

High-Resolution Stereo Matching

Sudipta N. Sinha

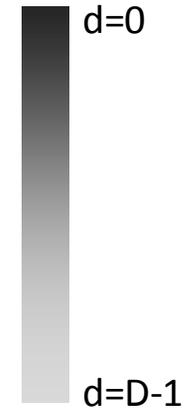
Interactive Visual Media Group
Microsoft Research

Stereo matching

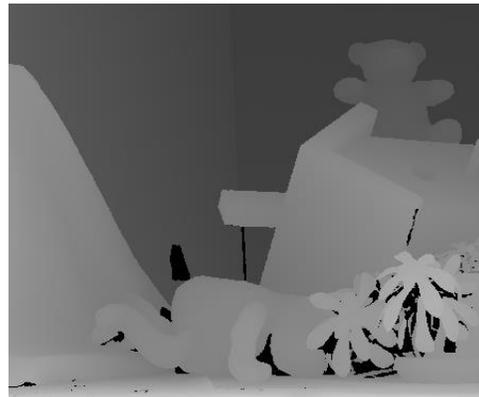
- Input: rectified image pair
- Output: disparity map



SGM [Hirschmuller 2005]



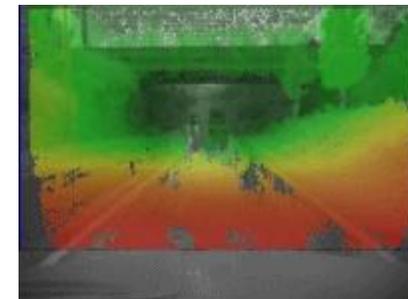
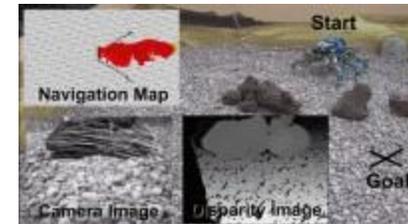
GT



6% errors
 $|\Delta d| \leq 1$ pixel
in non-occluded
regions

Stereo – applications

- Dense 3D reconstruction
- Robot navigation
- Automated driving
- Gaming, user interfaces
- Virtual viewpoint correction
- 3D movie editing



High-resolution Cameras

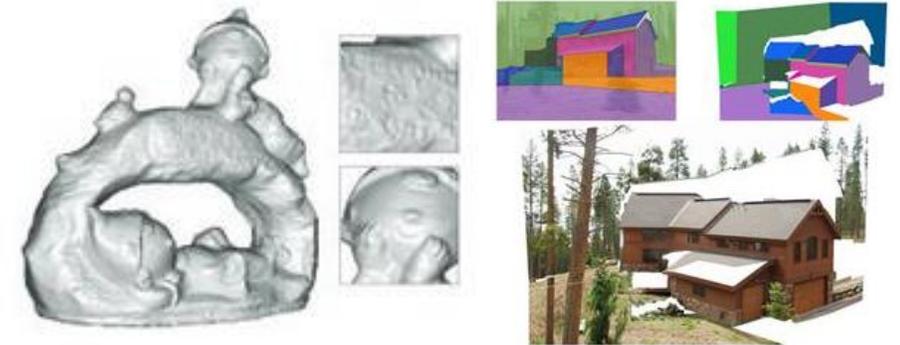
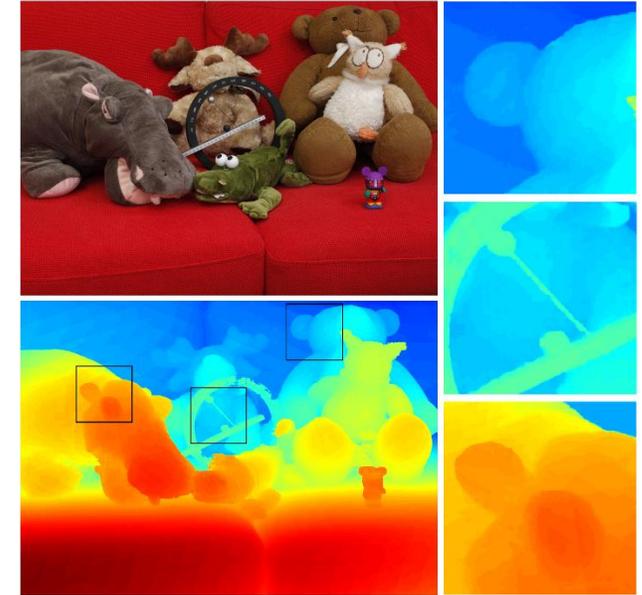


Nokia Lumia 1020
(41 MP sensor)

- High resolution cameras are everywhere
- 8+ MP on most commodity camera phones
- Impact on stereo matching techniques
 - Accuracy
 - Speed

Outline

- Stereo Matching
 - State of the art
 - High-resolution and Scalability
- High-Resolution Stereo Matching using Local Plane Sweeps
- Surface-based stereo matching
- Multi-view Stereo + Photometric Stereo



Middlebury Stereo benchmark

Images up to 450 x 375 (< 0.2 MP), D = 16 ... 60

vision.middlebury.edu

stereo • mview • MRF • flow • color

Stereo • Evaluation • Datasets • Code • Submit

Middlebury Stereo Evaluation - Version 2

[New features and main differences to version 1.](#)
[Submit and evaluate your own results.](#)

Open a new window for each link

Error Threshold = 1		Sort by nonocc			Sort by all			Sort by disc			Average percent of bad pixels (explanation)			
Error Threshold... ▾		▼			▼			▼						
Algorithm	Avg. Rank	Tsukuba ground truth			Venus ground truth			Teddy ground truth				Cones ground truth		
	▼	nonocc	all	disc	nonocc	all	disc	nonocc	all	disc	nonocc	all	disc	
ADCensus [82]	10.9	1.07 ¹⁶	1.48 ¹³	5.73 ¹⁹	0.09 ²	0.25 ⁹	1.15 ²	4.10 ¹³	6.22 ⁶	10.9 ¹¹	2.42 ¹⁴	7.25 ¹¹	6.95 ¹⁵	3.97
AdaptinqBP [16]	14.2	1.11 ¹⁹	1.37 ⁸	5.79 ²¹	0.10 ⁴	0.21 ⁶	1.44 ⁶	4.22 ¹⁵	7.06 ¹²	11.8 ¹⁵	2.48 ¹⁸	7.92 ²³	7.32 ²³	4.23
CoopRegion [39]	14.8	0.87 ⁴	1.16 ¹	4.61 ⁴	0.11 ⁵	0.21 ⁶	1.54 ¹⁰	5.16 ²⁸	8.31 ¹⁶	13.0 ²¹	2.79 ³⁵	7.18 ¹⁰	8.01 ⁴⁰	4.41
RDP [87]	19.2	0.97 ⁹	1.39 ¹⁰	5.00 ⁹	0.21 ³⁴	0.38 ²⁴	1.89 ¹⁹	4.84 ¹⁸	9.94 ²⁸	12.6 ¹⁸	2.53 ²¹	7.69 ¹⁶	7.38 ²⁴	4.57
MultiRBF [129]	19.6	1.33 ⁴¹	1.56 ¹⁷	6.02 ²⁸	0.13 ⁸	0.17 ²	1.84 ¹⁶	5.09 ²⁴	6.36 ⁷	13.4 ²⁵	2.90 ⁴²	6.76 ⁵	7.10 ²⁰	4.39
DoubleBP [34]	20.0	0.88 ⁶	1.29 ⁵	4.76 ⁷	0.13 ⁹	0.45 ⁴¹	1.87 ¹⁸	3.53 ¹⁰	8.30 ¹⁵	9.63 ⁶	2.90 ⁴¹	8.78 ⁵⁰	7.79 ³²	4.19
MDPM [140]	20.3	1.15 ²⁰	1.59 ²⁰	6.14 ³¹	0.14 ¹⁴	0.36 ²²	1.52 ⁹	3.79 ¹¹	5.78 ⁴	11.1 ¹³	2.74 ²⁹	8.38 ³⁶	7.91 ³⁵	4.22
OutlierConf [40]	20.5	0.88 ⁵	1.43 ¹²	4.74 ⁶	0.18 ²³	0.26 ¹¹	2.40 ³⁴	5.01 ²⁰	9.12 ²⁴	12.8 ²⁰	2.78 ³⁴	8.57 ⁴¹	6.99 ¹⁶	4.60
AdaptiveGF [127]	24.1	1.04 ¹²	1.53 ¹⁴	5.62 ¹⁴	0.17 ²²	0.41 ³²	1.98 ²²	5.71 ³⁵	11.3 ⁴¹	14.3 ³²	2.44 ¹⁶	8.22 ³⁰	7.05 ¹⁹	4.98

KITTI benchmark

Image size 1241 x 376 (< 0.5 MP), D ≈ 70...150

The KITTI Vision Benchmark Suite

A project of Karlsruhe Institute of Technology and Toyota Technological Institute at Chicago



home setup **stereo** flow odometry object tracking road raw data submit results jobs

Andreas Geiger (MPI Tübingen) | Philip Lenz (KIT) | Christoph Stiller (KIT) | Raquel Urtasun (University of Toronto)

Stereo Evaluation

Rank	Method	Setting	Code	Out-Noc	Out-All	Avg-Noc	Avg-All	Density	Runtime	Environment	Compare
1	SceneFlow			2.98 %	3.97 %	0.8 px	1.0 px	100.00 %	35 s	1 core @ 3.5 Ghz (C/C++)	<input type="checkbox"/>
Anonymous submission											
2	PCBP-SS			3.40 %	4.72 %	0.8 px	1.0 px	100.00 %	5 min	4 cores @ 2.5 Ghz (Matlab + C/C++)	<input type="checkbox"/>
K. Yamaguchi, D. McAllester and R. Urtasun: Robust Monocular Epipolar Flow Estimation . CVPR 2013.											
3	gtRF-SS			3.83 %	4.59 %	0.9 px	1.0 px	100.00 %	1 min	1 core @ 2.5 Ghz (Matlab + C/C++)	<input type="checkbox"/>
Anonymous submission											
4	StereoSLIC			3.92 %	5.11 %	0.9 px	1.0 px	99.89 %	2.3 s	1 core @ 3.0 Ghz (C/C++)	<input type="checkbox"/>
K. Yamaguchi, D. McAllester and R. Urtasun: Robust Monocular Epipolar Flow Estimation . CVPR 2013.											
5	PR-Sf+E			4.02 %	4.87 %	0.9 px	1.0 px	100.00 %	200 s	4 cores @ 3.0 Ghz (Matlab + C/C++)	<input type="checkbox"/>
C. Vogel, S. Roth and K. Schindler: Piecewise Rigid Scene Flow . International Conference on Computer Vision (ICCV) 2013.											
6	PCBP			4.04 %	5.37 %	0.9 px	1.1 px	100.00 %	5 min	4 cores @ 2.5 Ghz (Matlab + C/C++)	<input type="checkbox"/>
K. Yamaguchi, T. Hazan, D. McAllester and R. Urtasun: Continuous Markov Random Fields for Robust Stereo Estimation . ECCV 2012.											

Limitations

- 100+ new algorithms published (benchmarked) since 2002.
 - Middlebury: focus on accuracy
 - KITTI: focus on robust performance

Neither require high accuracy on hi-res images

- Most of the existing methods do not scale.

Limitations

- Searching full disparity space requires $O(P \cdot D)$ time
= $O(s^3)$ for image size s

Middlebury Teddy	KITTI
0.16MP x 60 = 10 Mdisp.	0.5MP x 80 = 40 Mdisp.
Middlebury New	Disney Mansion
6MP x 256 = 1.5 Gdisp.	19MP x 1000 = 19 Gdisp.

