar + Lecun 2015], Chen+ 2015, Luo+ 2016
- Continuous MRFs: [Taniai+ 2017] (Rank 1 on Middlebury 2014!)
- Deep stereo regression (end to end training)
  - FlowNet [Dosovitskiy+ 2015], DispNet [Mayer+ 2016]
- Return of “Correlation”
  - DispNetCorr [Mayer+ 2016]
  - GC-Net [Kendall+ 2017]
- Return of “CRFs” (Hybrid CNN-CRF models)
  - Seki and Pollefeys 2017, Knobelreiter+ 2017



# Stereo Benchmark Rankings

## Middlebury 2014

Mouseover the table cells to see the produced disparity map. Clicking a cell will link 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](#) [0.5](#) [bad](#) [1.0](#) [bad](#) [2.0](#) [bad](#) [4.0](#) [avgerr](#) [rms](#) [A50](#) [A90](#) [A95](#) [A99](#) [time](#) [time/MP](#) [time/GD](#)

Mask: [none](#) [cc](#) [all](#)

[plot selected](#)

[show invalid](#)

[Reset sort](#)

[Reference list](#)

Date	Name	Res	Avg	Austr		Austri		Bulg2		Class		Class2		Crossa		Crossa		Dambi		Dambi		Dambi		Houses		Liverp		Nubua		Plants		Stairs	
				MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP			
09/09	09/09	09/09	09/09	0.54	0.54	2.84	3.99	1.93	5.15	3.34	3.22	3.15	2.32	8.55	2.22	7.49	7.09	12.5	5.20	10.0	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54			
06/22/15	LocoExp	06/22/15	06/22/15	3.65	3.65	2.87	2.88	1.99	5.59	2.37	3.43	3.55	2.06	1.03	0.73	8.57	14.4	5.40	9.55	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54				
01/24/17	3DMS	01/24/17	01/24/17	5.92	5.92	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	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54				
09/10/17	MC-CNN-TDSR	09/10/17	09/10/17	F	0.95	0.45	1.45	0.80	2.46	10.7	0.05	0.50	5.19	2.82	10.8	0.62	6.59	11.4	6.01	7.04	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54			
09/12/15	PMSC	09/12/15	09/12/15	H	0.71	3.46	1.68	0.19	2.54	6.02	4.54	3.96	4.04	2.37	13.1	12.3	12.7	16.7	5.88	10.8	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54			
10/19/15	LW-CNN	10/19/15	10/19/15	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	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54		
04/12/16	MeshStereosEx	04/12/16	04/12/16	H	7.06	4.41	3.98	1.40	3.17	10.7	6.23	4.82	4.77	3.49	12.7	12.4	10.4	14.5	7.89	8.85	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54			
10/12/17	FEN-D2DR	10/12/17	10/12/17	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	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54			
05/28/16	APAP-Stereo	05/28/16	05/28/16	H	7.26	6.43	4.91	2.31	5.11	5.17	21.31	6.99	4.31	4.23	3.24	14.3	9.78	7.32	13.2	4.60	4.69	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54		
03/11/18	SGM-Forest	03/11/18	03/11/18	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	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54			
01/19/16	NTDE	01/19/16	01/19/16	H	7.44	6.72	4.76	4.50	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	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54		
02/28/18	FDR	02/28/18	02/28/18	H	7.89	5.41	4.22	4.20	2.73	10.2	5.40	6.40	4.50	4.75	11.2	5.5	13.4	16.5	5.22	13.0	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54			
09/25/15	MC-CNN-act	09/25/15	09/25/15	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	11.1	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54			
11/03/15	MC-CNN-RBS	11/03/15	11/03/15	H	8.42	6.05	5.10	5.27	6.24	3.27	11.4	4.36	1.87	9.83	3.21	15.1	15.9	12.8	13.5	7.04	9.99	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54		
09/30/16	SNP-RSM	09/30/16	09/30/16	H	8.75	5.46	4.85	21.60	5.37	10.4	7.31	17.87	9.27	3.58	14.3	14.7	14.9	12.8	10.1	10.8	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54			
12/11/17	CVOD	12/11/17	12/11/17	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	6.92	14.32	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54			
01/21/16	MCCNN_Layoff	01/21/16	01/21/16	H	8.94	5.53	5.63	5.06	3.59	12.6	7.23	7.53	8.99	5.79	23.0	12.8	15.0	14.7	5.85	10.4	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54			
04/26/16	MC-CNN-FM	04/26/16	04/26/16	H	8.47	7.35	24.50	7.18	4.71	16.8	24.47	7.37	6.97	2.82	20.7	21	17.4	15.4	1.51	7.90	12.6	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54	0.54		

## KITTI 2015

Environment	GPU	Time	Memory	Accuracy	Comments
1 core	1.5 GHz (C/C++)	2.04	1.00	99.9	KITTI 2015
	GeForce 980 Ti	2.04	1.00	99.9	KITTI 2015
	GeForce 970	2.04	1.00	99.9	KITTI 2015
	GeForce 960	2.04	1.00	99.9	KITTI 2015
	GeForce 950	2.04	1.00	99.9	KITTI 2015
	GeForce 940	2.04	1.00	99.9	KITTI 2015
	GeForce 930	2.04	1.00	99.9	KITTI 2015
	GeForce 920	2.04	1.00	99.9	KITTI 2015
	GeForce 910	2.04	1.00	99.9	KITTI 2015
	GeForce 900	2.04	1.00	99.9	KITTI 2015
	GeForce 880	2.04	1.00	99.9	KITTI 2015
	GeForce 870	2.04	1.00	99.9	KITTI 2015
	GeForce 860	2.04	1.00	99.9	KITTI 2015
	GeForce 850	2.04	1.00	99.9	KITTI 2015
	GeForce 840	2.04	1.00	99.9	KITTI 2015
	GeForce 830	2.04	1.00	99.9	KITTI 2015
	GeForce 820	2.04	1.00	99.9	KITTI 2015
	GeForce 810	2.04	1.00	99.9	KITTI 2015
	GeForce 800	2.04	1.00	99.9	KITTI 2015
	GeForce 790	2.04	1.00	99.9	KITTI 2015
	GeForce 780	2.04	1.00	99.9	KITTI 2015
	GeForce 770	2.04	1.00	99.9	KITTI 2015
	GeForce 760	2.04	1.00	99.9	KITTI 2015
	GeForce 750	2.04	1.00	99.9	KITTI 2015
	GeForce 740	2.04	1.00	99.9	KITTI 2015
	GeForce 730	2.04	1.00	99.9	KITTI 2015
	GeForce 720	2.04	1.00	99.9	KITTI 2015
	GeForce 710	2.04	1.00	99.9	KITTI 2015
	GeForce 700	2.04	1.00	99.9	KITTI 2015
	GeForce 690	2.04	1.00	99.9	KITTI 2015
	GeForce 680	2.04	1.00	99.9	KITTI 2015
	GeForce 670	2.04	1.00	99.9	KITTI 2015
	GeForce 660	2.04	1.00	99.9	KITTI 2015
	GeForce 650	2.04	1.00	99.9	KITTI 2015
	GeForce 640	2.04	1.00	99.9	KITTI 2015
	GeForce 630	2.04	1.00	99.9	KITTI 2015
	GeForce 620	2.04	1.00	99.9	KITTI 2015
	GeForce 610	2.04	1.00	99.9	KITTI 2015
	GeForce 600	2.04	1.00	99.9	KITTI 2015
	GeForce 590	2.04	1.00	99.9	KITTI 2015
	GeForce 580	2.04	1.00	99.9	KITTI 2015
	GeForce 570	2.04	1.00	99.9	KITTI 2015
	GeForce 560	2.04	1.00	99.9	KITTI 2015
	GeForce 550	2.04	1.00	99.9	KITTI 2015
	GeForce 540	2.04	1.00	99.9	KITTI 2015
	GeForce 530	2.04	1.00	99.9	KITTI 2015
	GeForce 520	2.04	1.00	99.9	KITTI 2015
	GeForce 510	2.04	1.00	99.9	KITTI 2015
	GeForce 500	2.04	1.00	99.9	KITTI 2015
	GeForce 490	2.04	1.00	99.9	KITTI 2015
	GeForce 480	2.04	1.00	99.9	KITTI 2015
	GeForce 470	2.04	1.00	99.9	KITTI 2015
	GeForce 460	2.04	1.00	99.9	KITTI 2015
	GeForce 450	2.04	1.00	99.9	KITTI 2015
	GeForce 440	2.04	1.00	99.9	KITTI 2015
	GeForce 430	2.04	1.00	99.9	KITTI 2015
	GeForce 420	2.04	1.00	99.9	KITTI 2015
	GeForce 410	2.04	1.00	99.9	KITTI 2015
	GeForce 400	2.04	1.00	99.9	KITTI 2015
	GeForce 390	2.04	1.00	99.9	KITTI 2015
	GeForce 380	2.04	1.00	99.9	KITTI 2015
	GeForce 370	2.04	1.00	99.9	KITTI 2015
	GeForce 360	2.04	1.00	99.9	KITTI 2015
	GeForce 350	2.04	1.00	99.9	KITTI 2015
	GeForce 340	2.04	1.00	99.9	KITTI 2015
	GeForce 330	2.04	1.00	99.9	KITTI 2015
	GeForce 320	2.04	1.00	99.9	KITTI 2015
	GeForce 310	2.04	1.00	99.9	KITTI 2015
	GeForce 300	2.04	1.00	99.9	KITTI 2015
	GeForce 290	2.04	1.00	99.9	KITTI 2015
	GeForce 280	2.04	1.00	99.9	KITTI 2015
	GeForce 270	2.04	1.00	99.9	KITTI 2015
	GeForce 260	2.04	1.00	99.9	KITTI 2015
	GeForce 250	2.04	1.00	99.9	KITTI 2015
	GeForce 240	2.04	1.00	99.9	KITTI 2015
	GeForce 230	2.04	1.00	99.9	KITTI 2015
	GeForce 220	2.04	1.00	99.9	KITTI 2015
	GeForce 210	2.04	1.00	99.9	KITTI 2015
	GeForce 200	2.04	1.00	99.9	KITTI 2015
	GeForce 190	2.04	1.00	99.9	KITTI 2015
	GeForce 180	2.04	1.00	99.9	KITTI 2015
	GeForce 170	2.04	1.00	99.9	KITTI 2015

# Stereo Benchmark Rankings

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 x 2 (no longer active).

