

Depth From Stereo Camera

Enabling Machines to see as Human

Compiled by

Pankaj Kumar Bajpai

Samsung R&D Institute India, Bangalore

Dec 14, 2018

SAMSUNG

01 | Introduction

02 | Problem Formulation

03 | Computer Vision Approaches

04 | Deep Learning Approaches

05 | Challenges

06 | SRIB Achievements

SAMSUNG

An Introduction to

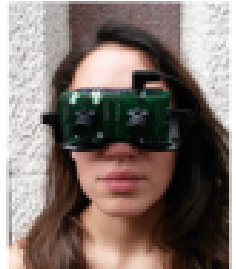
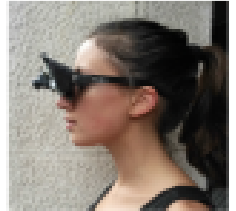
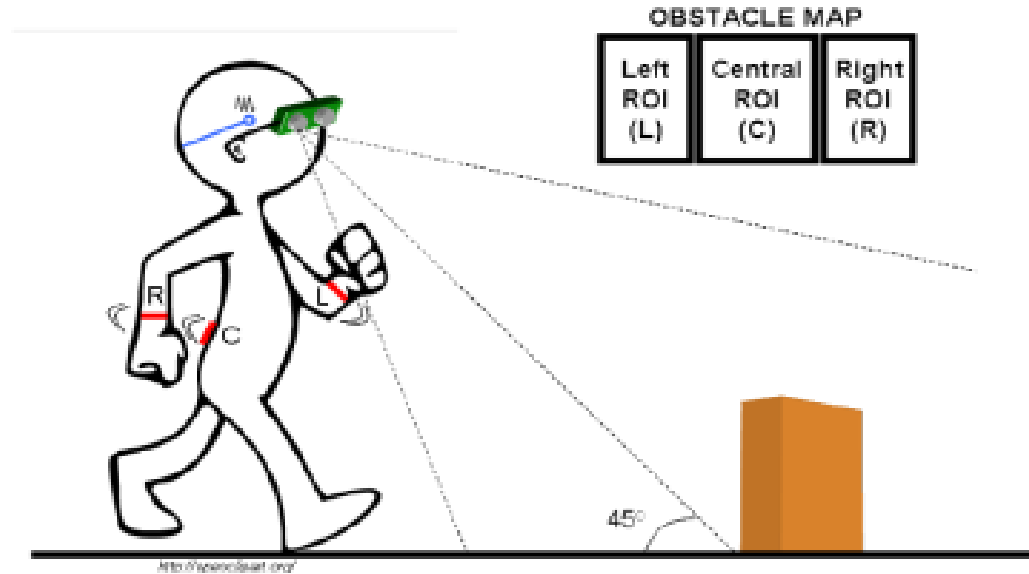
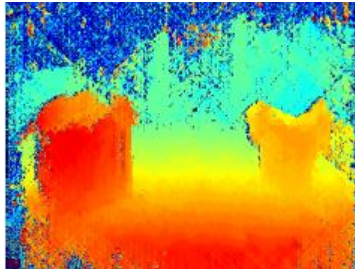
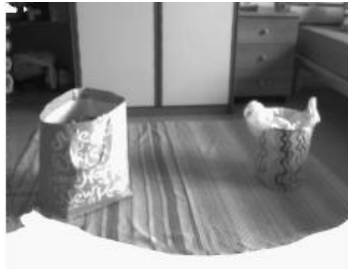
STEREO VISION

SAMSUNG

Why ? Importance of Depth



Autonomous Robot/
vehicle navigation

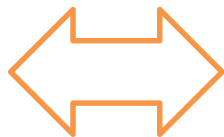


What ? Stereo Vision

➤ Close one of your eyes and try complicated tasks like tossing an object and catching it. Ask yourselves the following questions

- Can I perceive depth with one eye closed?
- If so, what cues does my eye use?
- Will it work well under all circumstances (like playing sports)?