New datasets and benchmarks

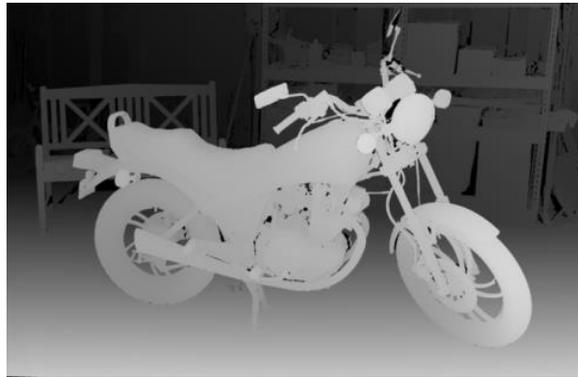
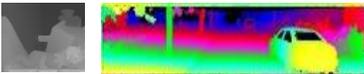
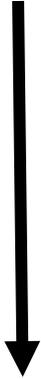
- Need more challenging datasets for algorithm design
- Middlebury Stereo Eval v.3 is underway *
- 30 new datasets
 - a subset is discussed here -- “Middlebury New 7”

* The benchmark is not public yet.

If interested, contact Daniel Scharstein (schar@middlebury.edu)

Middlebury
Teddy

Middlebury
New



KITTI

Disney →

Scalability

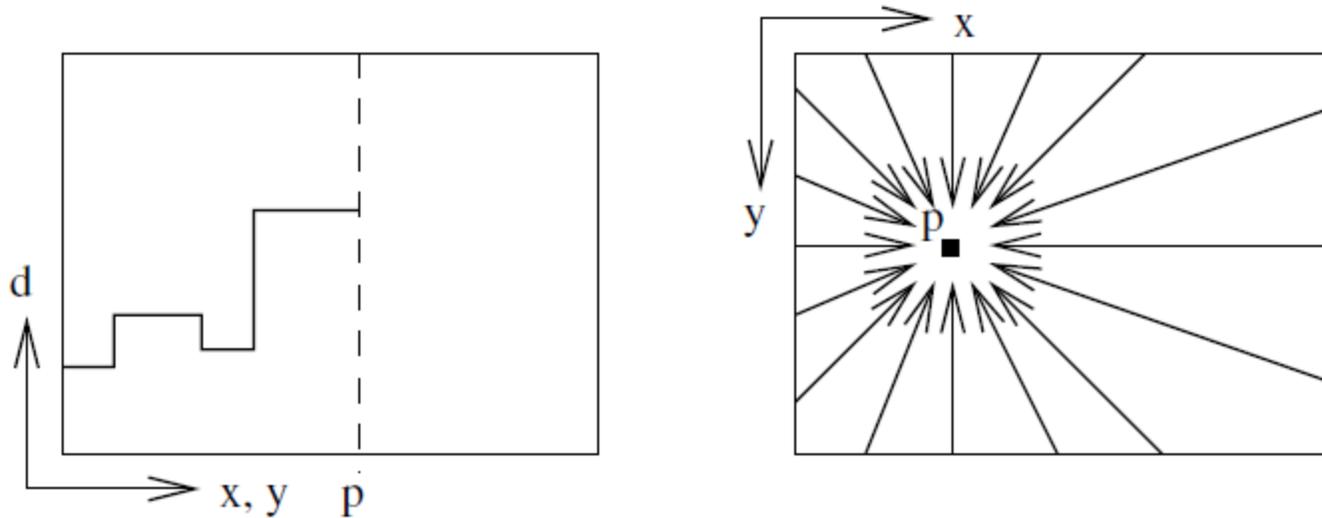
- (Most) existing methods are $O(P \cdot D)$, or $O(P \cdot D^2)$
 - These methods do not scale
 - Ideally, want $O(P)$
- P: pixels
D: disparities
- Key Issues:
 - Does higher resolution even help?
 - Does it make sense to enumerate disparities ?
 - Coarse to fine strategies ?
 - Practical optimization techniques

Promising Approaches

- Efficient approximate energy minimization
 - Only run at low-res; Upsample and refine disparities
[Ferst et. al. 2013, Ma et. al. 2013]
 - Semi-global Matching (SGM)
[Hirschmüller 2005]
- Avoid exploring the whole DSI
 - Coarse-to-fine [long tradition]
 - Seed & Grow [Cech & Sara 2007, ...]
 - PatchMatch stereo [Bleyer et al. 2011]
 - ELAS [Geiger et al. 2010]
 - Local Plane Sweep Stereo [this talk]

Semi-global Matching

[Hirschmüller 2005]



- Aggregate along 1D minimum cost paths ending at pixel p
- Only cost of this path needed; not the path itself
- Efficiently computed via message-passing.
- Winner-take-all disparity selection

Efficient High-Resolution Stereo Matching using Local Plane Sweeps

Will be presented at CVPR 2014



Daniel Scharstein

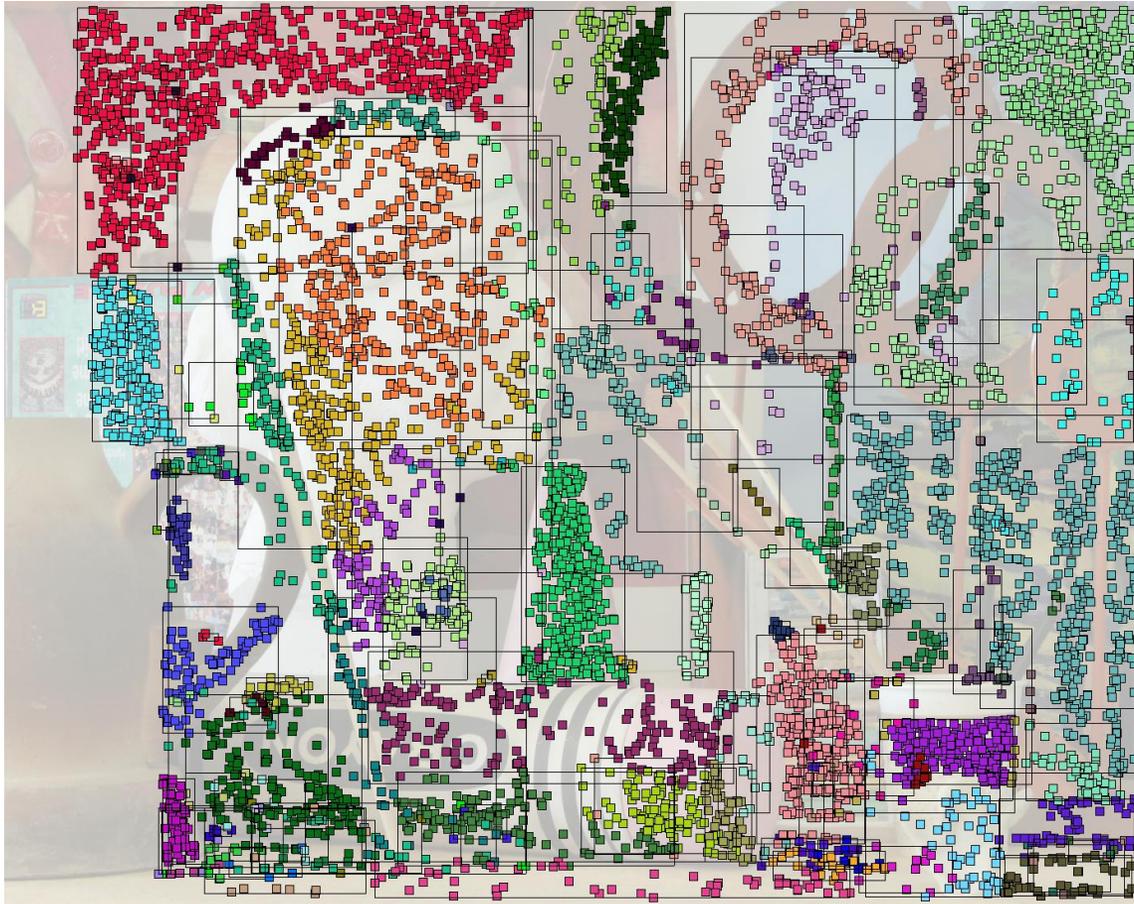


Rick Szeliski

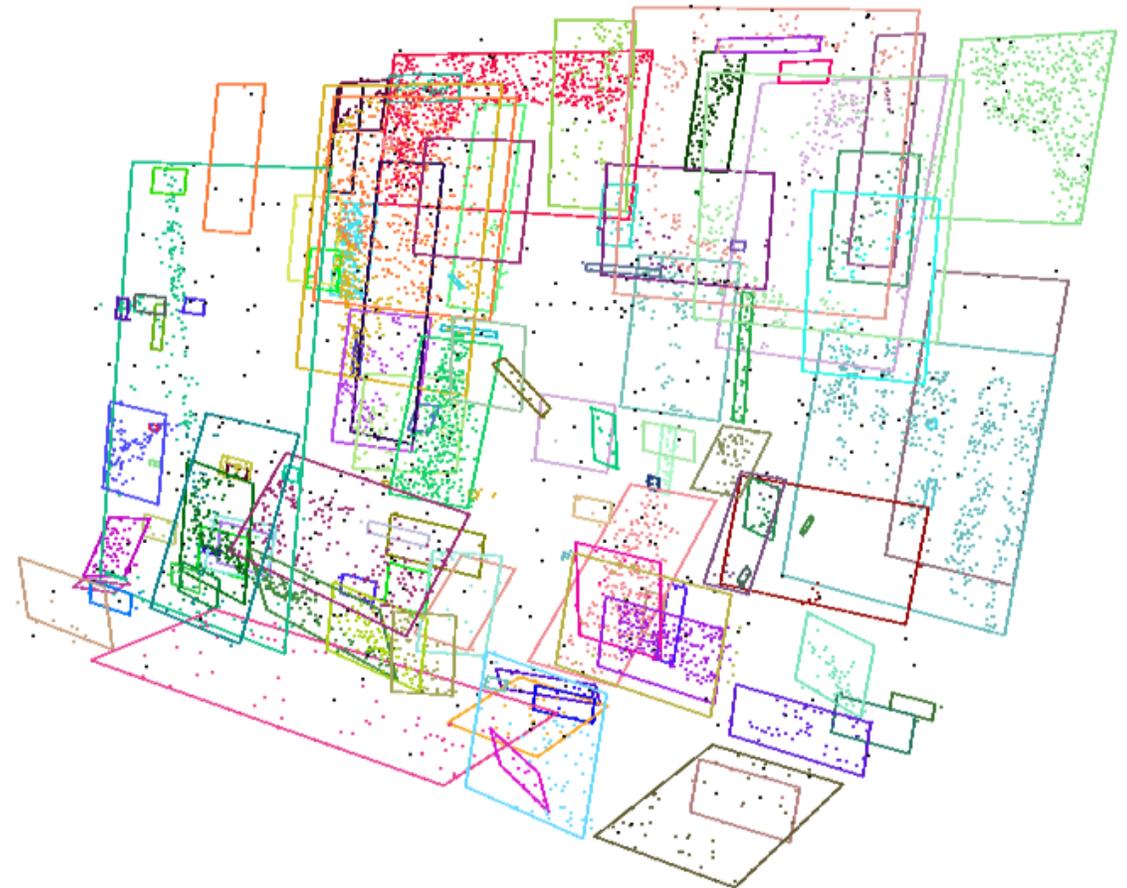
Local Plane Sweep Stereo

- Sparse feature matching; refine vertical disparities
- Generate plane hypotheses (with unknown extents)

Plane hypothesis generation



Plane Labels



Local Plane Sweep Stereo

- Sparse feature matching; refine vertical disparities
- Generate plane hypotheses (with unknown extents)
- Perform local plane sweeps (LPS) around hypothesized planes
 - local stereo problem with narrow disparity range; solved using SGM