Set: test-dense test-sparse training-dense training-sparse

Medic.	Bas 0.5	Bas 1.0	Bas 2.0	Bas 4.0	avgell	115
Mask:	none	all				

state selected  show legend **Reset** **Print** **Reference**

Selected		Show Invalid		Reset Selection		Reset Instances		Performance Metrics											
Date	Name	Res	Avg	Austr1	Austr2	Bicy2	Class1	Class2	Compu	Cross1	Cross2	Djemb1	Djemb2	Hoops	Lvgrm	Niuba	Plants	Stairs	
		MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	MP	
		8.0	8.0	8.0	8.0	8.0	8.0	8.0	8.0	8.0	8.0	8.0	8.0	8.0	8.0	8.0	8.0	8.0	
03/06/16	NOSS	H	8.04	3.57	2.84	2.39	1.93	1.51	3.34	3.34	3.19	2.32	8.55	7.45	7.06	12.53	5.21	10.01	
06/22/17	LogiExp	H	8.43	3.65	2.87	2.98	1.97	5.59	3.37	3.49	3.36	2.03	1.51	8.75	9.44	5.40	9.65	8.01	
01/24/17	3DMST	H	9.21	3.71	2.78	4.76	2.72	7.36	4.28	3.44	3.76	2.93	12.64	11.51	8.94	14.49	5.25	8.87	
03/19/17	MC-CNN+TDS	H	9.21	3.54	4.15	4.65	3.61	10.47	6.05	5.01	5.19	2.67	10.84	9.62	6.91	11.44	1.01	7.04	
05/12/19	PLM	H	7.11	3.48	3.48	3.48	3.48	5.24	0.92	4.54	3.96	4.04	2.37	13.1	12.33	12.27	16.2	5.88	10.81
05/12/19	LW-CNN	H	7.04	4.65	7.04	7.04	7.04	7.04	7.04	7.04	7.04	7.04	12.27	12.34	12.36	14.18	4.72	12.00	
04/12/19	MemShredExt	H	7.26	4.41	4.41	4.41	4.41	4.41	4.41	4.41	4.41	4.41	12.73	12.74	10.48	14.49	7.05	8.85	
10/12/17	FeN-D2OR	H	7.23	4.68	4.11	5.00	3.03	5.42	0.95	1.92	3.02	4.41	15.11	14.1	13.44	13.97	5.02	14.33	
05/26/16	APAF-Net	H	7.26	5.43	4.91	3.43	5.11	5.71	21.69	4.99	4.31	4.23	3.24	13.43	9.78	7.32	13.44	6.28	8.46
03/11/18	SM-Forest	H	7.37	4.71	3.69	4.93	3.18	11.11	5.57	5.57	5.81	2.65	14.5	13.24	13.11	14.18	5.93	11.12	
01/19/16	NTDE	H	7.44	5.72	1.46	5.92	2.30	10.43	5.71	5.30	5.54	11.20	13.5	14.1	12.6	13.98	6.38	12.20	
09/29/18	ESD	H	7.44	5.41	4.73	4.93	3.74	5.5	5.40	4.40	4.40	4.70	13.23	13.24	14.49	16.24	5.75	11.90	
02/28/15	MC-CNN-act	H	8.01	6.59	4.56	5.96	2.82	11.47	5.81	8.82	6.88	2.77	16.3	14.1	13.2	13.0	6.40	11.11	
11/13/15	MC-CNN+RBS	H	8.42	6.05	6.16	6.24	3.27	11.11	6.26	8.87	6.83	3.21	17.0	12.5	13.5	7.94	9.99	1.01	
09/13/16	SNP-RSM	H	8.75	5.46	4.85	5.05	3.37	10.4	7.31	8.73	9.37	3.56	14.3	14.7	14.9	12.83	10.1	10.81	
12/11/17	COVD	H	8.87	5.44	5.74	5.81	5.11	4.82	12.8	6.51	9.49	1.91	9.66	3.15	16.6	14.4	14.5	6.27	13.24
02/28/15	MC-CNN_Layout	H	9.04	5.53	5.53	5.66	3.98	12.6	7.23	7.53	8.66	5.79	12.0	13.6	15.0	14.7	5.85	10.43	

# Group A

#13

## MC-CNN acrt

#15

- Group A and B have no methods in common!
- Group A entries all use MC-CNN acrt but no other “deep learning” technique!
- Group B methods do NOT use MC-CNN acrt; they use ResNet, 3D convolutions, 3D deconvolutions, U-shaped Nets, RNNs; End to end learning is very popular!

KITTI 2015

Author and T. LaCapra: Some Remarks on Twistor & Conformal Theory, Remarks on Conformal Gauge Theories, Submitted to JHEP

# Conclusions



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

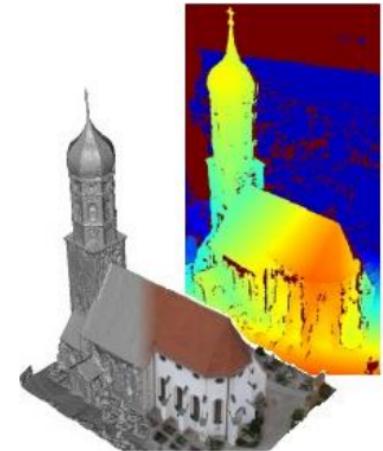
	Stereo	MVS	Flow	Depth	Semantic	Instance
Middlebury	X	X	X			
KITTI	X		X	X	X	X
MPI Sintel				X		
ETH3D	X	X				
HD1K				X		
ScanNet				X	X	X
Cityscapes					X	X
WildDash					X	X

# Outline

- Stereo Matching: New Trends
- **Semi-Global Stereo Matching (SGM)**
  - SGM with Surface Orientation Priors
- Stereo Scene Flow with Motion Segmentation
- Trajectory Planning for Aerial Multi View Stereo
- Deep 6D 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 in real-time on FPGAs, GPUs ...



# Scansline Optimization (1D)

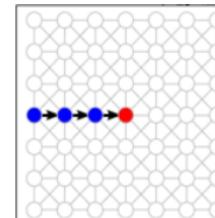
Minimize:

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

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

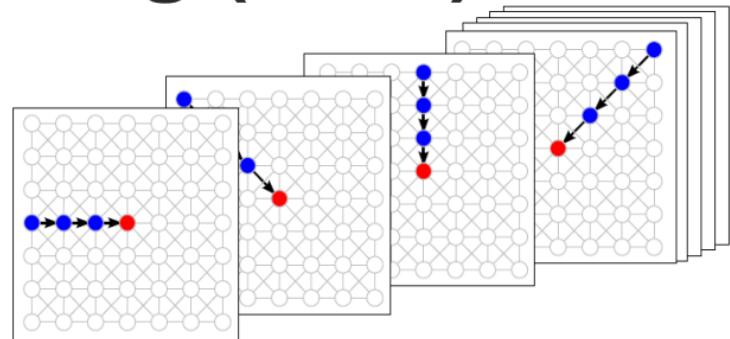
$$L_{\mathbf{r}}(\mathbf{p}, d) = C_{\mathbf{p}}(d) + \min_{d' \in \mathcal{D}} (L_{\mathbf{r}}(\mathbf{p} - \mathbf{r}, d') + V(d, d')).$$