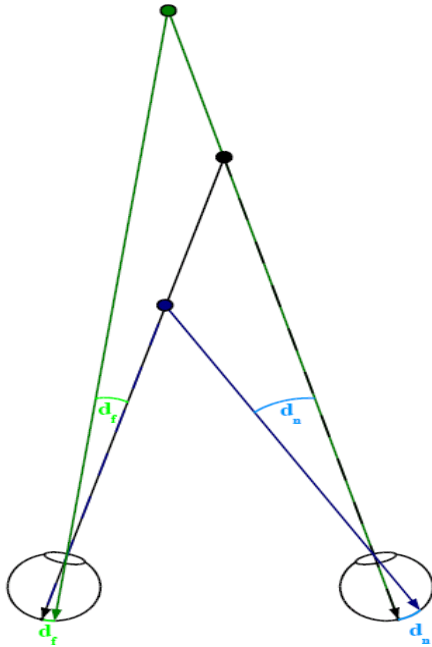
Stereo
Vision



Human
Eye



What ? Stereo Vision aka Binocular Vision



(a) Left eye image



(b) Right eye image

Two images of a stereoscopic photograph. The difference between the two images, such as the distances between the front cactus and the window in the two views, creates retinal disparity. This creates a perception of depth when (a) the left image is viewed by the left eye and (b) the right image is viewed by the right eye.

- Binocular **Disparity** : Relative 2D displacement of the image of the same point in space when projected on two different focal planes (i.e. two different eyes)
 - Objects closer to eye ➔ higher retinal disparity
 - Disparity inversely proportional to depth

Depth Perception : Other Modalities



- Focus
- Atmosphere

- Perspective
- Occlusion

- Motion based
- Past learning?

Scene analysis and 3D reconstruction

REAL WORLD SCANNING FOR AR AND VR

- Perfect virtual object integration (scale, occlusion, and lighting)
- Mixed reality experience by integrating real objects



<https://www.sony-deptsensing.com/DepthSense/Markets/HMD>

WORLD-FACING APPLICATIONS

- Mixed reality
- 3D object reconstruction
- 3D room reconstruction
- Indoor 3D navigation
- Metrology
- DSLR quality photography



<https://www.sony-deptsensing.com/DepthSense/Markets/HMD>

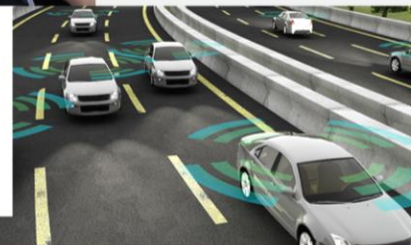
NEXT-GEN INFOTAINMENT CONTROL

- Hands-on wheel micro gestures
- HUD perspective correction (parallax) based on the position of the driver's head
- Augmented reality HUD



EXTERIOR SAFETY & COMFORT

- Detection of pedestrians, obstacles, nearby cars, bicycles and other hazards
- Automatic door/trunk release
- Autonomous parking
- Autonomous driving



<https://www.sony-deptsensing.com/DepthSense/Markets/Automotive>

Face modelling



USER-FACING APPLICATIONS

- Mixed reality
- Face authentication
- Touchless interaction
- DSLR quality photography