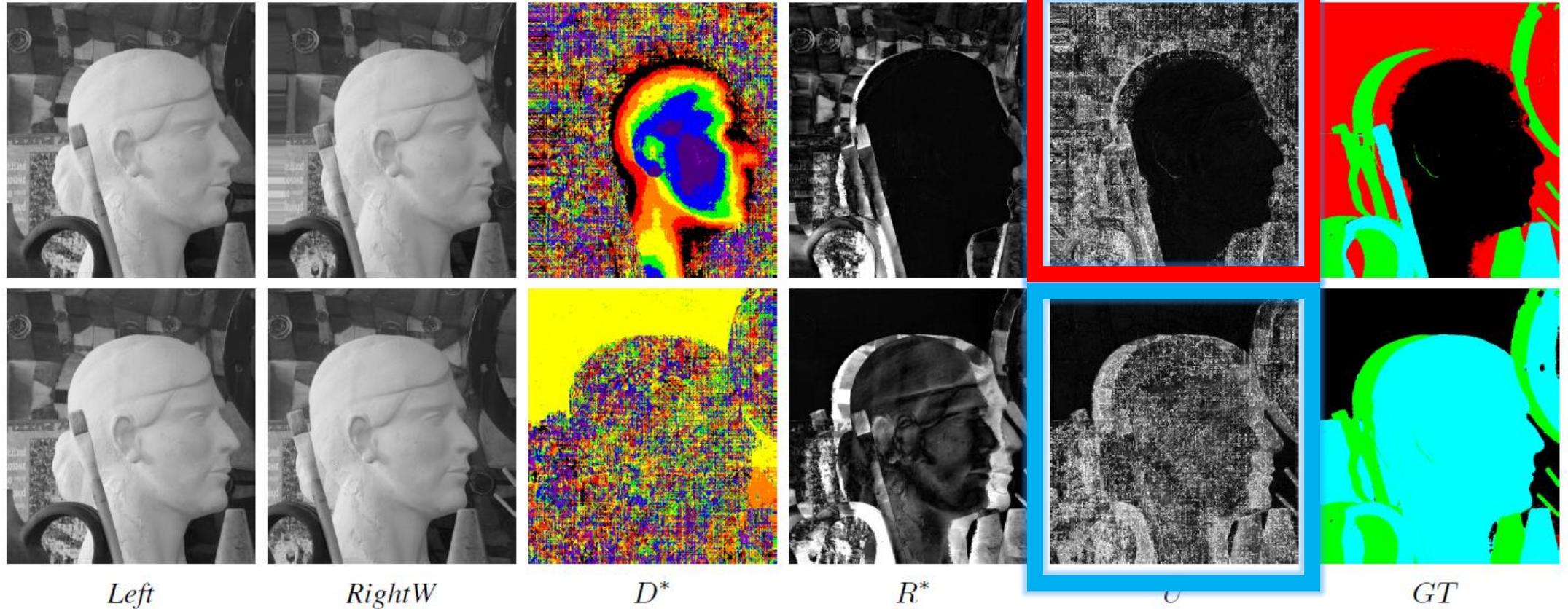
Local Plane Sweeps

Plane 2

Local Plane Sweep solution

Cost Map

In-range
disparities
(ground truth)



Plane 1

Local Plane Sweep solution

d

-3

-2

-1

0

1

2

3

Local Plane Sweep Stereo

- Sparse feature matching; refine vertical disparities
- Generate plane hypotheses (with unknown extents)
- Perform local plane sweeps (LPS) around hypothesized planes
 - local stereo problem with narrow disparity range; solved using SGM
- Tile structure
 - Perform LPS on tiles and propagate planes to adjoining tiles

Proposal Propagation

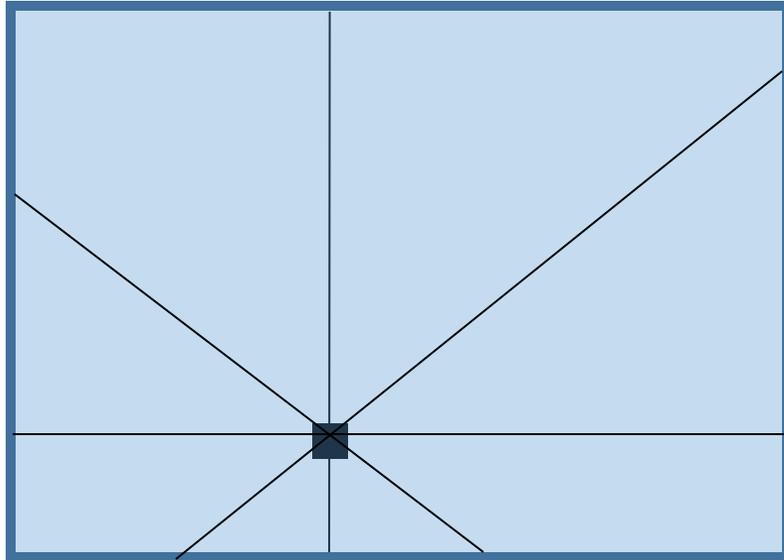
- Initial proposals:
 - planes π with feature points per tile
- Repeat $nR=3$ times:
 - Compute surfaces via local plane sweeps per tile
 - Update winner-take-all label map (L_{WTA})
 - Use L_{WTA} to predict well-supported plane proposals
 - Propagate these planes to the neighboring tiles

Local Plane Sweep Stereo

- Sparse feature matching; refine vertical disparities
- Generate plane hypotheses (with unknown extents)
- Perform local plane sweeps (LPS) around hypothesized planes
 - local stereo problem with narrow disparity range; solved using SGM
- Tile structure
 - Perform LPS on tiles and propagate planes to adjoining tiles
- Global optimization
 - Assign pixels to surface proposals
 - Fast approximate energy minimization (via SGM)
 - Extend SGM to exploit tile structure and sparse label sets

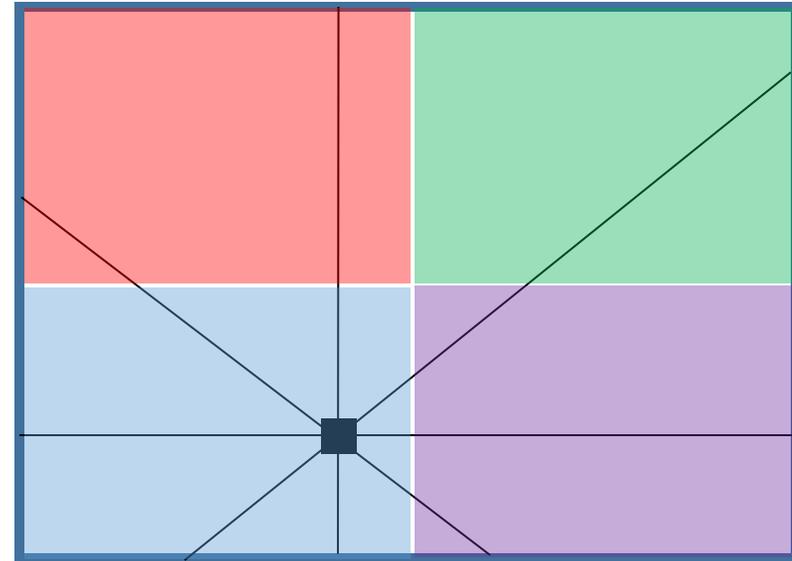
Global Optimization (via SGM)

- Message passing on 1D paths (8 directions)



SGM

- fixed label set at all pixels



LPS

- Label sets vary across tiles

Disparity Selection

- Original SGM:
 - Aggregated costs
 - WTA at every pixel
- LPS:
 - Aggregated costs
 - Top-m candidates at every pixel ($m = 2$)
 - Median filter on candidates within a small window.

Experiments

- Evaluation:

- - PatchMatch Stereo [Bleyer et al. 2011]

- - SGM (our impl.)

- - SGM-HH [Hirschmüller 2005]

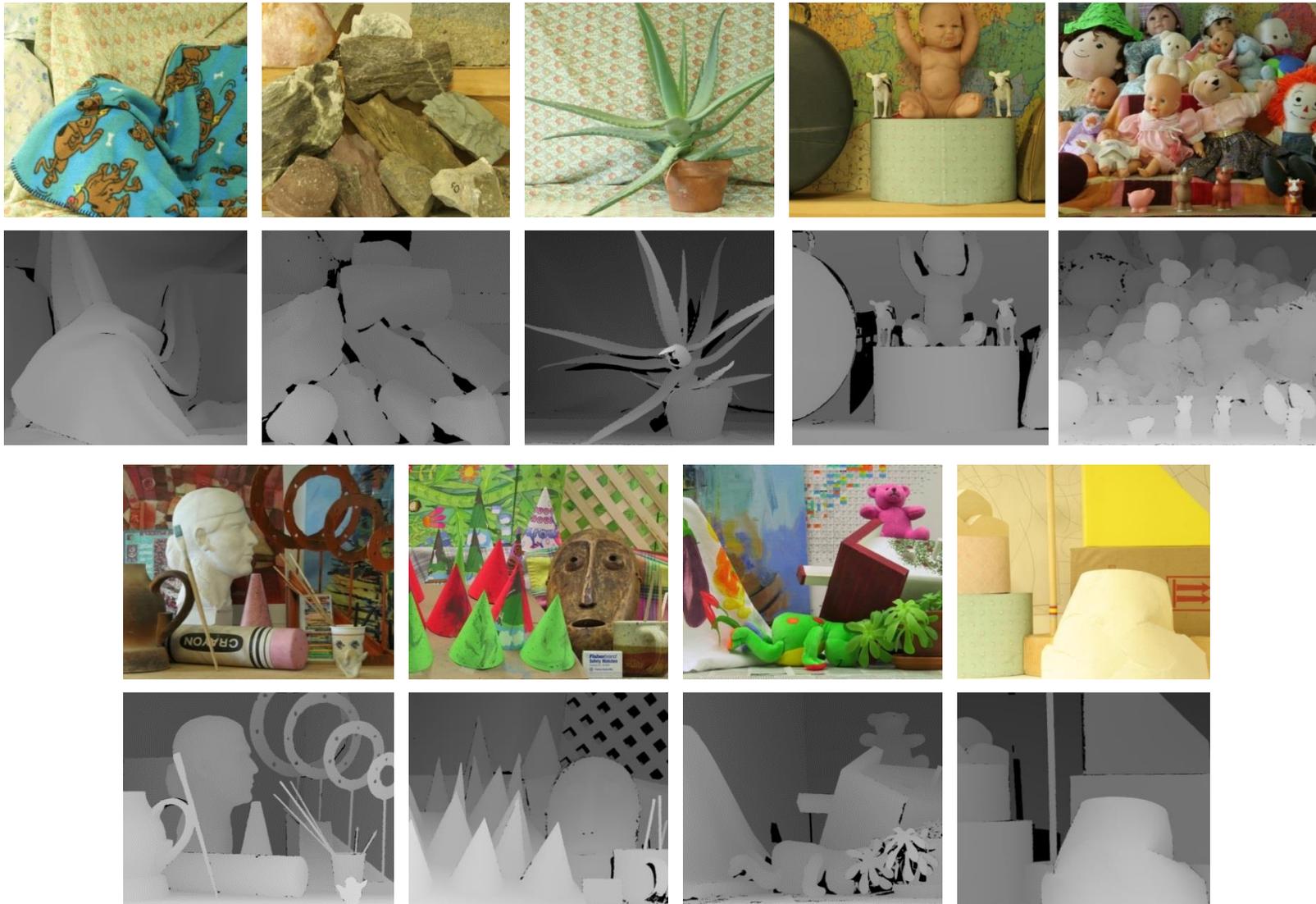
- - ELAS [Geiger et al. 2010]

- - LPS

- Metric:

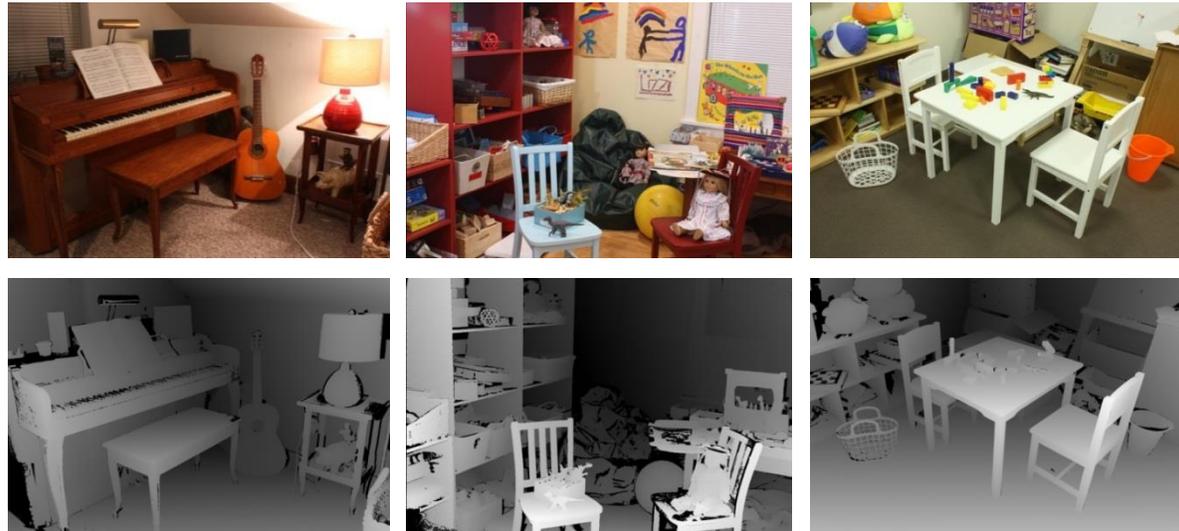
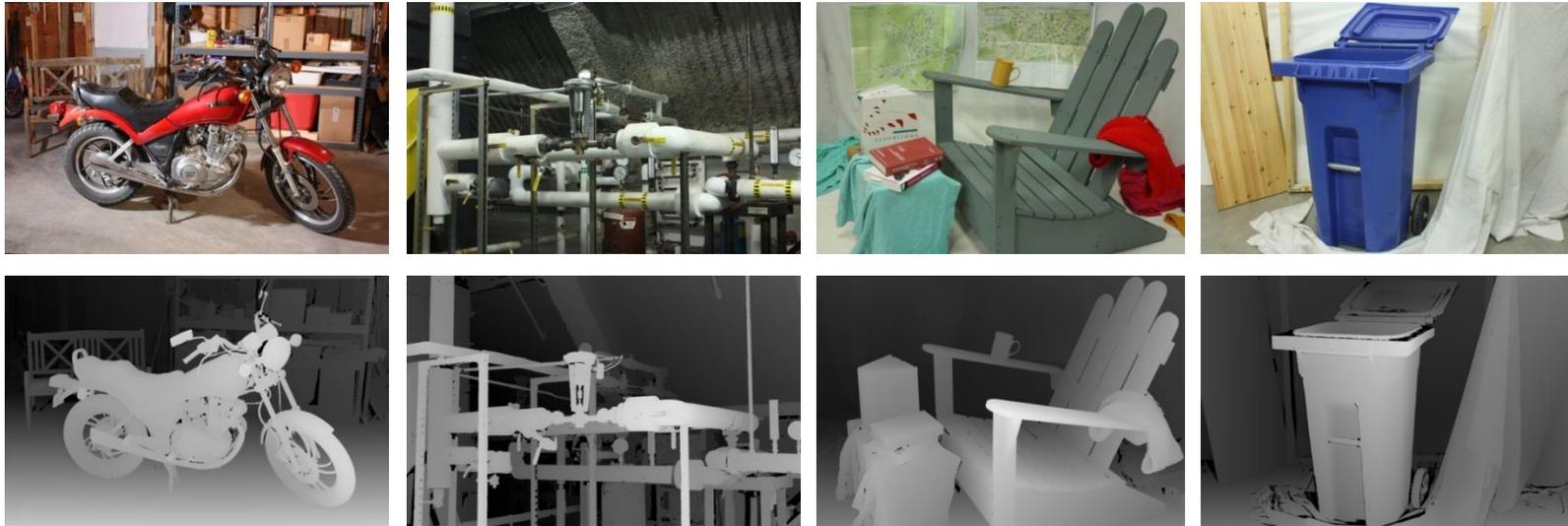
- 1 and 2 pixel disparity error at non-occluded pixels.

Midd9 (1.4–2.7 MP)



Subset of full-resolution 2003-2006 Middlebury datasets

MidNew7 (5.1–6.0 MP)



Subset of new 2011-2014 Middlebury training sets

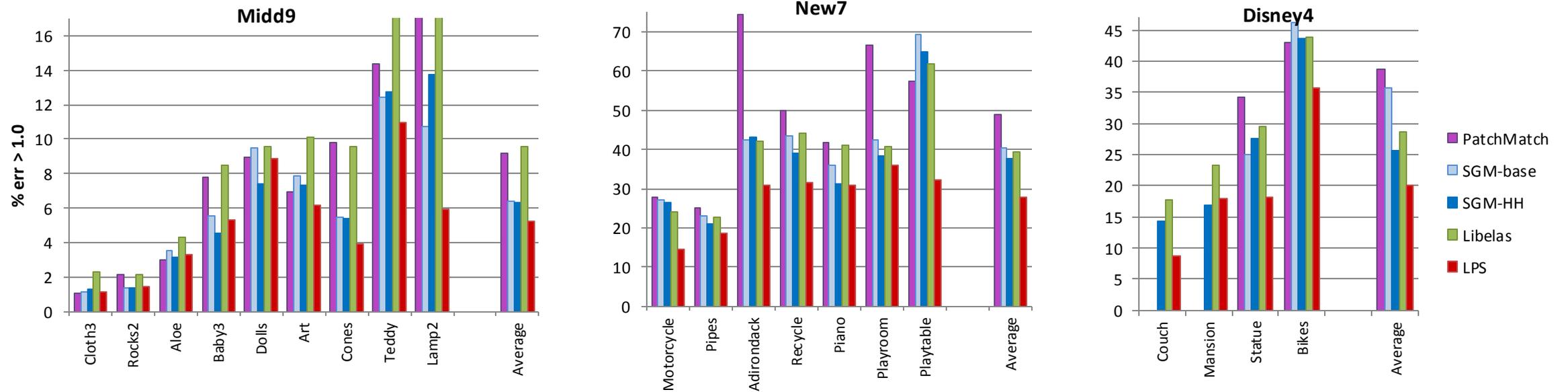
Disney4 (4.5–19 MP)



C. Kim, H. Zimmer, Y. Pritch, A. Sorkine-Hornung, and M. Gross
Scene reconstruction from high spatio-angular resolution light fields
SIGGRAPH 2013

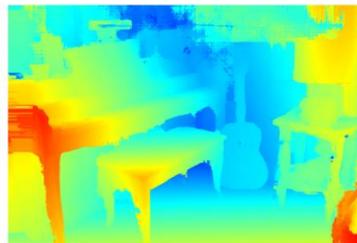
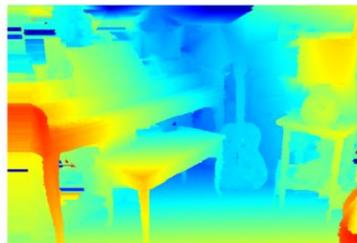
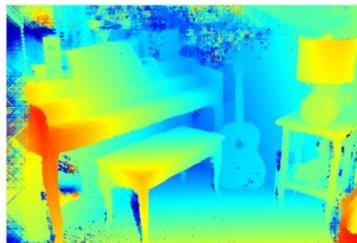
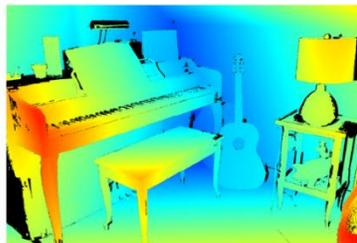
* We treat their results (computed from 100 images) as GT

Results – accuracy

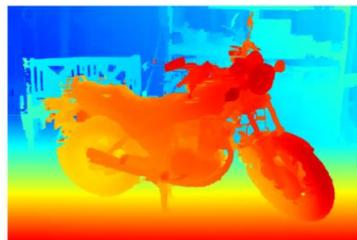
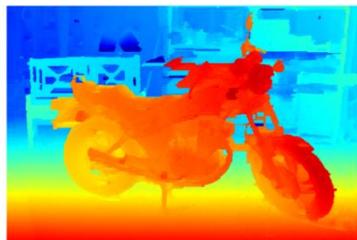
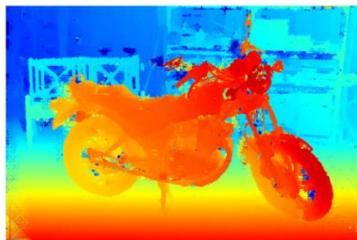


Bad pixels (%), thresh = 1.0 pixel

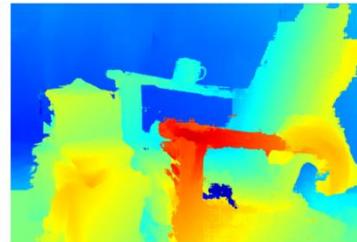
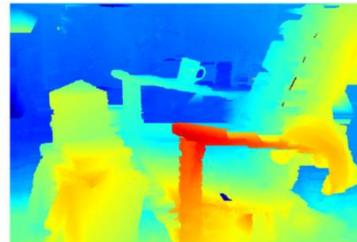
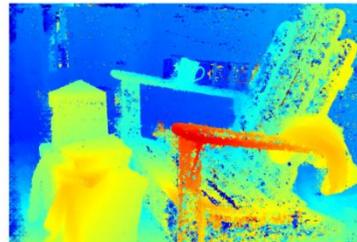
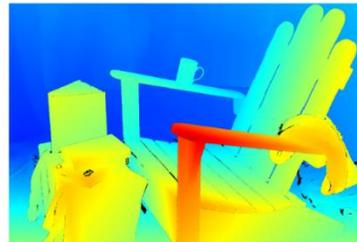
Piano



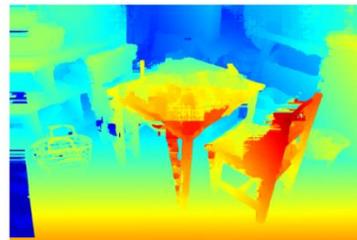
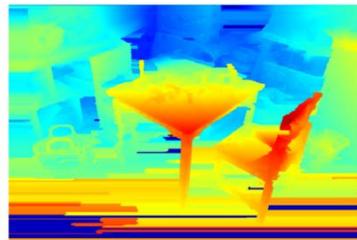
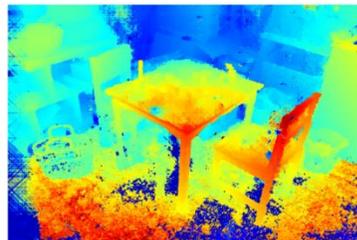
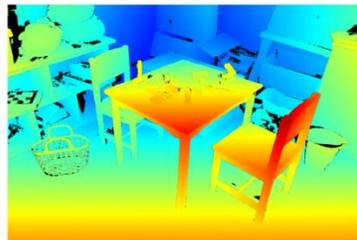
Motorcycle



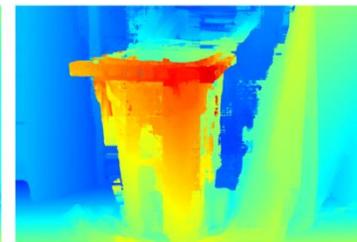
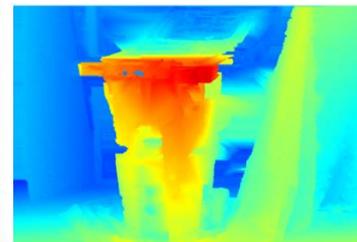
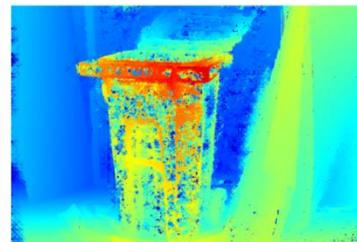
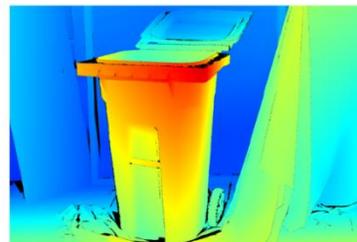
Adirondack



Playtable



Recycle



Left Image Ground Truth

SGM

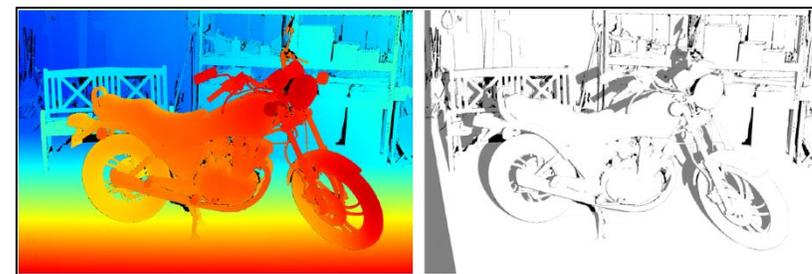
ELAS

LPS (ours)