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



# Semi Global Matching (SGM)

- For 8 directions
  - calculate aggregated costs



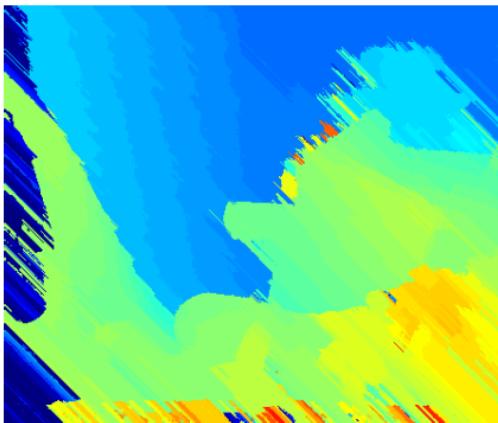
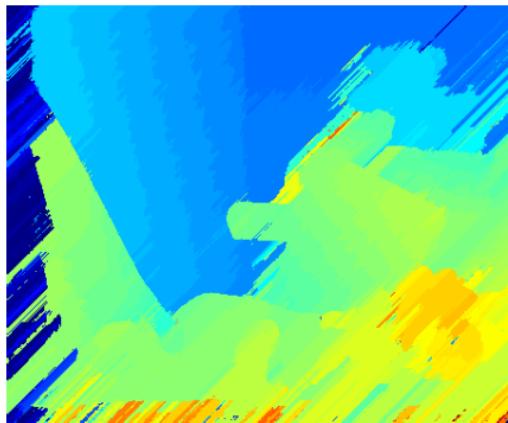
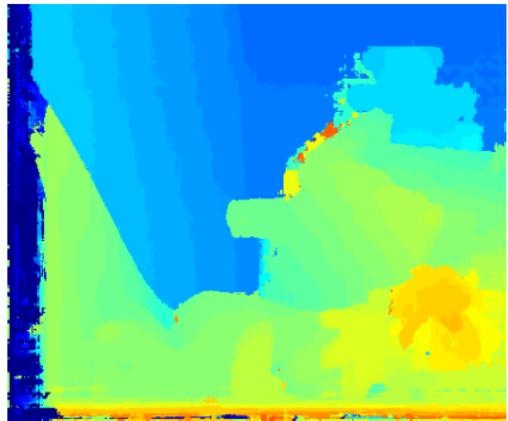
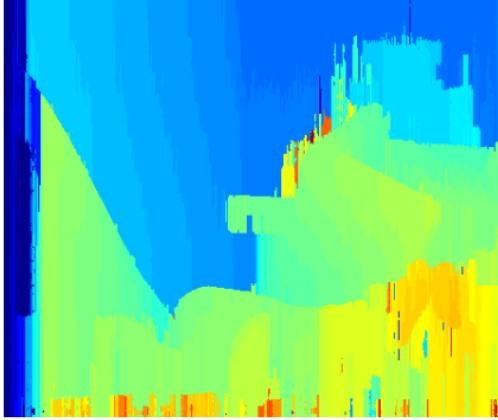
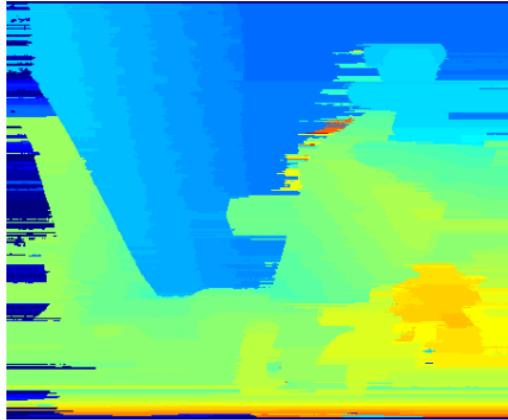
$$L_{\mathbf{r}}(\mathbf{p}, d) = C_{\mathbf{p}}(d) + \min_{d' \in \mathcal{D}} (L_{\mathbf{r}}(\mathbf{p} - \mathbf{r}, d') + V(d, d')).$$

- 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)

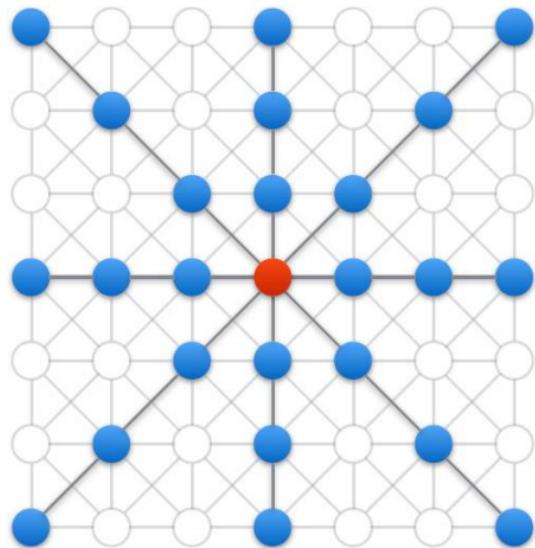


# Semi Global Matching [Hirschmüller 2005]

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

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

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



# Semi Global Stereo Matching with Surface Orientation Priors

3DV 2017

Daniel Scharstein

Middlebury College

Tatsunori Taniai

RIKEN, Tokyo

Sudipta Sinha

Microsoft Research

# Semi Global Matching [Hirschmüller 2005]

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

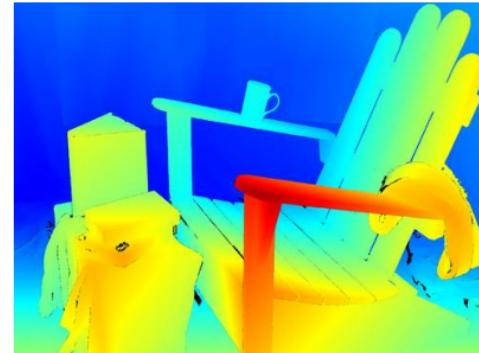
$$E(D) = \sum_{\mathbf{p}} C_{\mathbf{p}}(d_{\mathbf{p}}) + \sum_{\mathbf{p}, \mathbf{q} \in \mathcal{N}} V(d_{\mathbf{p}}, d_{\mathbf{q}})$$

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

- Fronto parallel bias
- Inaccurate on slanted untextured surfaces

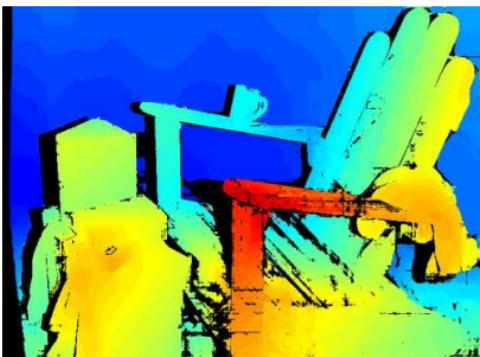


Left Image

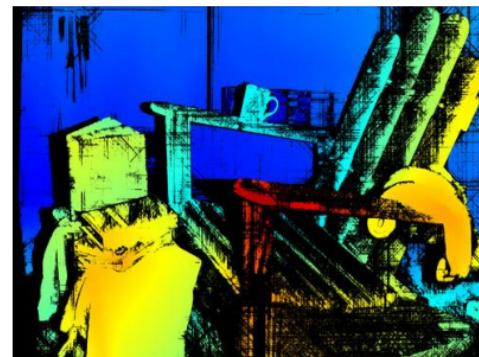


Ground Truth

SGM\* (quarter resolution)



SGM\* (full resolution (6 MP))



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

# SGM-P: SGM with orientation priors

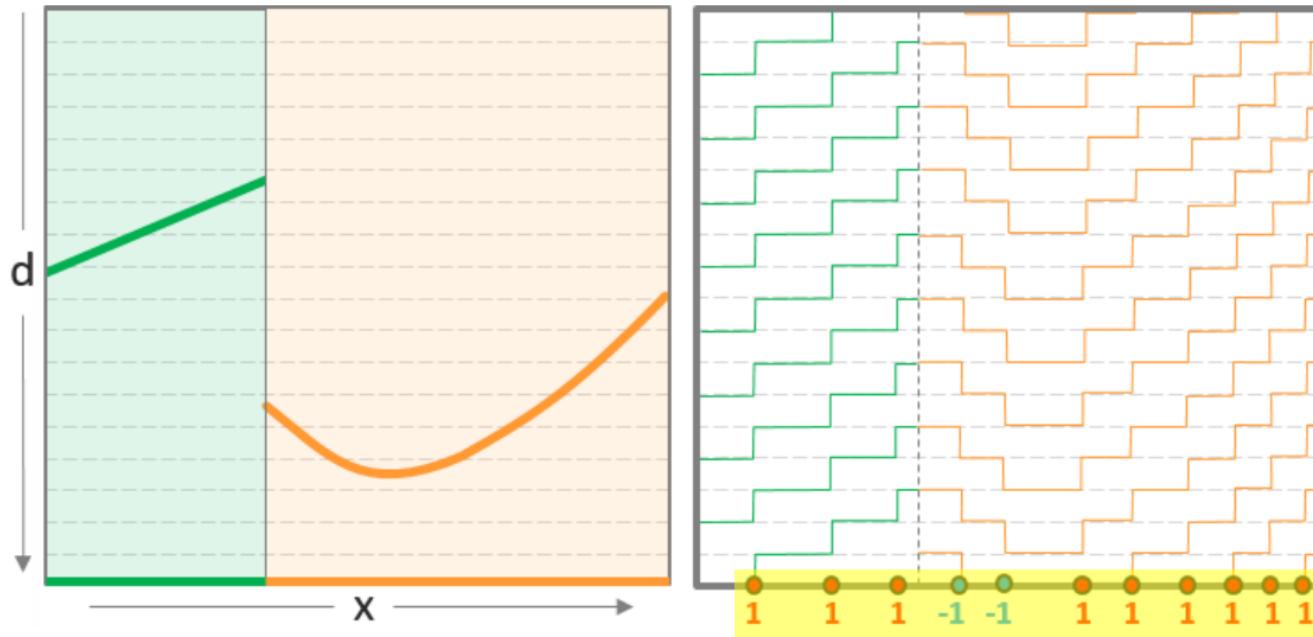
[Scharstein, Taniai, Sinha, 3DV 2017]

- What if we knew the surface slant?
- Replace fronto-parallel bias with bias parallel to surface

Idea:

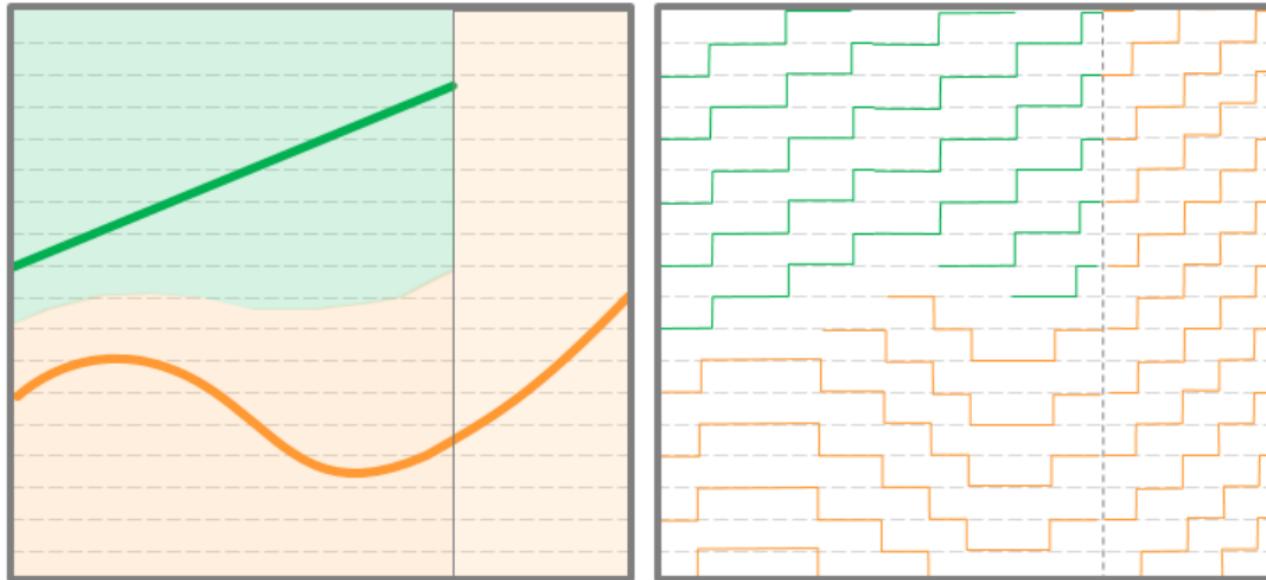
- *Rasterize* disparity surface prior (at arbitrary depth)
- Adjust  $V(d, d')$  to follow discrete disparity “steps”