<https://www.sony-depthsensing.com/Depthsense/Markets/Mobile>

SECURITY MONITORING

- Biometric recognition
- Behavior analytics
- Intrusion detection

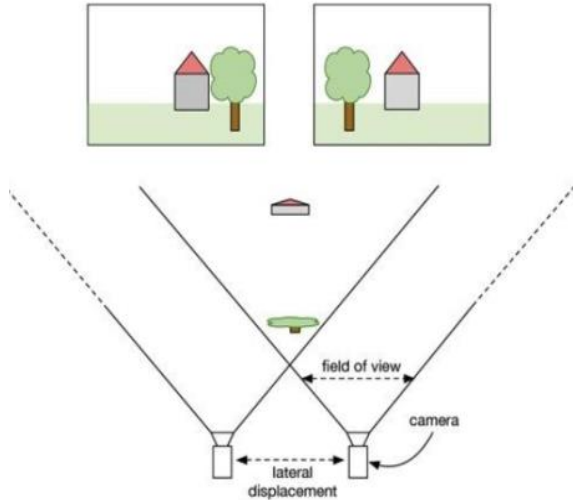


SAMSUNG

PROBLEM FORMULATION

SAMSUNG

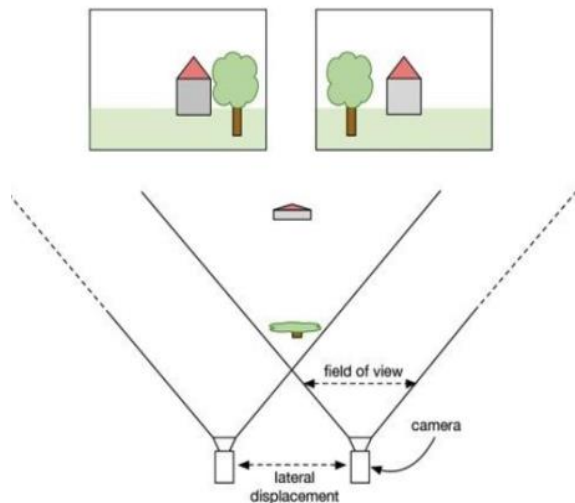
How ? Stereo Disparity & Depth



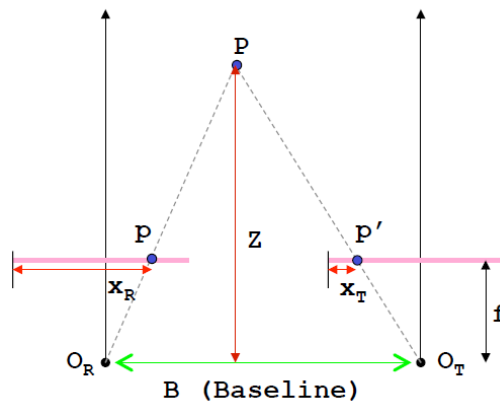
- Rectified pairs: 1D search space



How ? Stereo Disparity & Depth



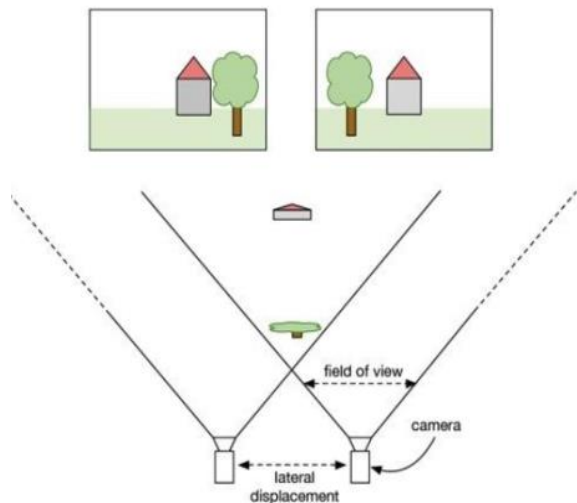
- Rectified pairs: 1D search space



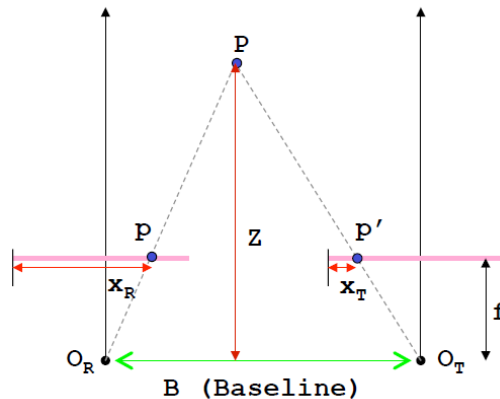
$$D(x, y) = x_R - x_T$$

$$Z(x, y) = \frac{B * f}{D(x, y)}$$

How ? Stereo Disparity & Depth



