

CVPR 2015 PerMeCop Workshop, June 11, 2015



The Middlebury Stereo Evaluation Version 3

Daniel Scharstein

Department of Computer Science Middlebury College



This work is supported by NSF grants IIS-0413169 and IIS-0917109

Joint work with ...



Heiko Hirschmüller, DLR / roboception



Duncan Levear '15



Nera Nesic '13



Greg Krathwohl '14



Xi Wang '14



York Kitajima '15



Porter Westling '12

...building on joint work with



Rick Szeliski, Microsoft Research

Alan Lim '09





Anna Blasiak '07

Sarri Al-Nashashibi '08





Jeff Wehrwein '08





Gonzalo Alonso '06



Jiaxin (Lily) Fu '03

Brad Hiebert-Treuer '07



Alexander Vandenberg-Rodes



Padma Ugbabe '03



Motivation

 StereoEval v.2 unchanged since 2005, virtually solved, overfitting

 KITTI (2012) was important addition, but only provides partial new challenges

Increasing need for more datasets

Goals

 Provide new hi-res datasets to propel research in stereo matching

 Address common problems with existing benchmarks

• Provide useful UI

Middlebury Stereo Page

(Scharstein & Szeliski – CVPR 2001, IJCV 2002)

- vision.middlebury.edu/stereo
- Evaluator with web interface

v.1 (2002) by Lily Fu '03 v.2 (2005) by Anna Blasiak '07



Middlebury stereo benchmark v.2

Currently 162 entries

vision.middlebury.edu														
										stereo	• mv	iew •	MRF •	flow • color
Stereo Evaluation • Datasets • Code • Submit Stereo Table Version 3 (beta) 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 Error Threshold ▼				Sort by all				Sort by disc						
Algorithm	Avg.	Tsukuba ground truth			Venus Teddy ground truth ground truth			Teddy ground truth	th ground truth			Average percent of bad pixels (<u>explanation</u>)		
	Rank	nonocc	all	<u>disc</u>	nonocc	all	<u>disc</u>	nonocc	all	<u>disc</u>	nonocc		<u>disc</u>	
LCU [156]	11.6	<u>1.06</u> 17	1.34 8	5.50 15	<u>0.07</u> 2	0.26 15	1.03 2	<u>3.68</u> 15	9.95 35	10.4 14	<u>1.63</u> 2	6.87 12	4.82 2	3.89
TSG0 [143]	11.7	<u>0.87</u> 4	1.13 1	4.66 8	<u>0.11</u> 8	0.24 9	1.47 11	<u>5.61</u> 41	8.09 18	13.8 34	<u>1.67</u> з	6.16 ₂	4.95 3	4.06
JSOSP+GCP [151]	13.5	<u>0.74</u> 1	1.34 9	3.98 1	<u>0.08</u> 3	0.16 1	1.15 з	<u>3.96</u> 17	10.1 38	11.8 20	<u>2.28</u> 18	7.91 32	6.74 21	4.18
ADCensus [82]	15.8	<u>1.07</u> 20	1.48 18	5.73 23	<u>0.09</u> 4	0.25 12	1.15 з	<u>4.10</u> 19	6.22 8	10.9 16	<u>2.42</u> 24	7.25 18	6.95 25	3.97
AdaptingBP [16]	19.7	<u>1.11</u> 23	1.37 11	5.79 25	<u>0.10</u> 6	0.21 8	1.44 10	<u>4.22</u> 21	7.06 16	11.8 21	<u>2.48</u> 28	7.92 34	7.32 33	4.23
CoopRegion [39]	20.0	<u>0.87</u> 6	1.16 2	4.61 5	<u>0.11</u> 7	0.21 6	1.54 15	<u>5.16</u> 33	8.31 22	13.0 28	<u>2.79</u> 48	7.18 17	8.01 53	4.41
CCRADAR [152]	24.3	<u>1.15</u> 26	1.42 18	6.23 39	<u>0.15</u> 20	0.27 16	1.89 25	<u>5.39</u> 38	10.6 41	14.7 45	<u>2.01</u> 4	7.37 20	5.88 4	4.75
<u>RDP [87]</u>	25.8	<u>0.97</u> 11	1.39 13	5.00 11	<u>0.21</u> 42	0.38 33	1.89 25	<u>4.84</u> 25	9.94 34	12.6 25	<u>2.53</u> 32	7.69 24	7.38 34	4.57
MultiRBF [129]	25.8	<u>1.33</u> 50	1.56 23	6.02 34	<u>0.13</u> 12	0.17 3	1.84 22	<u>5.09</u> 31	6.36 9	13.4 32	<u>2.90</u> 54	6.76 10	7.10 30	4.39
DoubleBP [34]	26.4	<u>0.88</u> 8	1.29 6	4.76 9	0.13 13	0.45 51	1.87 24	<u>3.53</u> 14	8.30 21	9.63 9	<u>2.90</u> 53	8.78 64	7.79 45	4.19