# 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

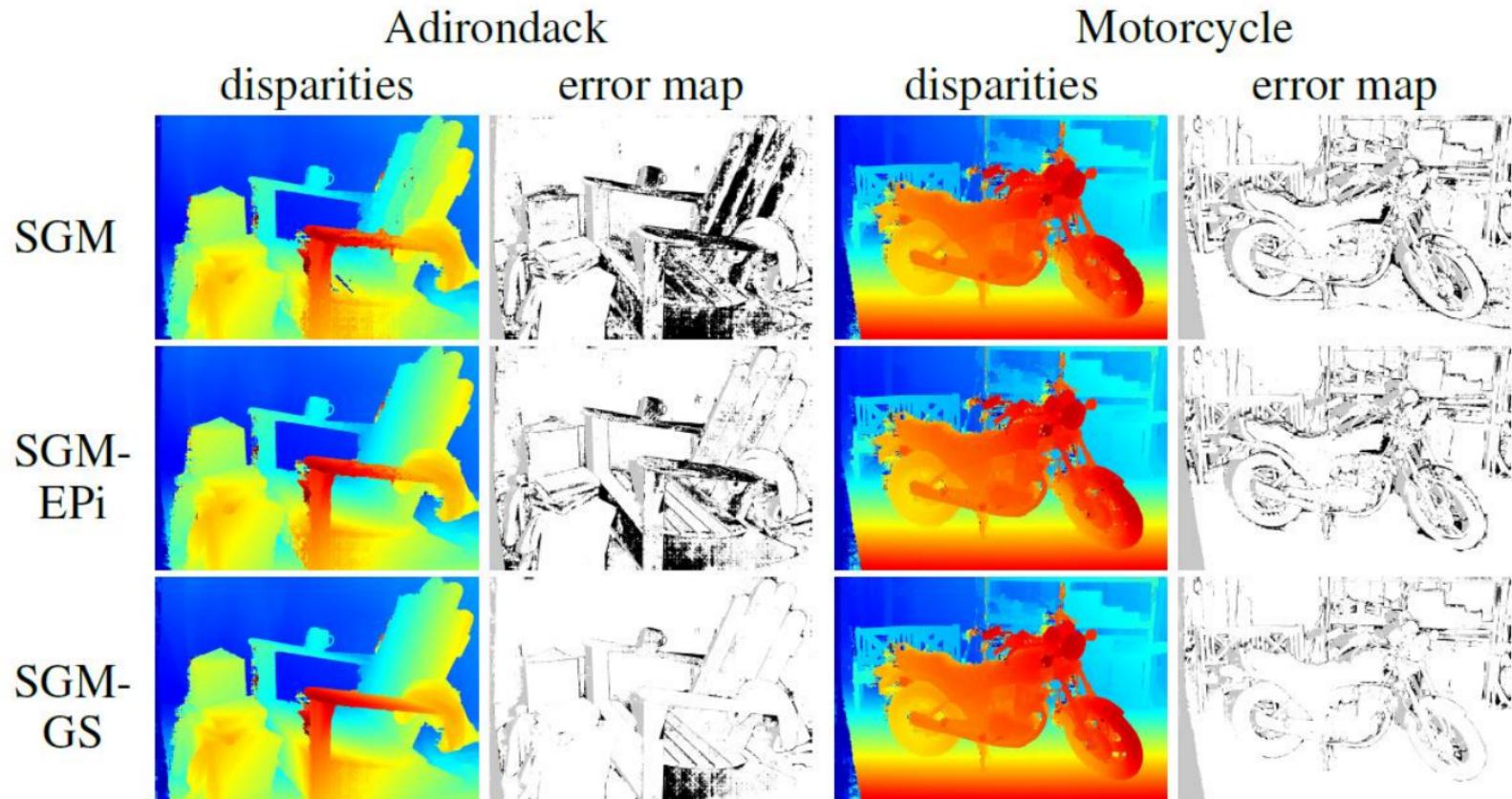


Input

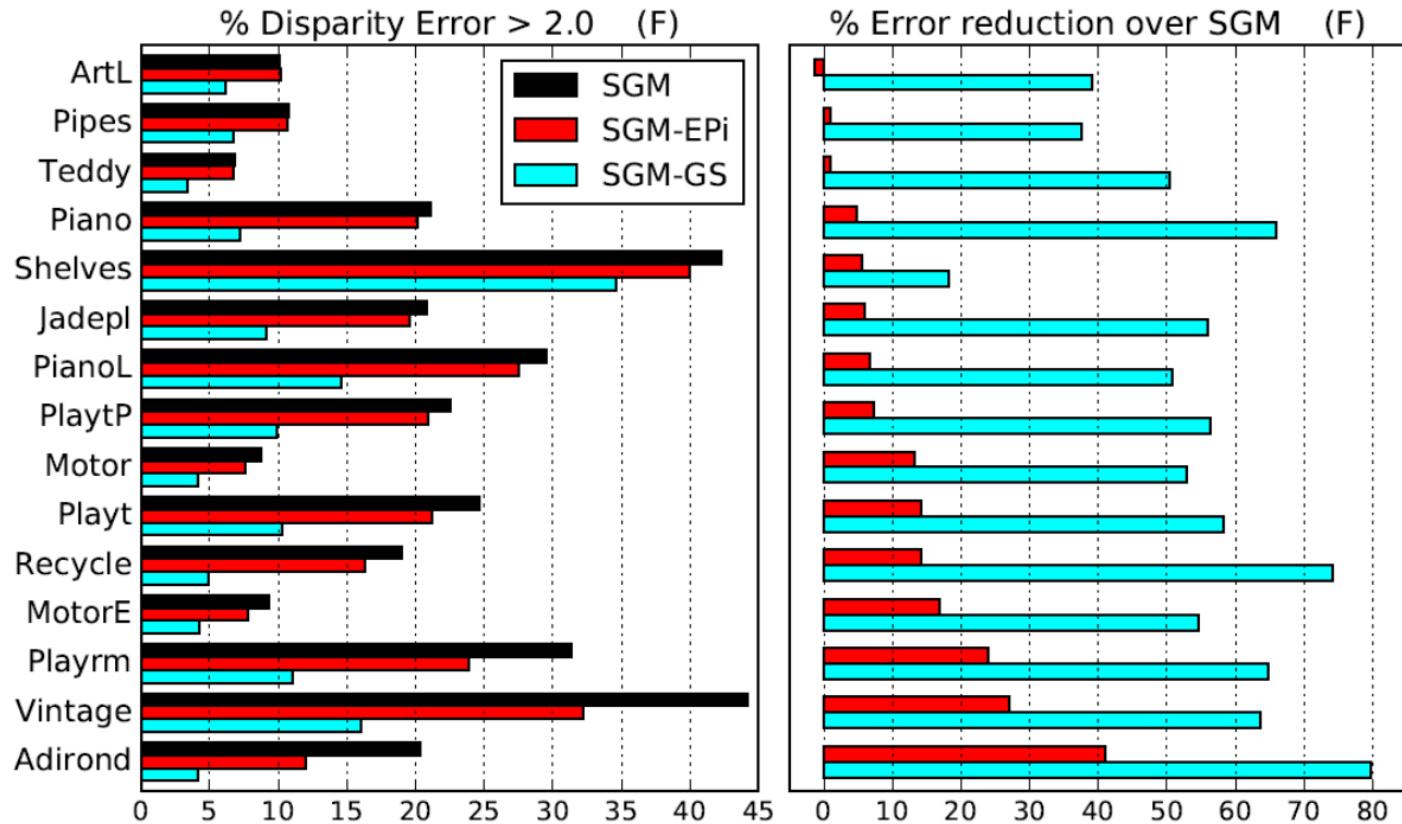
SGM-EPi

SGM-GS

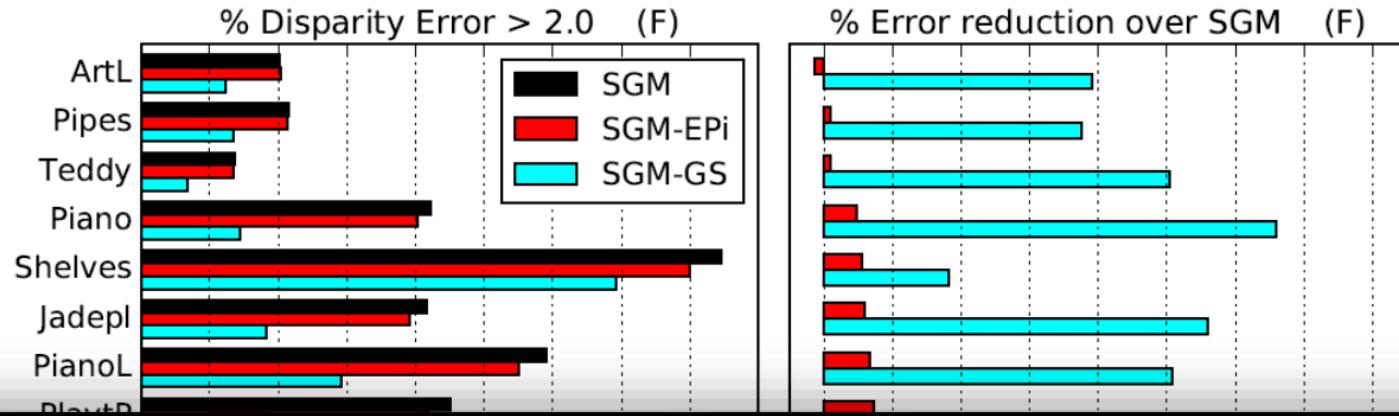
# SGM-P: Results



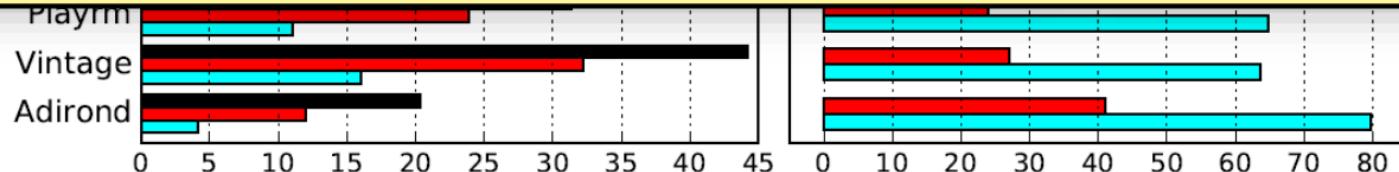
# SGM-P: Results



# SGM-P: Results



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



# Outline

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

# Fast Multi-frame Stereo Scene Flow with Motion Segmentation

CVPR 2017

Tatsunori Taniai

RIKEN, Tokyo

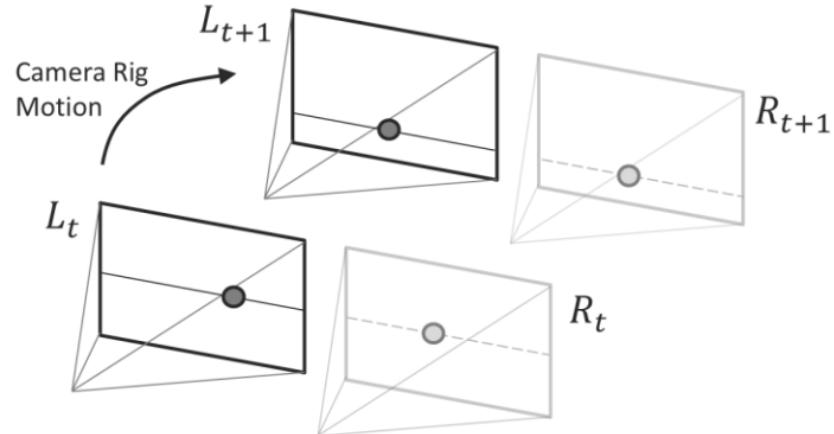
Sudipta Sinha

Microsoft Research

Yoichiro Sato

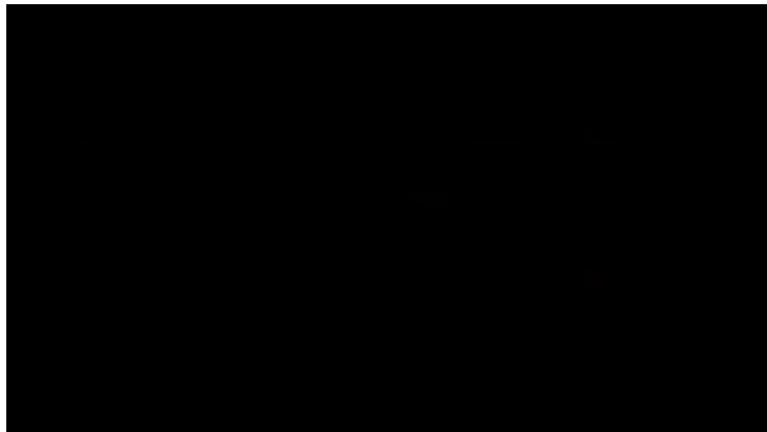
Univ. of Tokyo

# 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

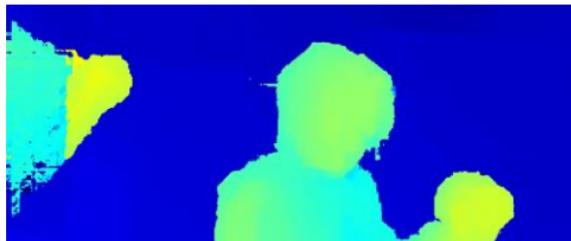
# Motivation

Efficient, unified method for

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



Disparity Map



Optical Flow



Moving Object Segmentation



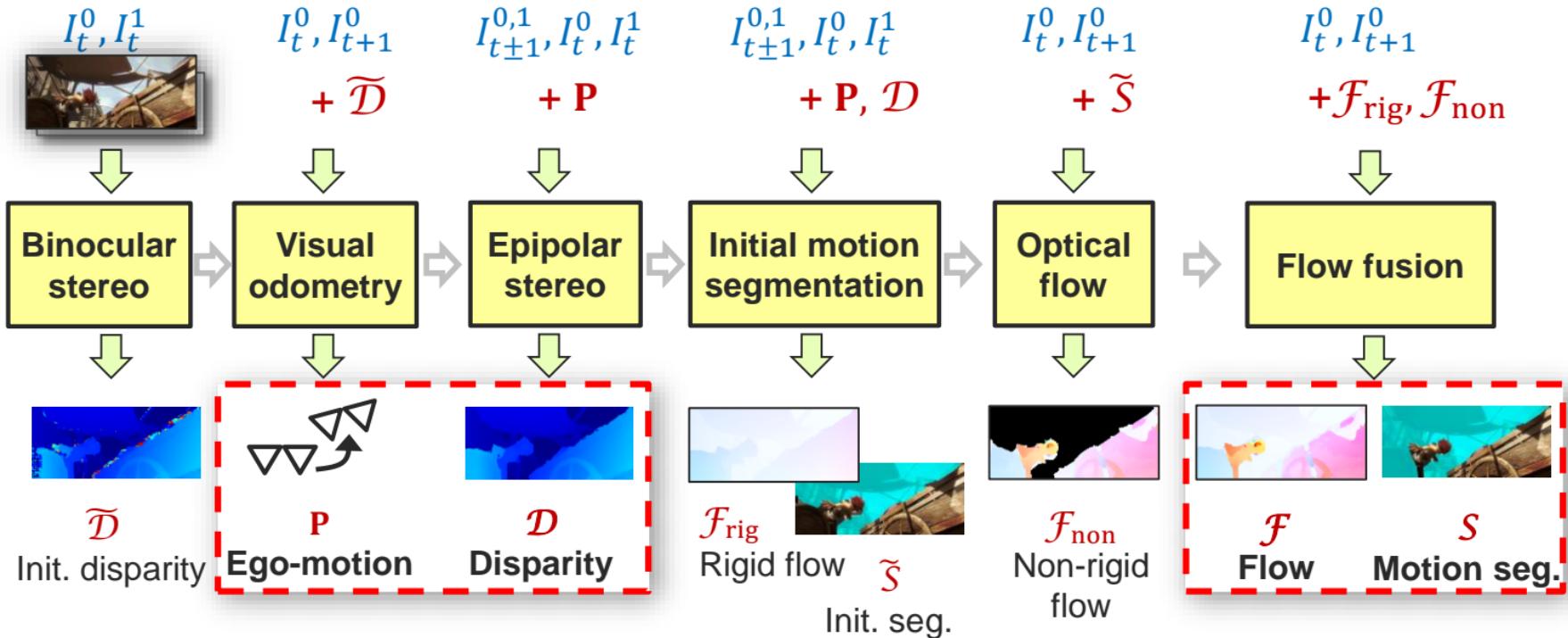
# Main Idea: Dominant Rigid Scene Assumption



- 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]	<b>3.02</b>	<b>10.52</b>	<b>4.27</b>	<b>5.13</b>	<b>15.11</b>	<b>6.79</b>	<b>5.33</b>	<b>17.02</b>	<b>7.28</b>	<b>6.61</b>	<b>23.60</b>	<b>9.44</b>	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	<b>FSF+MS (ours)</b>	<b>5.72</b>	<b>11.84</b>	<b>6.74</b>	<b>7.57</b>	<b>21.28</b>	<b>9.85</b>	<b>8.48</b>	<b>29.62</b>	<b>12.00</b>	<b>11.17</b>	<b>37.40</b>	<b>15.54</b>	<b>2.7 s</b>
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