- Rectified pairs: 1D search space



$$D(x, y) = x_R - x_T$$

$$Z(x, y) = \frac{B * f}{D(x, y)}$$

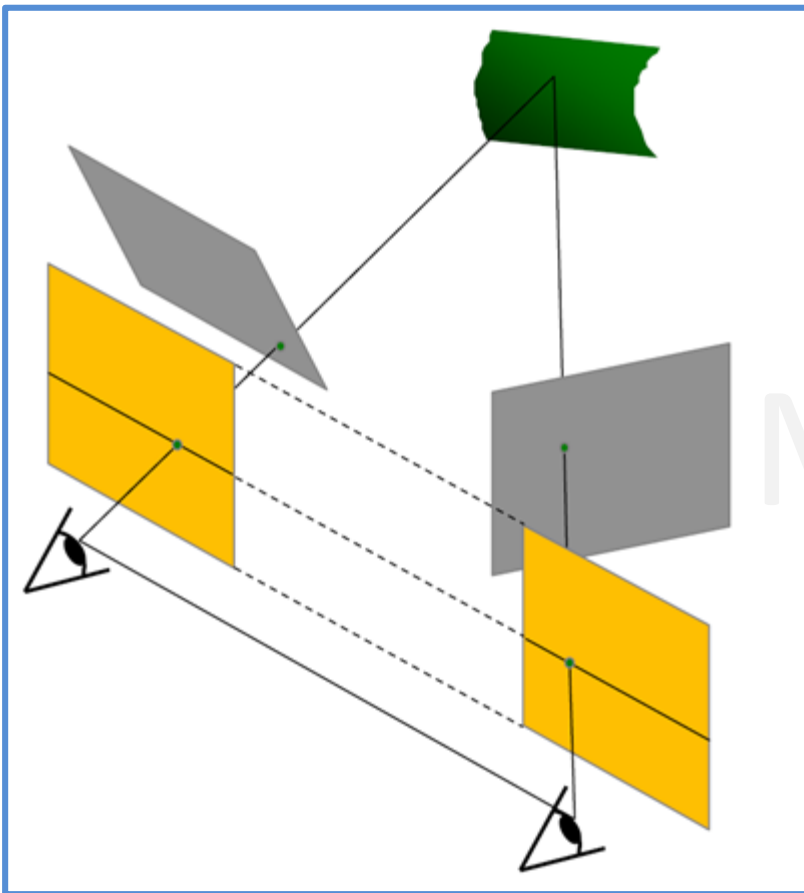
➤ Assumptions

- f, B are known (through camera calibration)
- The epipolar lines run horizontally
- The points p and p' are visible in both views

➤ Challenges

- For every point p in left image, how to find p' ?
➔ Stereo correspondence problem

Assumption : Rectified Stereo Image Pair



Input to Stereo Disparity Algorithm is considered to be Rectified Stereo Image Pair

- C. Loop and Z. Zhang. [Computing Rectifying Homographies for Stereo Vision](#). IEEE Conf. Computer Vision and Pattern Recognition, 1999.