Limitations of existing stereo benchmarks

	Midd v.2	KITTI stereo
# of image pairs	4	194 + 195
Image size	< 0.2 MP	< 0.5 MP
Disparity range	16 60	70 150
Scene complexity	Limited	Limited
Scene variety	Limited	Limited
Realism	Limited	Good
Radiometric challenges	None	Some
Ground truth accuracy	Good	Limited
Ground truth coverage	Good	Limited
Control against overfitting	None	Some

Comparison with Midd v.3

	Midd v.2	KITTI stereo	Midd v.3
# of image pairs	4	194 + 195	15 + 15
Image size	< 0.2 MP	< 0.5 MP	5-6 MP
Disparity range	16 60	70 150	250 800
Scene complexity	Limited	Limited	Good
Scene variety	Limited	Limited	Good*
Realism	Limited	Good	Good*
Radiometric challenges	None	Some	Some
Ground truth accuracy	Good	Limited	Very good
Ground truth coverage	Good	Limited	Good
Control against overfitting	None	Some	Yes

* indoor only

New Datasets

- 2011-2013 collected 33 new datasets
 - multi exposure, multi lighting
 - floating-point disparities (PFM)

• 2013-2014 improved processing at DLR

 Build on structured lighting method by Scharstein & Szeliski [CVPR 2003]

New structured lighting system 2010-2014

- Portable rig
 - 2 DSLRs, 2 consumer cameras
- Improved Gray codes
- Natural scenes
- Specular surfaces



[Gupta et al., CVPR 2011]



2011: 5 datasets



2012: 7 datasets



2013: 21 datasets

• 2014: Improved processing























Improved processing

- How to get subpixel accuracy at 6 MP?
- Contributions:
 - Robust interpolation of codes
 - Fast 2D subpixel correspondence search
 - Improved calibration via bundle adjustment
 - Improved self-calibration of projectors
- Results:
 - Rectification and disparity accuracy of < 0.2 pixels
 - "perfect" and "imperfect" versions of datasets
- Best paper award at GCPR 2014

Processing pipeline old system



Rectification errors









Processing pipeline

new system



Bundle adjustment

- Minimize residual y-disps using nonlinear opt.
- Refine subset of camera params:



Accuracy

- No absolute measurements available, but
- Can check consistency of disparity estimates from P different projectors:
 - avg sample stddev s = 0.20
 - avg # of samples n = 7.7
 - we provide these as PFM images
- Can check residuals in planar regions

Averag	e absolute re	Improvement			
int disps	no subpix	ours			
\overline{r}_0	\overline{r}_1	\overline{r}_2	$\overline{r}_0/\overline{r}_2$	$\overline{r}_1/\overline{r}_2$	
0.252	0.135	0.032	7.9	4.2	





























-





Demo sv / plyv (part of svkit software)

Benchmark Design – Goals

- Integrate lessons learned (Middlebury, KITTI, Sintel, ...)
- Provide both overview and enable detailed interactive analysis
- Prevent overfitting; enable periodic recalibration of difficulty

Middlebury Benchmark v.3

- Datasets:
 - 15 training pairs, 15 test pairs (hidden GT)
 - Full, Half, and Quarter resolution
 - Some pairs with varying exposure, illumination
 - Most pairs with imperfect rectification
- Evaluation:
 - Multiple performance metrics
 - Weighted average for overall ranking
 - Allows evaluation of "sparse" results

Metrics

- Bad pixel % (t=0.5, 1.0, 2.0, 4.0)
- Avg abs and RMS error
- Quantiles 50, 90, 95, 99
- Runtimes, also normalized by # of pixels (t/MP) and # of disp hypotheses (t/Gdisp)
- Evaluate both dense and sparse results
- Region masks: nonocc or all regions

Middlebury Benchmark v.3

- Web interface by Duncan Levear '15
 - Automatic upload and evaluation
 - Interactive sorting and plotting
 - Selections persist between tables
 - Collect and display metadata
 - Supports history of snapshots
 - Supports periodic changes of weights

http://vision.middlebury.edu/stereo/eval3/



Overall Ranking

- Use weighted average instead of ranks
- Down-weight datasets too hard or too easy
- Allows adjusting of challenge as state of the art progresses
 - Adjust weights periodically
 - Could add more region masks in future
 - Could change default metric in future

Some personal insights

- "It's the UI, stupid" 🙂
- Huge value in compact representation and visualization of results
- Participation must be easy
- Strive for automatic scripts, but referee / moderator always needed
- Unless "one-shot contest," can never completely avoid overfitting
- Significant time commitment

Wishlist

- Accurate GT for outdoor images
- "Internet vision" dataset w/ dense GT
- Datasets for scene flow [Menze & Geiger CVPR 2015]
- More efforts by others (Sintel & KITTI good start)
- More work on synthetic images

Conclusion

- Benchmarks are important, stimulate research
- Creating ground-truth data is challenging, fun
- Hi-resolution images require new level of accuracy
- Stereo is not a solved problem!