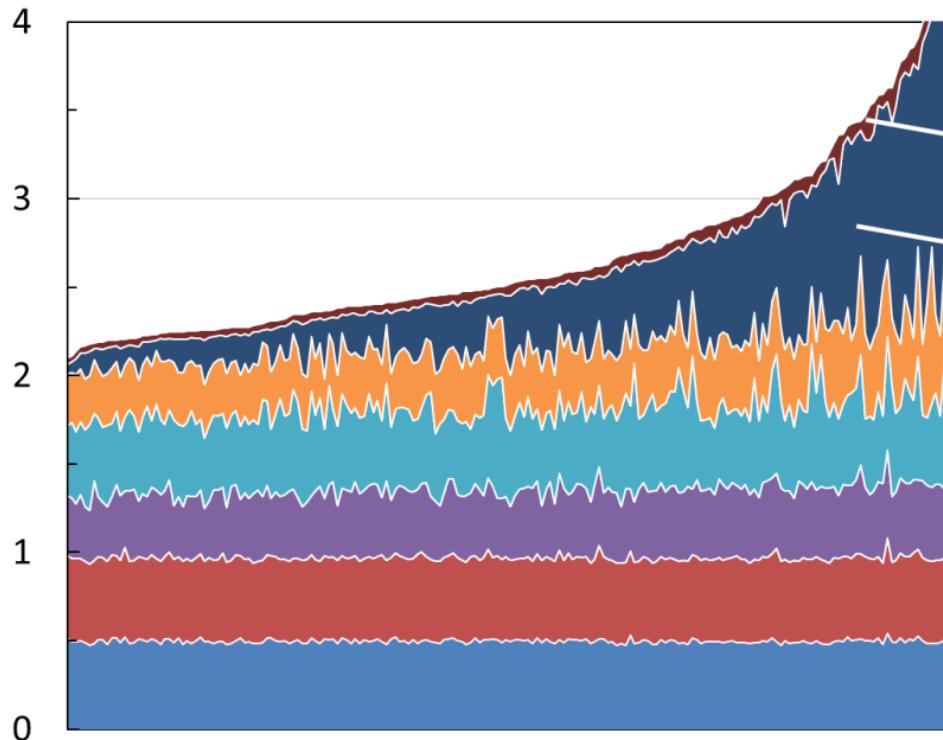


200 road scenes with multiple moving objects

Rank	SF-all	Time
1	<b>9.44</b>	300 s
2	10.63	50 min
3	<b>15.54</b>	<b>2.7 s</b>
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)



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
- Trajectory Planning for Aerial Multi View Stereo
- Deep 6D Object Pose Prediction

# Submodular Trajectory Optimization for Aerial 3D Scanning

ICCV 2017

Mike Roberts<sup>1,2</sup> Debadatta Dey<sup>2</sup> Anh Truong<sup>3</sup> Sudipta Sinha<sup>2</sup>  
Shital Shah<sup>2</sup> Ashish Kapoor<sup>2</sup> Pat Hanrahan<sup>1</sup> Neel Joshi<sup>2</sup>

<sup>1</sup>Stanford University

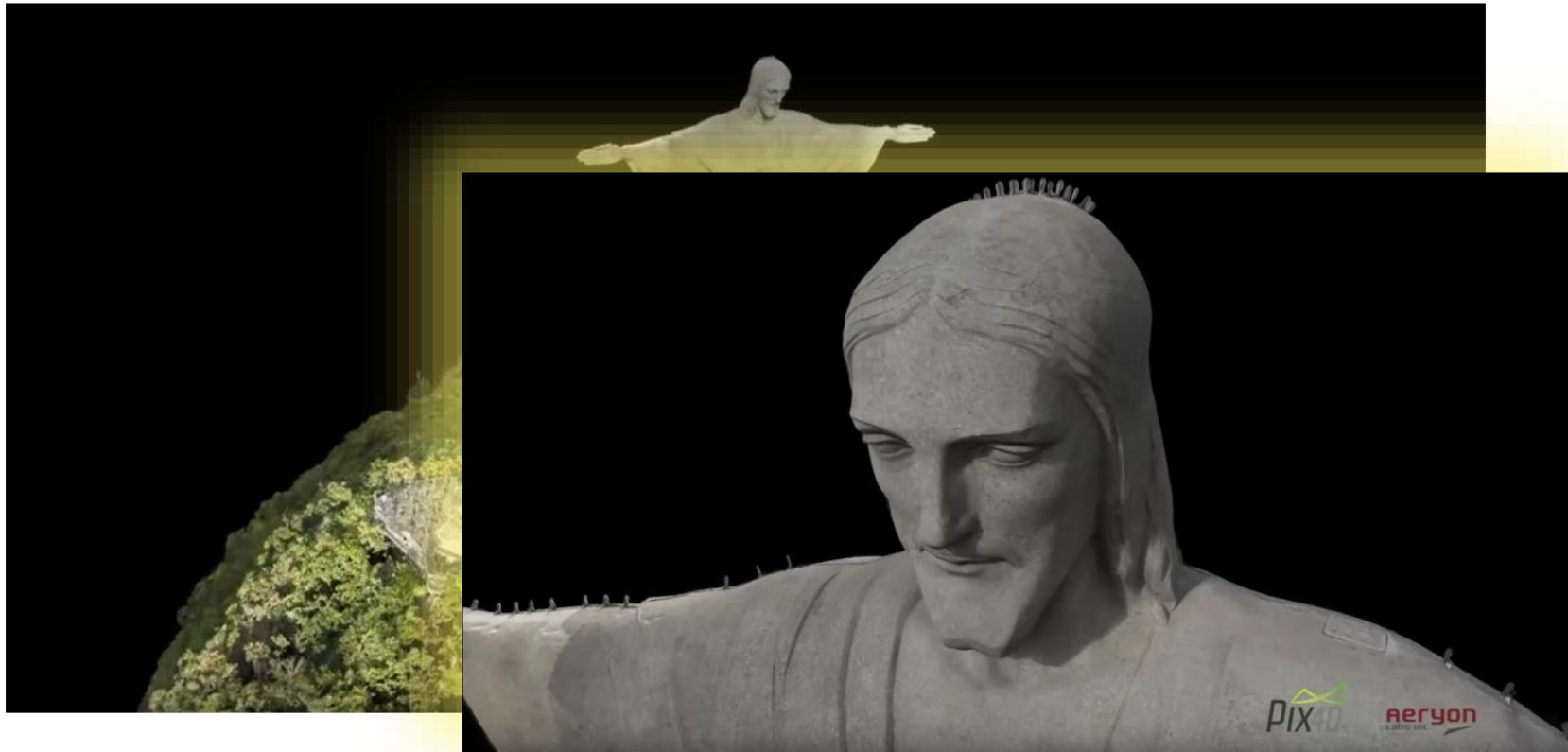
<sup>2</sup>Microsoft Research

<sup>3</sup>Adobe Research

# Acquiring imagery using drones



# Multi-view Stereo Reconstruction



PIX4D

aeryon  
Labs Inc

# Manual Planning Prior to Capture



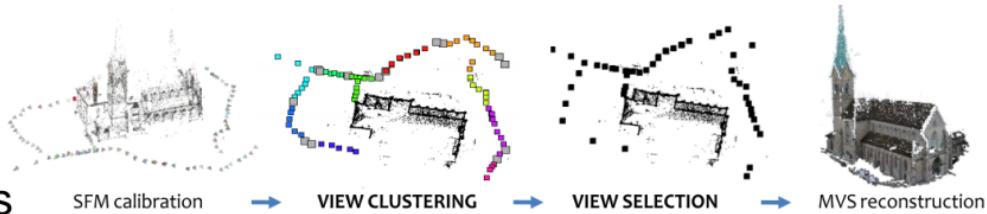
- Waypoints planned by human experts ...
  - Several redundant flight trajectories were flown
- 3,500 images from 6 days with 19 ten-minute flights
- Projeto Redentor (Pix4D whitepaper, 2015)

# Our Goal

- Automatically generate optimized trajectories for 3D scanning using drones, such that
  1. the acquired images will produce an accurate 3D model when processed using a Multi View Stereo (MVS) algorithm.
  2. the UAV makes best use of its limited flight time budget.
- Processing happens post flight.
- Battery typically lasts 15—20 minutes.