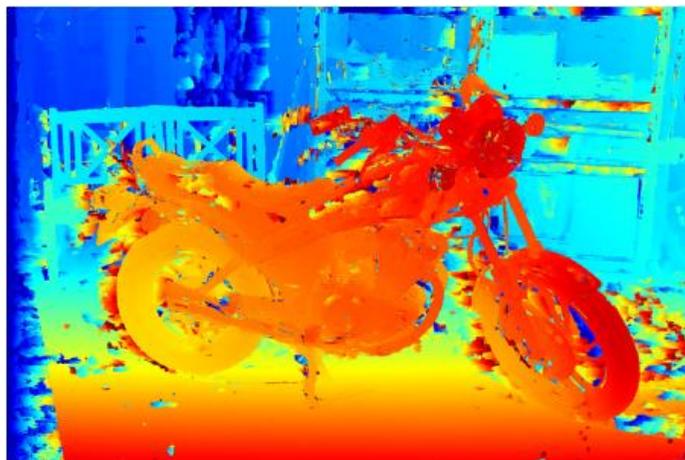
Error maps (Motorcycle)

Ground Truth

Occlusion Mask



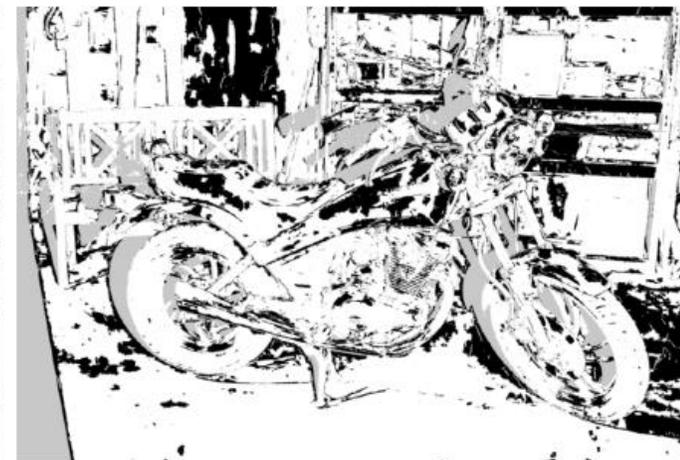
PatchMatch



3330 seconds



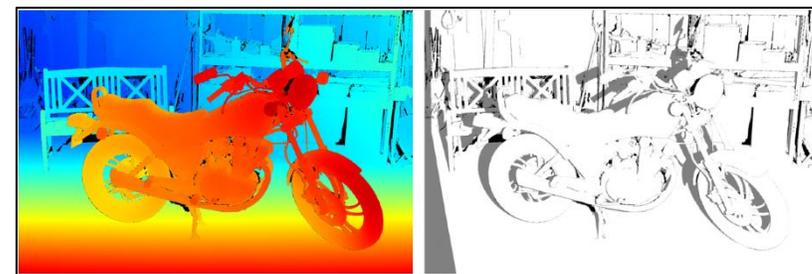
err1 = 33.8 %



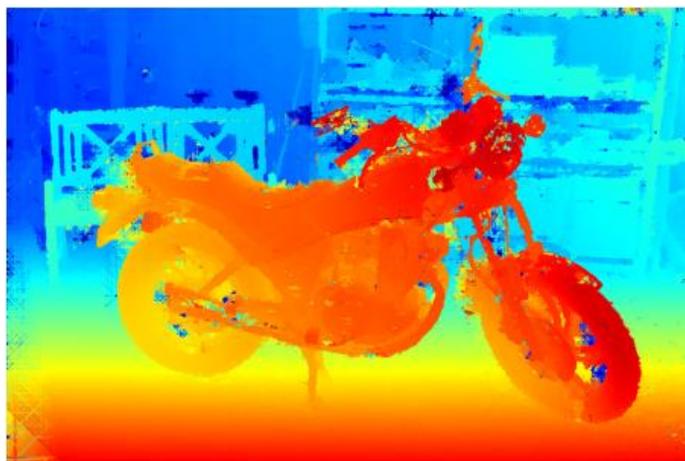
err2 = 24.2 %

Error maps (Motorcycle)

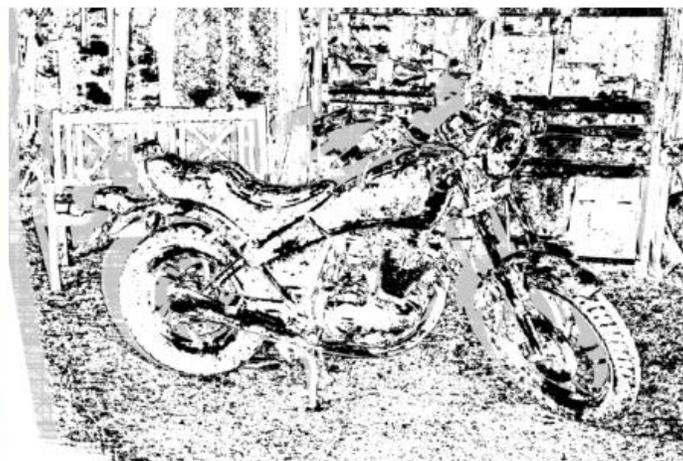
Ground Truth Occlusion Mask



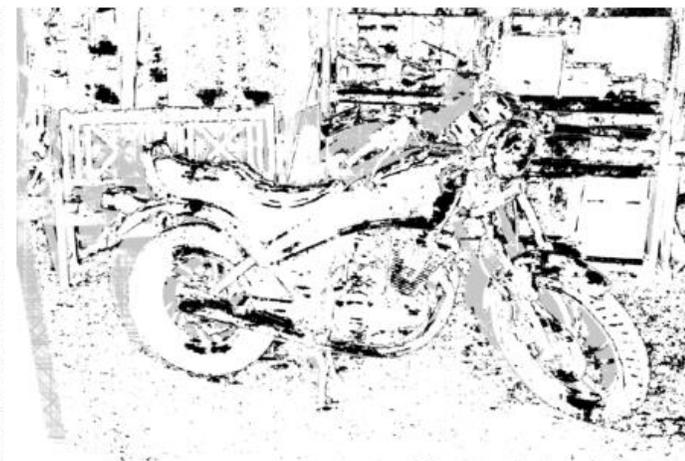
SGM



51.4 seconds



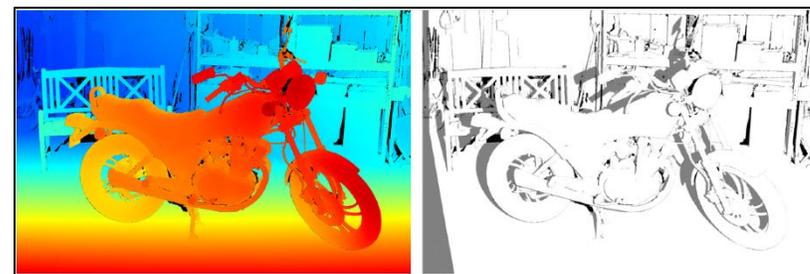
err1 = 29.3 %



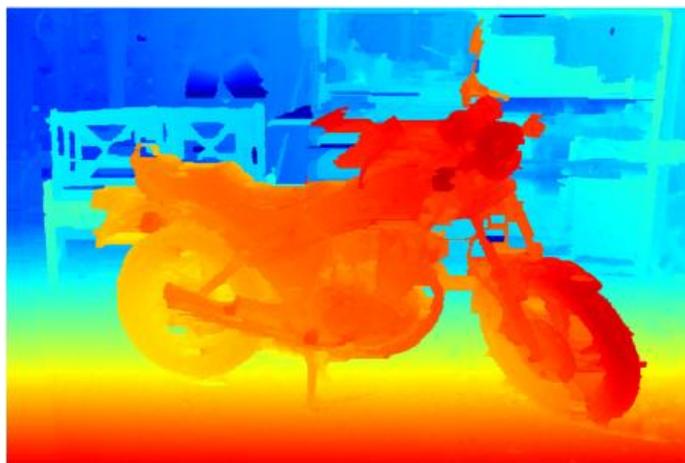
err2 = 15.1 %

Error maps (Motorcycle)

Ground Truth Occlusion Mask



ELAS



5.0 seconds



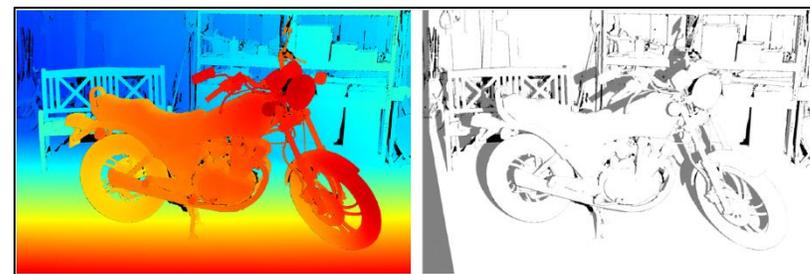
err1 = 34.0 %



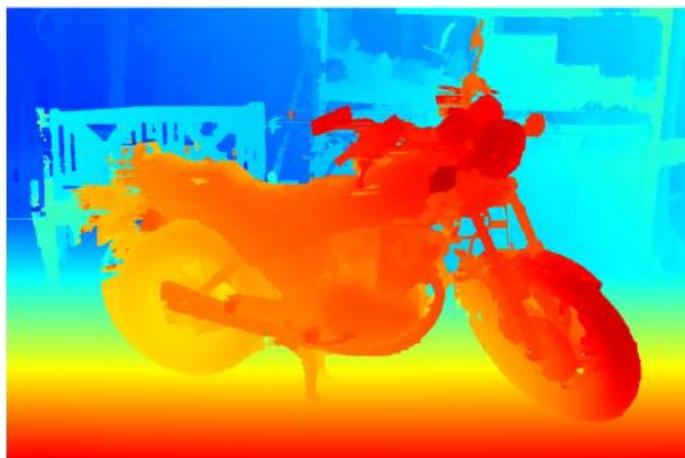
err2 = 19.1 %

Error maps (Motorcycle)