# Related Work

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



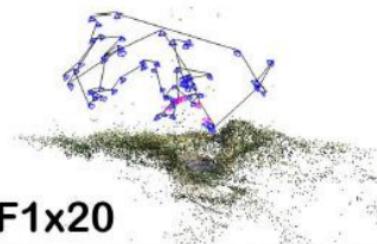
- Next-best-view planning
  - Information-gain maximization [Isler+ 2016]
  - Robotic RGB-D 3D scanning [Wu+ 2014]
  - No travel budget constraints



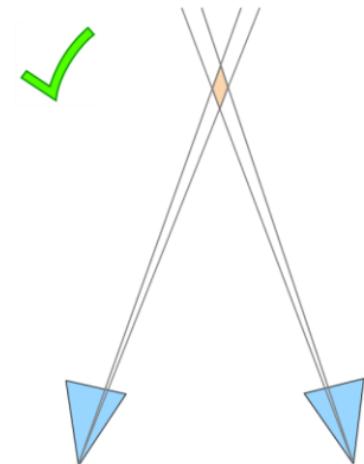
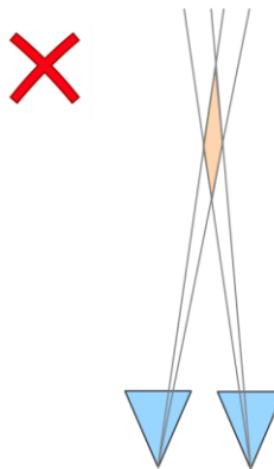
- Real-time Drone view planning

[Mostegel+ 2016]

- Greedy technique; heuristic-based



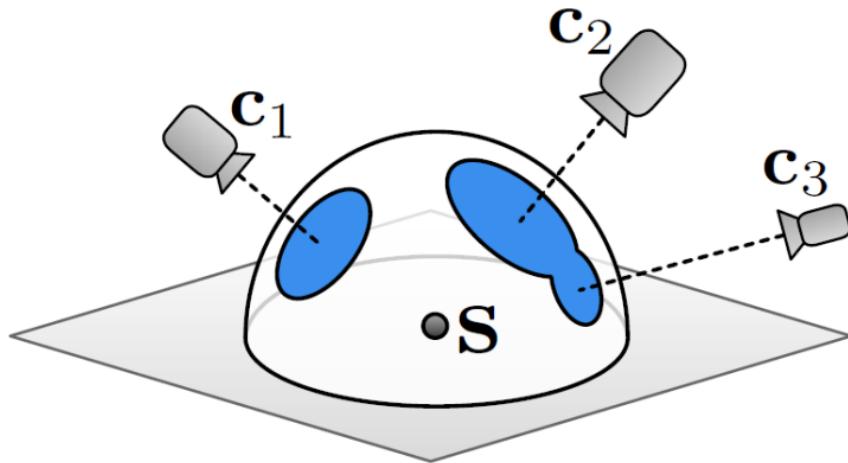
# Diverse Viewpoints help Multi View Stereo



Preference for

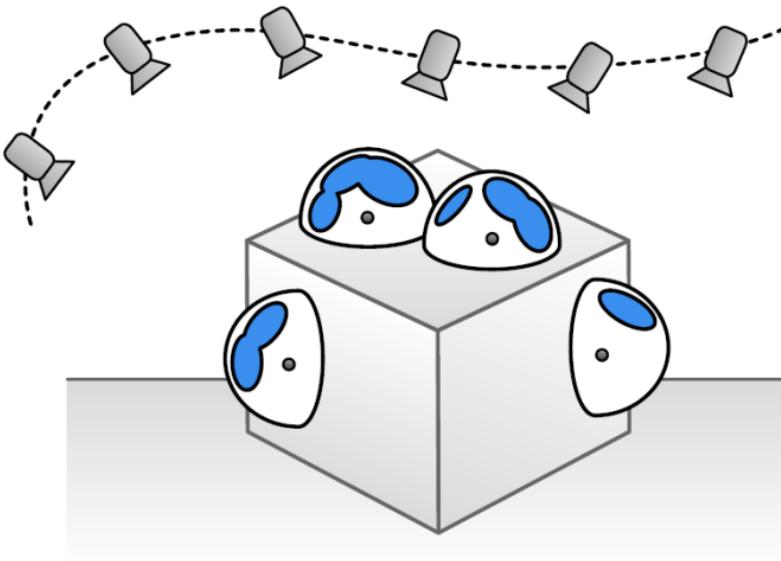
- diverging viewing angles
- close-up views
- fronto parallel views of surfaces

# 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

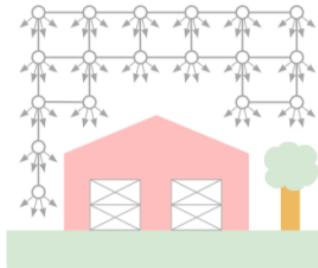


Similarly, we define coverage for multiple surface points observed from multiple camera viewpoints.

# Method

- Evaluating coverage function requires knowledge of scene geometry
- Thus, we follow a two-staged procedure.
  1. Fly an easy-to-generate trajectory;
  2. Compute coarse reconstruction (SFM → MVS → meshing)
  3. Plan optimized trajectory based on mesh from step 2.
  4. Fly trajectory computed in step 3.
  5. Run SFM + MVS on images from step 1 and 4.

# Planning Optimized Trajectories



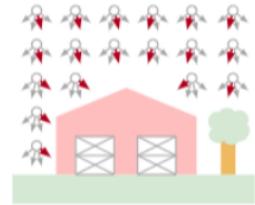
Graph of all possible camera location (and orientation); edge weights are Euclidean distances between locations.

Propose to solve the problem in two steps.

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

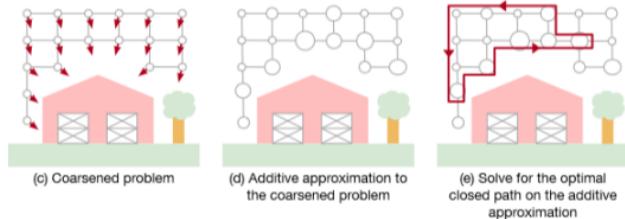
# Solving for Camera Orientations

- 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



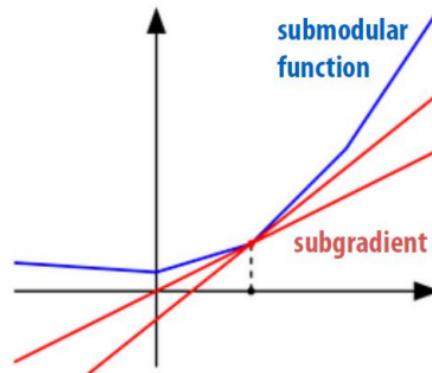
# Solving for Camera Positions

- 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*.



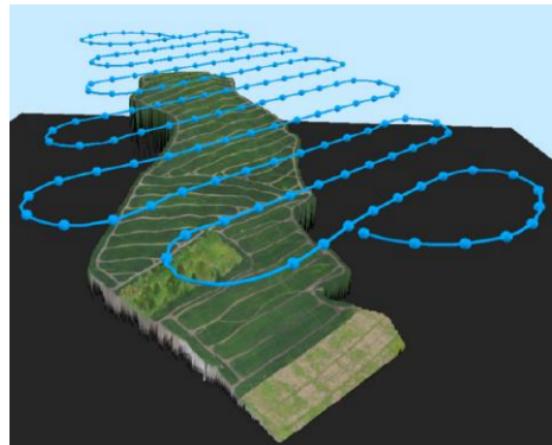
# Solving for Camera Positions

- Choose a good sub-gradient (additive approximation) for our submodular function
- Approximation yields an instance of the orienteering problem  
$$\begin{aligned} & \text{maximize} && \sum c(s) \\ & \text{subject to} && T(S) < B \end{aligned}$$
- Solve as an integer linear program (ILP)



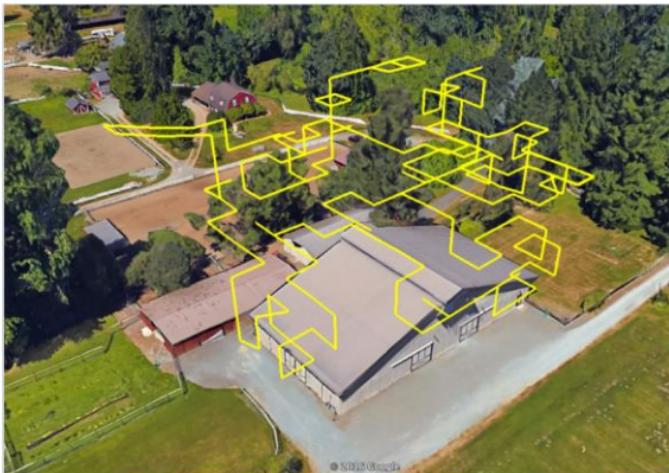
# Results

- Pix4D for 3D reconstruction
  - Outputs texture-mapped 3D model
- Baselines:
  - **Overview:** Lawn-mower pattern
  - **Random:**
    - recover coarse 3d model; estimate free space.
    - select random points in free space.
    - compute TSP tour.



# Results

Our computed trajectories visualized in Google Earth



Barn



MSR Redmond

# Results



*Insert video (ICCV supplementary video here)*

# Outline

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

# Real-Time Seamless Single Shot 6D Object Pose Prediction

CVPR 2018

Bugra Tekin

EPFL

Sudipta Sinha

Microsoft Research

Pascal Fua

EPFL

# 3D Recognition, 2D-3D Model Alignment

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

Classical methods:



Lowe 2001



Rothganger+ 2005



Lepetit+ 2005

- Recognizing Image Patches
- Scale, Affine invariant features
- Geometric verification (rigid scenes)
- Worked for textured, distinctive objects
- Required a small # of training images

# Object 6D Pose Estimation

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

## RGB-D methods:

- Lai+ 2010
- Hinterstoisser+ 2012
- Brachmann+ 2014, 2016

- Classical Object Recognition
- Fast Image Retrieval

## CNN methods:

- Rad + Lepetit 2017
- Kehl+ 2017
- Xiang+ 2017

- Global deep feature representations
- Not much use of geometry
- Promising for small, texture-less objects
- Huge training set needed

# Texture-less Object 6D Pose Datasets



LINEKIT [2012]  
15 objects



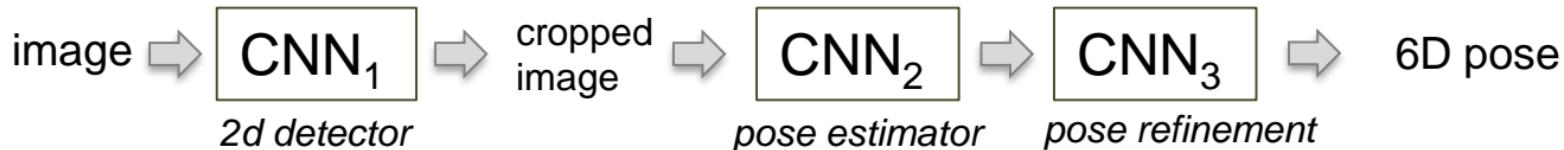
T-LESS [2017]  
30 objects



YCB-VIDEO [2018]  
21 objects

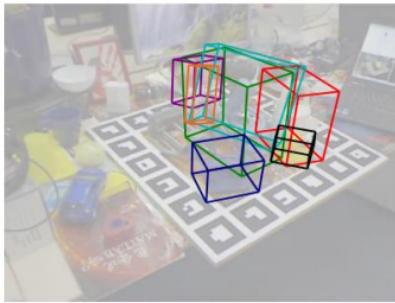
# Deep 6D object pose estimation

- BB8 [Rad and Lepetit 2017]