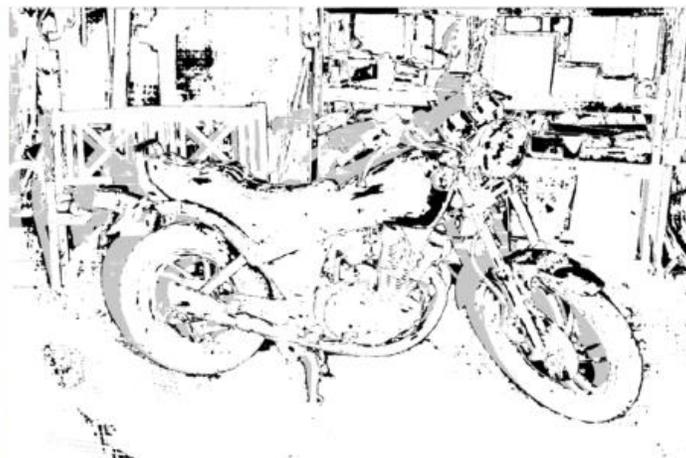
Ground Truth Occlusion Mask



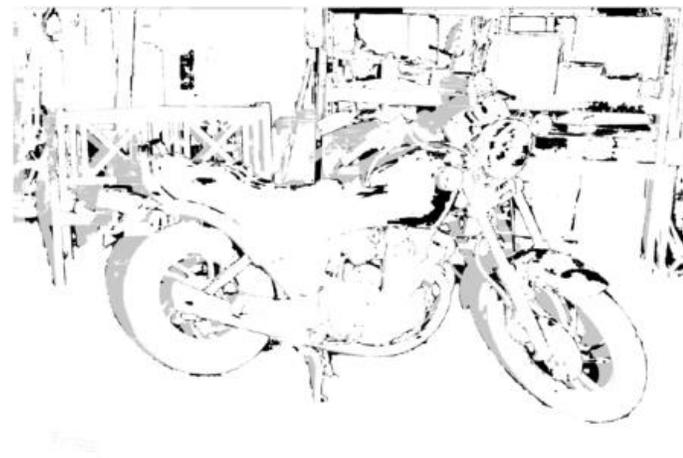
LPS



9.6 seconds

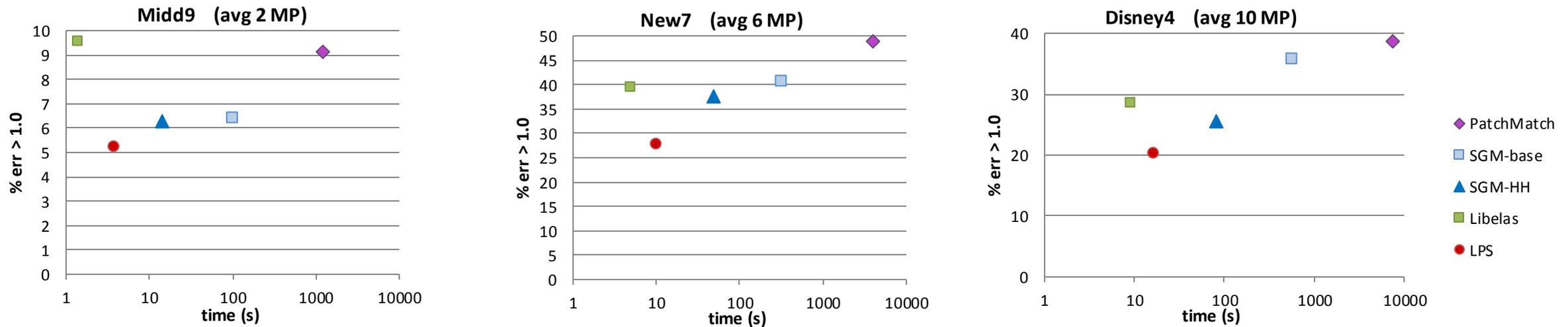


err1 = 12.2 %



err2 = 6.5 %

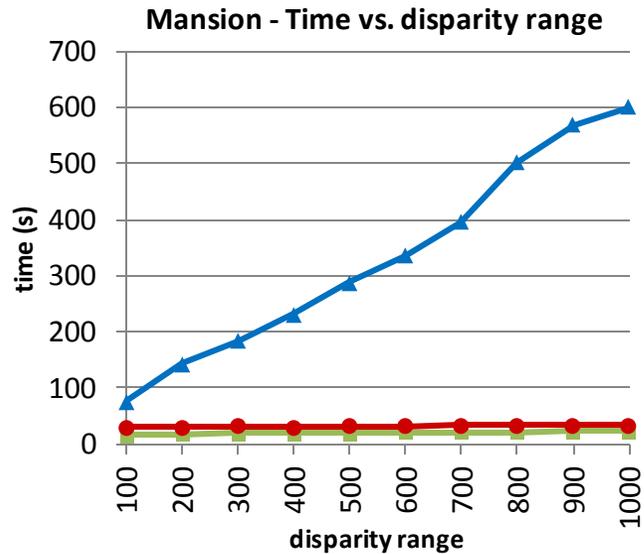
Results – Accuracy vs. Runtime



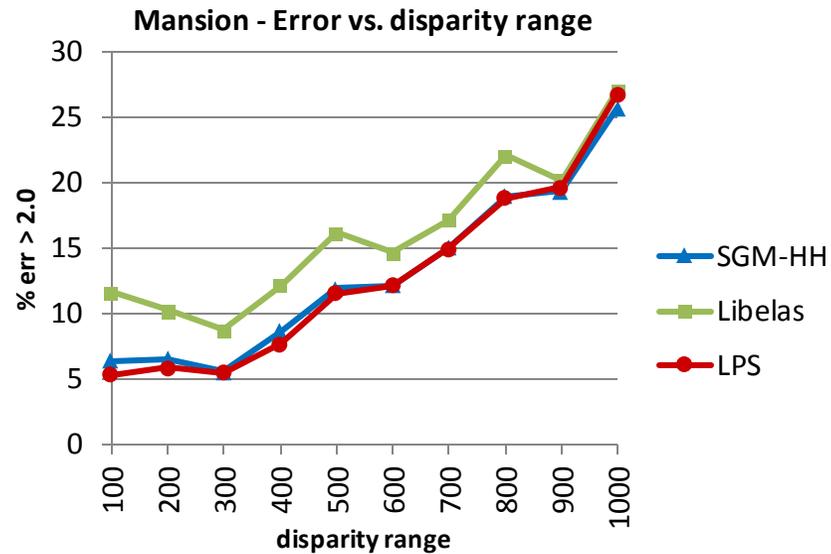
Avg. error vs. runtime, thresh = 1.0 pixel

Results – Scalability

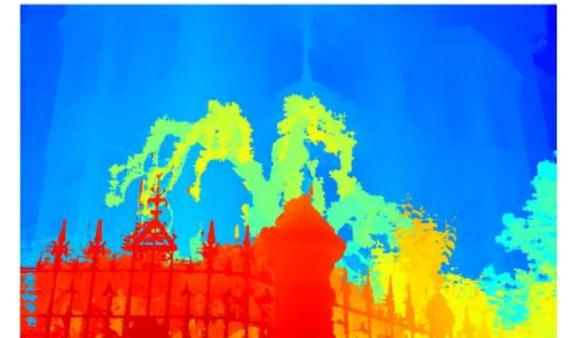
- With increasing disparity range



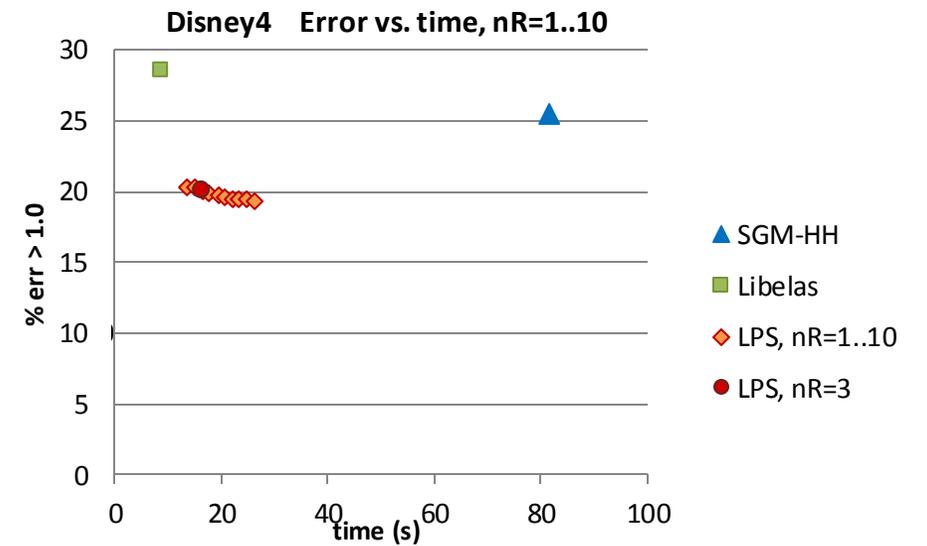
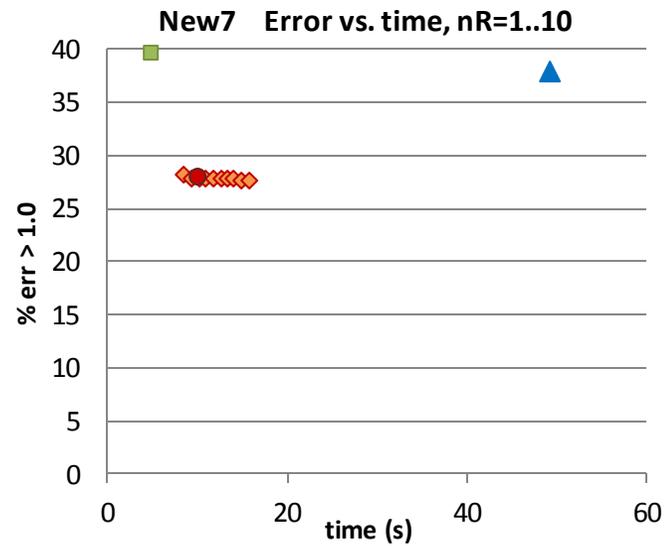
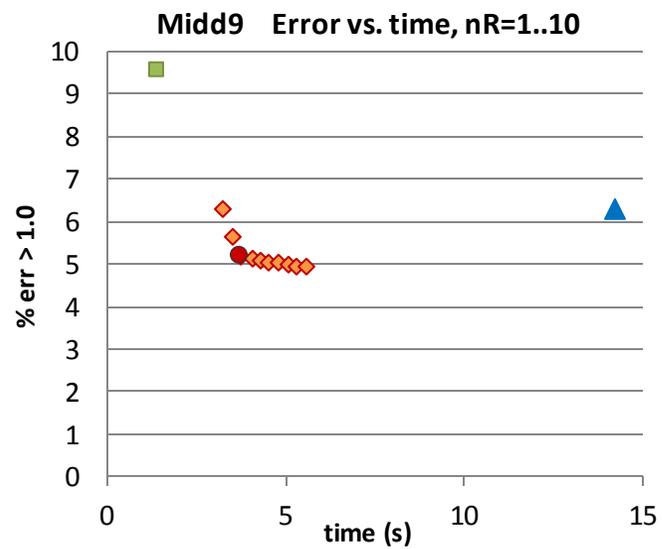
Runtime



Error



Results – more proposals



Avg. error vs. number of rounds

LPS – Summary

- Plane proposals from matched features
- Refine into surfaces using LPS
- Propagate promising planes
- Final pixel assignment to surfaces using global optimization

LPS – Benefits

- Doesn't explore full search space
- Runtime independent of disparity range
- No fronto-parallel bias
- Excellent recovery of slanted surfaces, even with weak texture
- Easy integration of vertical disparity correction

LPS – Limitations

- Can miss surfaces if not among initial proposals
- Cannot handle completely untextured surfaces
- Need “stopping criterion” for proposal generation
- No occlusion reasoning

Future work

- Other types of proposals
 - Hierarchical (coarse to fine)
 - Line, edge features
 - New proposals via “residual analysis”
- Add color models and occlusion reasoning into final pixel assignment
- Better modeling of calibration errors

Promising directions

- Moving away from monolithic optimization
- Global pixel-to-surface assignment (like segmentation; not matching)
- Local stereo matching for proposal generation
- Residual analysis to guide additional search

Surface-based stereo matching

- Piecewise planar stereo



Birchfield and Tomasi 2001

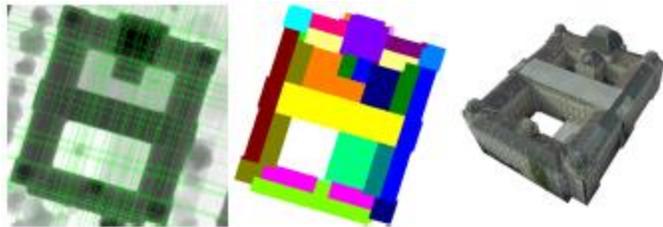


Furukawa et al. 2008



Sinha et al. 2009

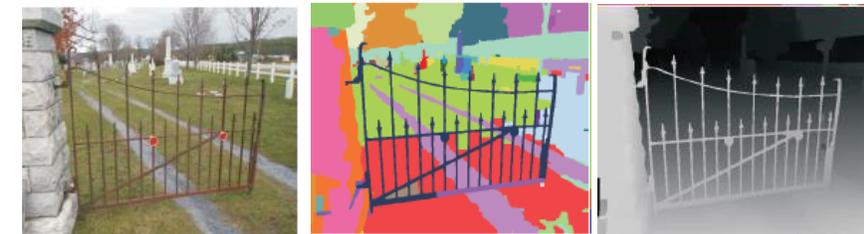
- Surface stereo



Zebedin et al. 2008



Gallup et al. 2010



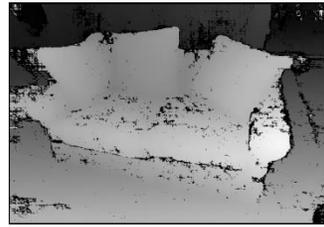
Bleyer et al. 2010, 2011

Multiple View Object Cosegmentation using Appearance and Stereo

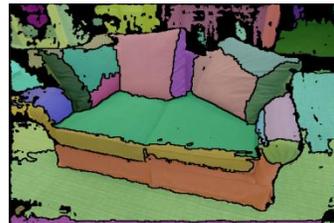
Adarsh Kowdle, Sudipta N. Sinha and Rick Szeliski (ECCV 2012)