# 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.



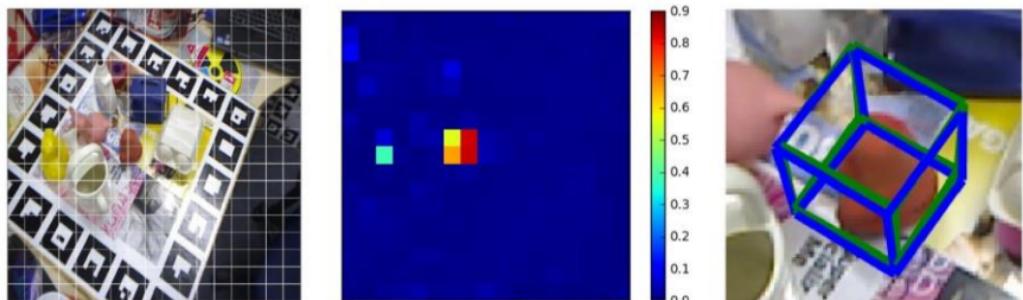
# Our Method

## Training:

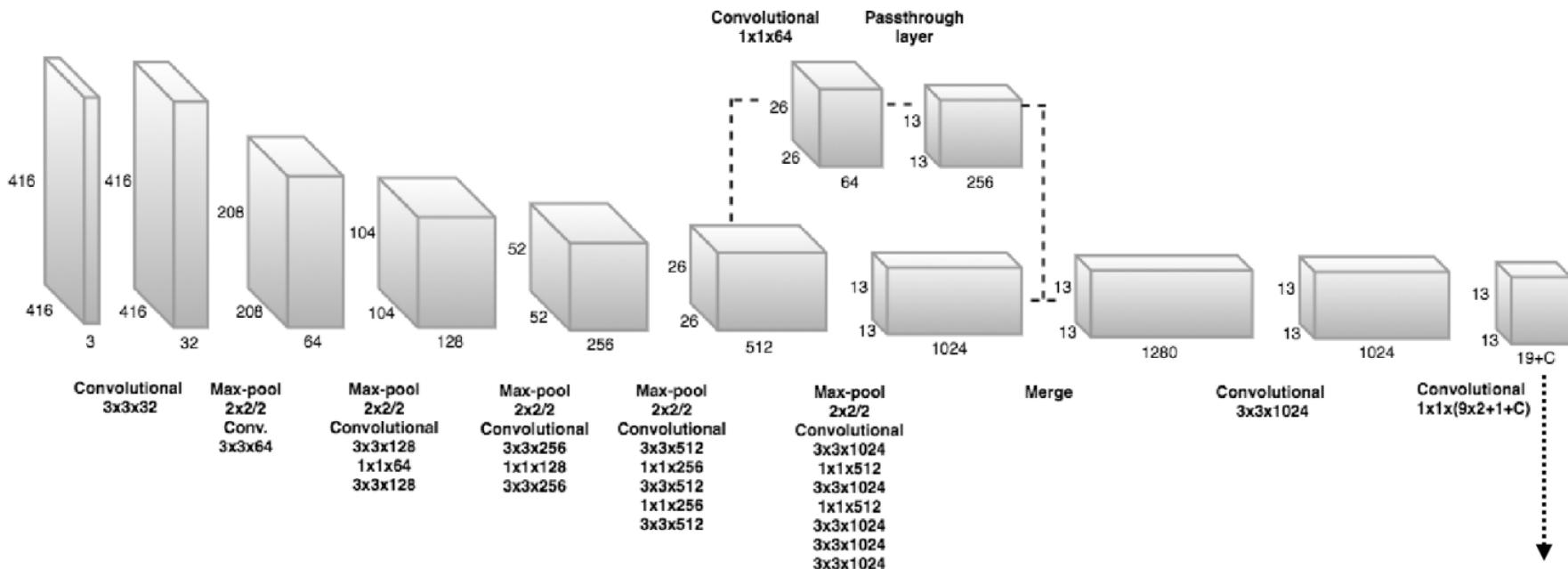
- ground truth 2D coordinates of the 9 control points are the targets
- modify YOLO loss function (for confidence estimation)
- data augmentation

## Testing:

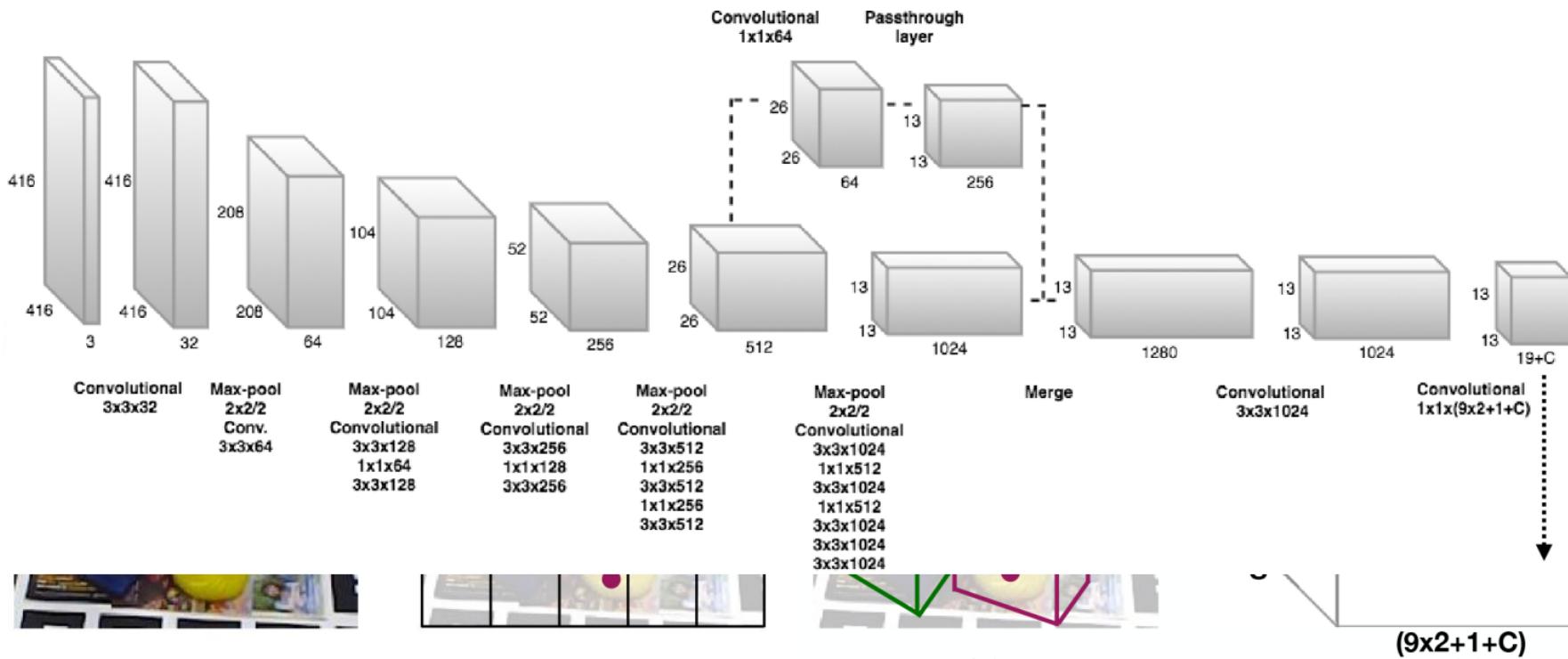
- Subpixel refinement
- PnP (RANSAC, least squares)



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

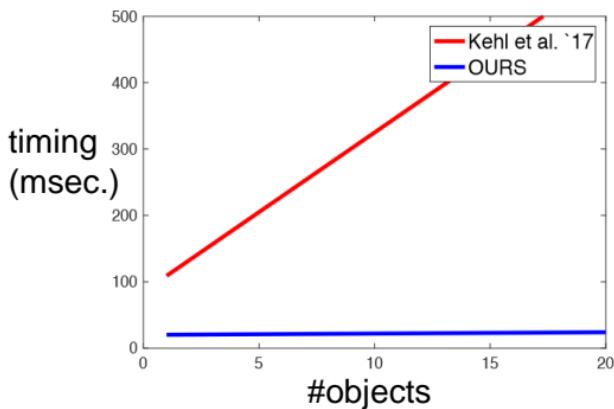
3D metric

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

# 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
$416 \times 416$	89.71	94 fps
$480 \times 480$	90.00	67 fps
$544 \times 544$	90.37	50 fps
$688 \times 688$	90.65	43 fps



When input image is resized, our method remains accurate and runs much faster

# Conclusions

- State of the art in stereo matching; new challenges
- Improvements to Semi Global Matching
  - Incorporating soft surface orientation priors
- Fast scene flow with motion segmentation
- Camera path planning for improved multi-view stereo
- Deep single shot 6D object pose estimation
  - CNN architecture conceptually simpler (~YOLO architecture) and faster