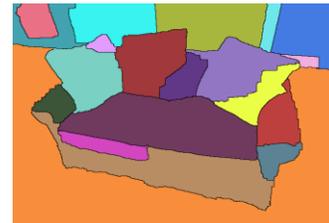
Input images



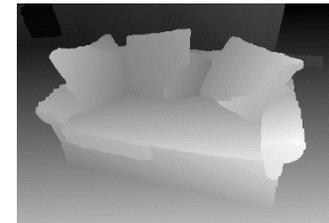
Stereo matching



Plane hypotheses



Plane labels



Depth map

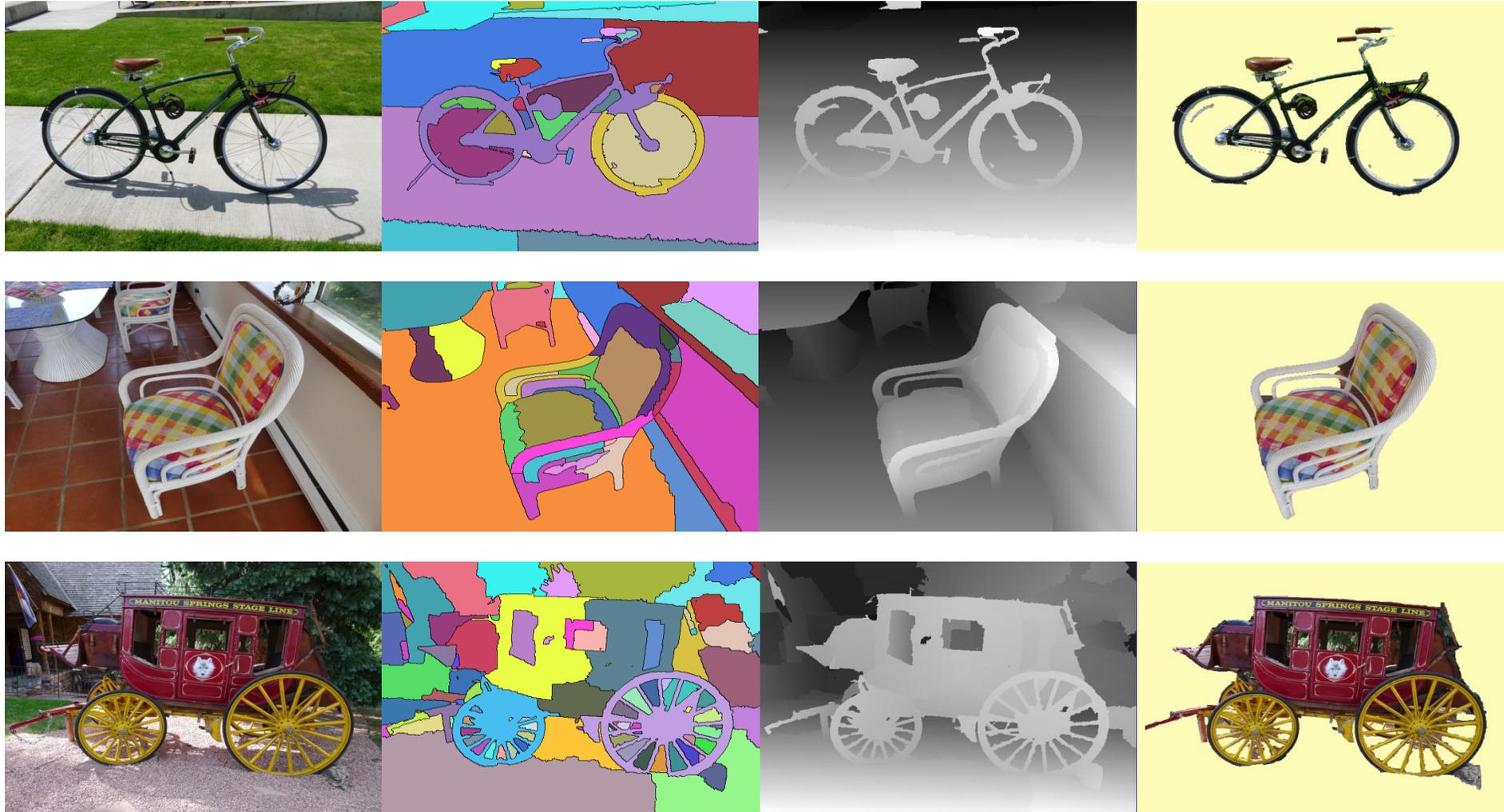


Using multiple views, infer what constitutes the foreground *object*



Multiple View Object Cosegmentation using Appearance and Stereo

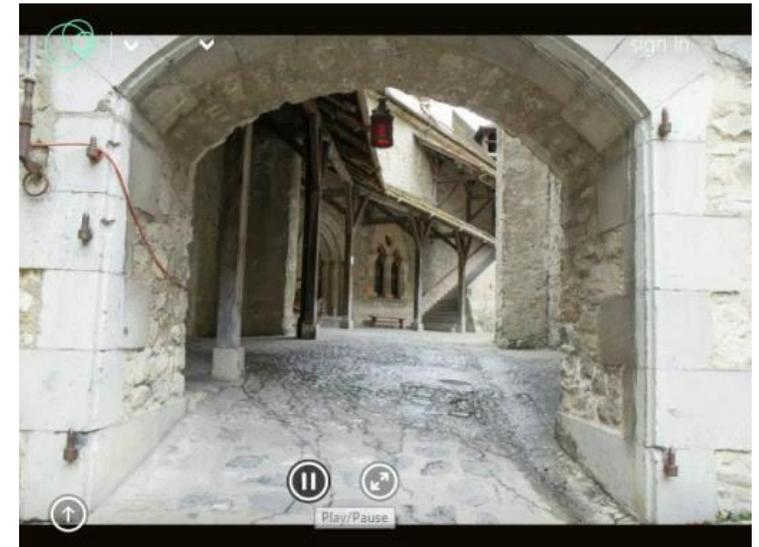
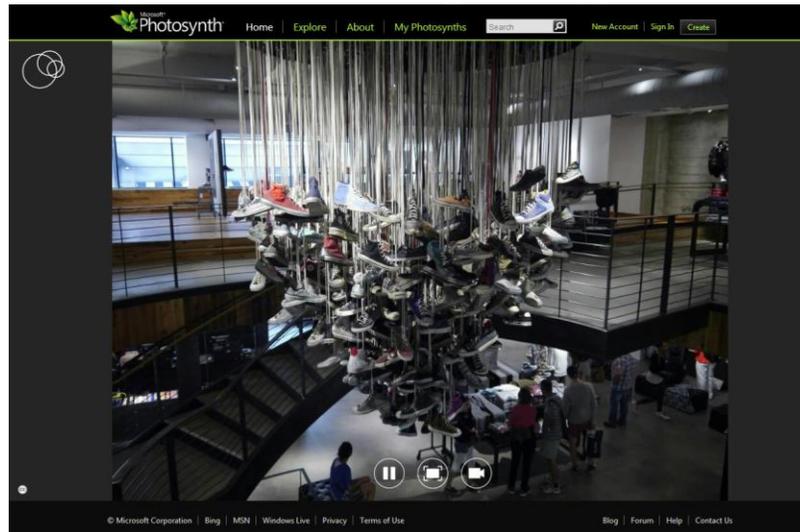
Adarsh Kowdle, Sudipta N. Sinha and Rick Szeliski (ECCV 2012)



Application: Image-based Rendering

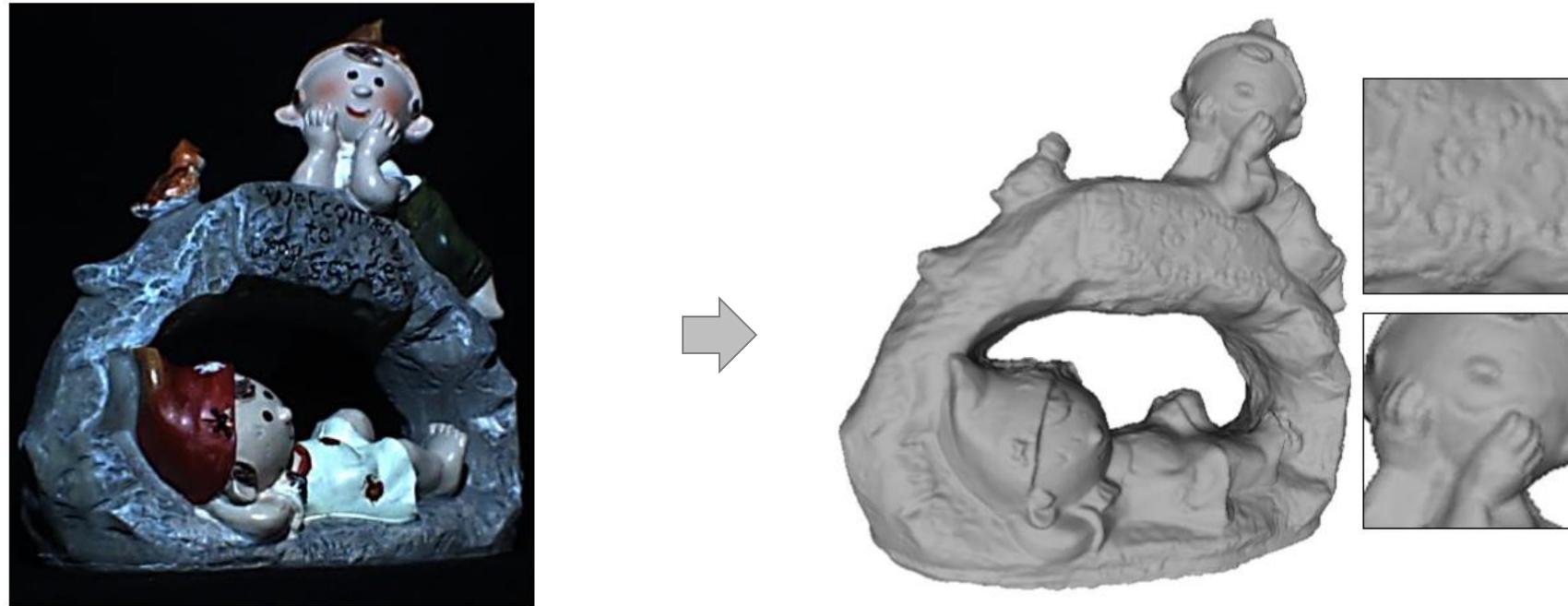
Photosynth 2 (www.photosynth.net)

- Capture the world in 3D;
- Novel view synthesis
- Interactive viewer



Multiview Photometric Stereo using Planar Mesh Parameterization

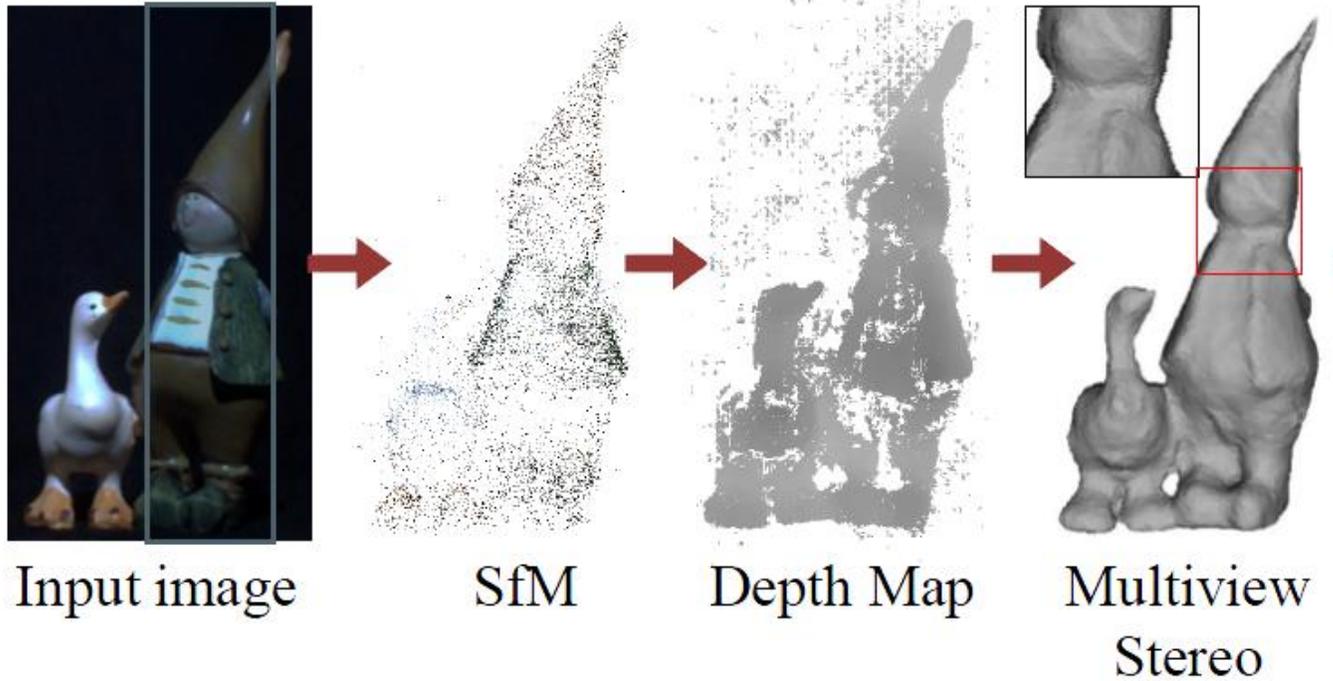
Jaesik Park, Sudipta N. Sinha, Yasuyuki Matsushita, Yu-Wing Tai and In So Kweon (ICCV 2013)



- Automatic 3D reconstruction from RGB images
- Multi-view stereo gives coarse shape
- Multi-view Photometric stereo refines shape

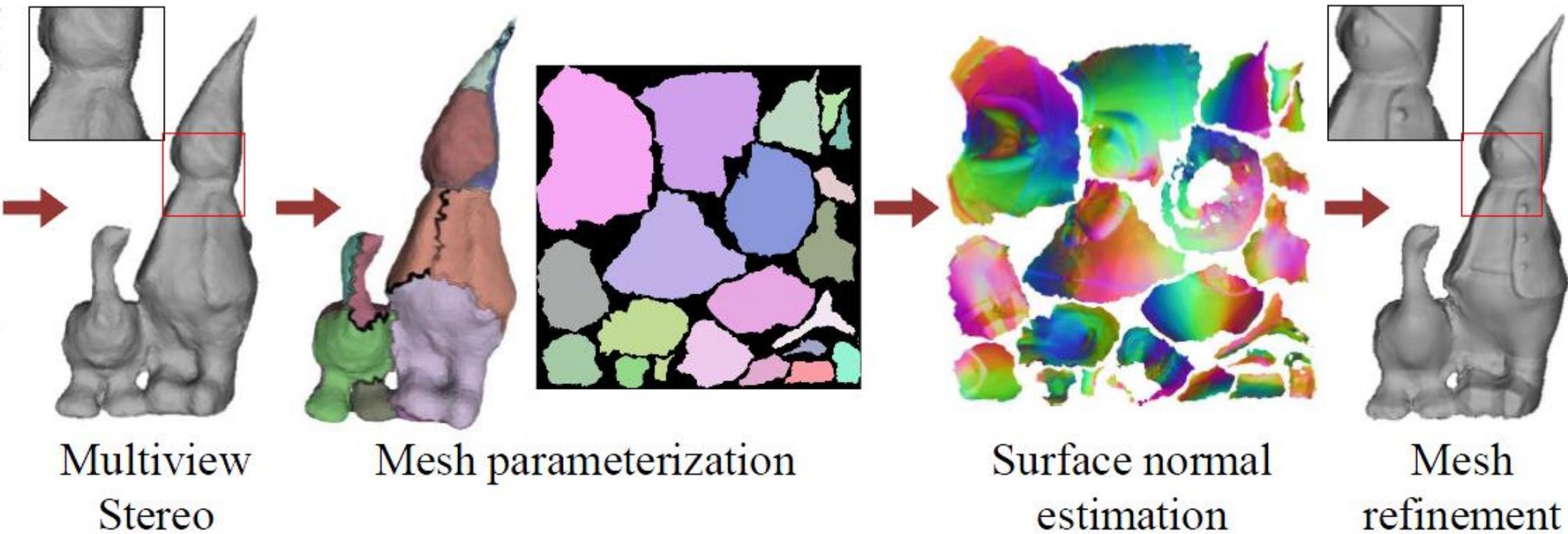
Multiview Photometric Stereo using Planar Mesh Parameterization

Jaesik Park, Sudipta N. Sinha, Yasuyuki Matsushita, Yu-Wing Tai and In So Kweon (ICCV 2013)



Multiview Photometric Stereo using Planar Mesh Parameterization

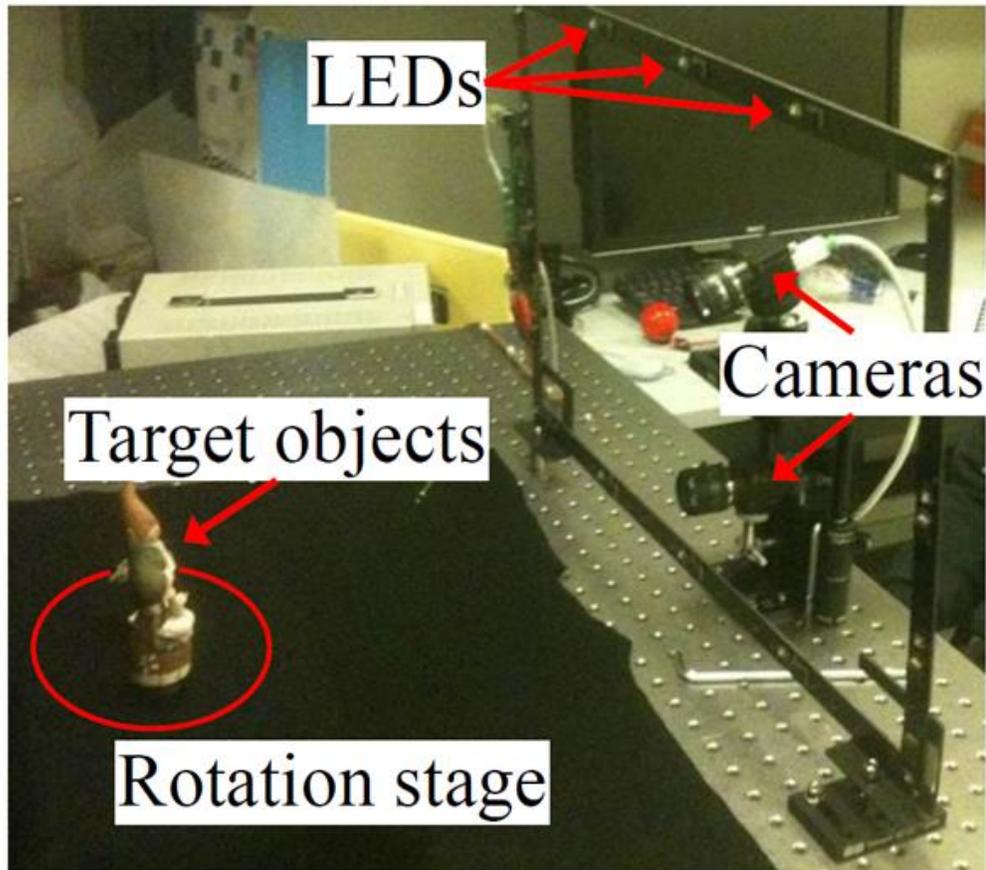
Jaesik Park, Sudipta N. Sinha, Yasuyuki Matsushita, Yu-Wing Tai and In So Kweon (ICCV 2013)



Multiview Photometric Stereo using Planar Mesh Parameterization

Jaesik Park, Sudipta N. Sinha, Yasuyuki Matsushita, Yu-Wing Tai and In So Kweon (ICCV 2013)

Acquisition Setup

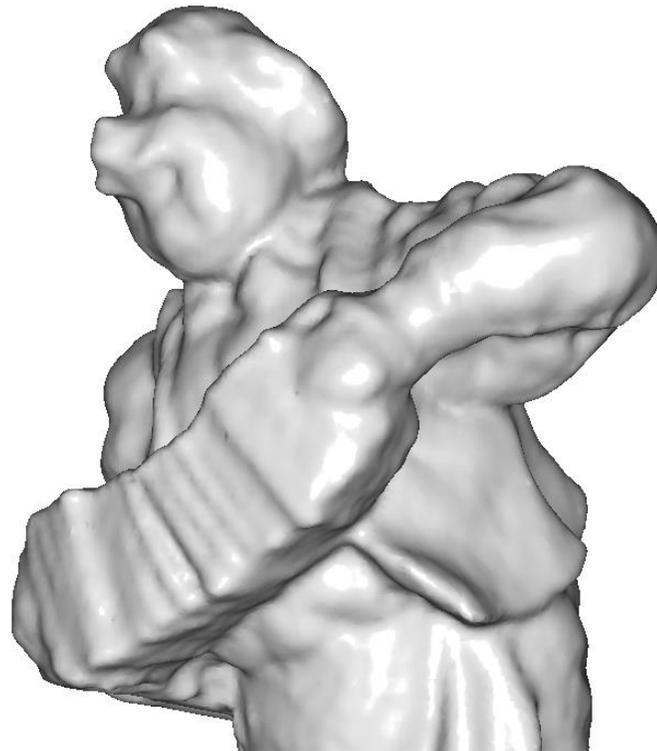


Buddha Statue

Multiview Photometric Stereo using Planar Mesh Parameterization



Input image



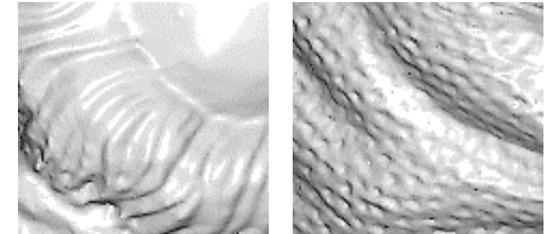
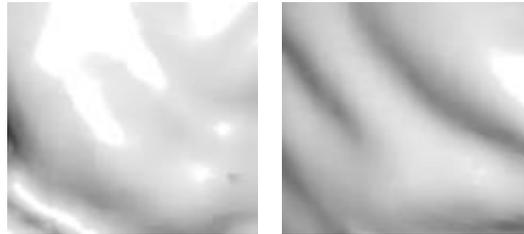
Base Mesh (Multiview Stereo)



Final 3D model

Fine geometric details recovered from hi-res images (10+ MP)

Multiview Photometric Stereo using Planar Mesh Parameterization



Input image



Base Mesh (Multiview Stereo)



Final 3D model

Conclusion

- New directions for hi-res stereo matching
- New datasets and benchmarks are coming !
- Local Plane Sweep stereo
 - avoid exploring the whole DSI
 - moving away from monolithic optimization
- Surface-based stereo matching
 - robust; allows a range of priors to be incorporated
- Combining multi-view stereo and photometric stereo

Acknowledgements

Daniel Scharstein, Rick Szeliski, Adarsh Kowdle,
Michael Bleyer, Carsten Rother, Pushmeet Kohli,
Jaesik Park, Yasuyuki Matsushita, Yu-Wing Tai, In So Kweon

Thanks!