SAMSUNG

COMPUTER VISION APPROACHES

SAMSUNG

Computer Vision (CV) based Stereo Disparity Pipeline

Cost computation

- Pixels dis-similarity
- Lower matching cost, more likely the match

Cost aggregation

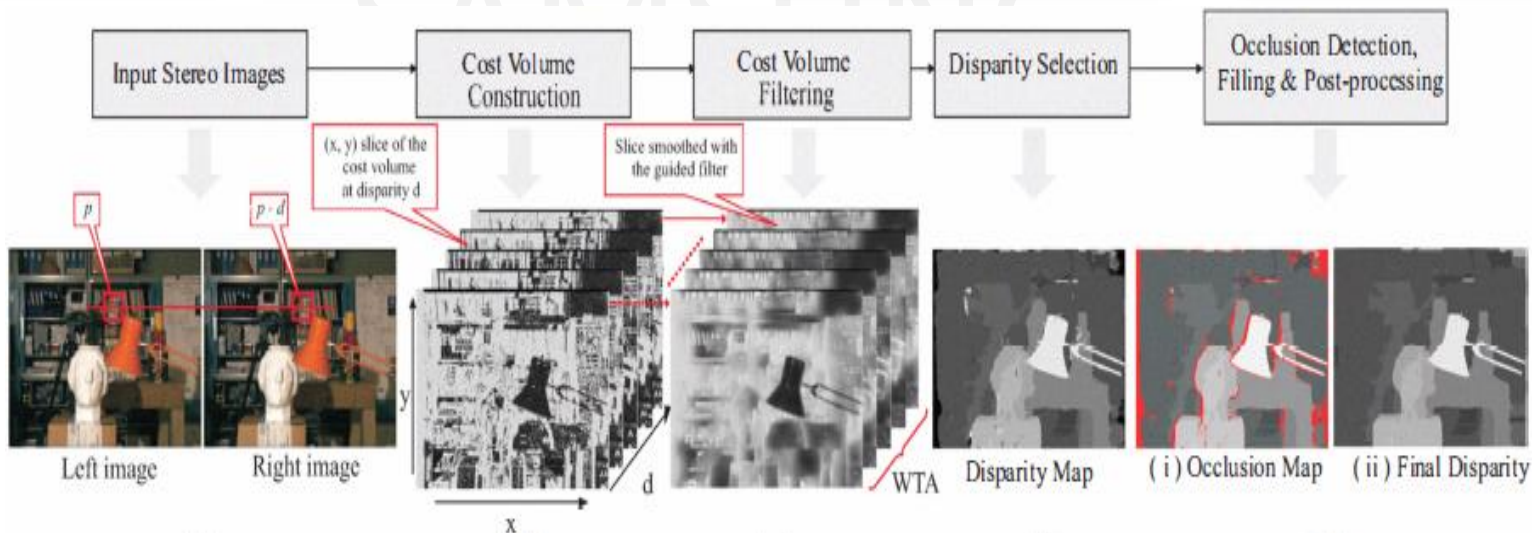
- Find suitable neighborhood
- Aggregation can be weighted

Optimization

- Strategy to decide disparity based on aggregated cost

Post Processing

- Handle occlusions
- Smooth and dense depth within objects



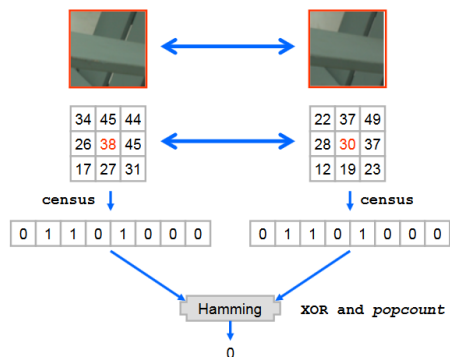
CV based Stereo Disparity – Simple Algorithm

Cost computation

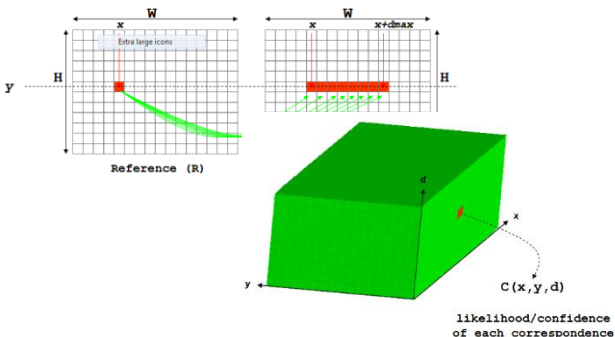
Cost aggregation

Optimization

Post Processing



SAMSUNG



CV based Stereo Disparity – Simple Algorithm

Cost computation

Cost aggregation

Optimization

Post Processing

- SAD
- Census transform
- Feature based :SIFT

SAMSUNG

CV based Stereo Disparity – Simple Algorithm

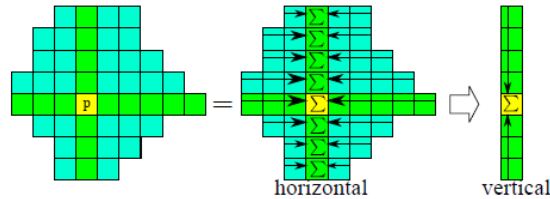
Cost computation

Cost aggregation

Optimization

Post Processing

- SAD
- Census transform
- Feature based :SIFT



CV based Stereo Disparity – Simple Algorithm

Cost computation

- SAD
- Census transform
- Feature based :SIFT

Cost aggregation

- Average over local region
- Cross-arm
- Guided/Bilateral filter

Optimization

Post Processing

CV based Stereo Disparity – Simple Algorithm

Cost computation

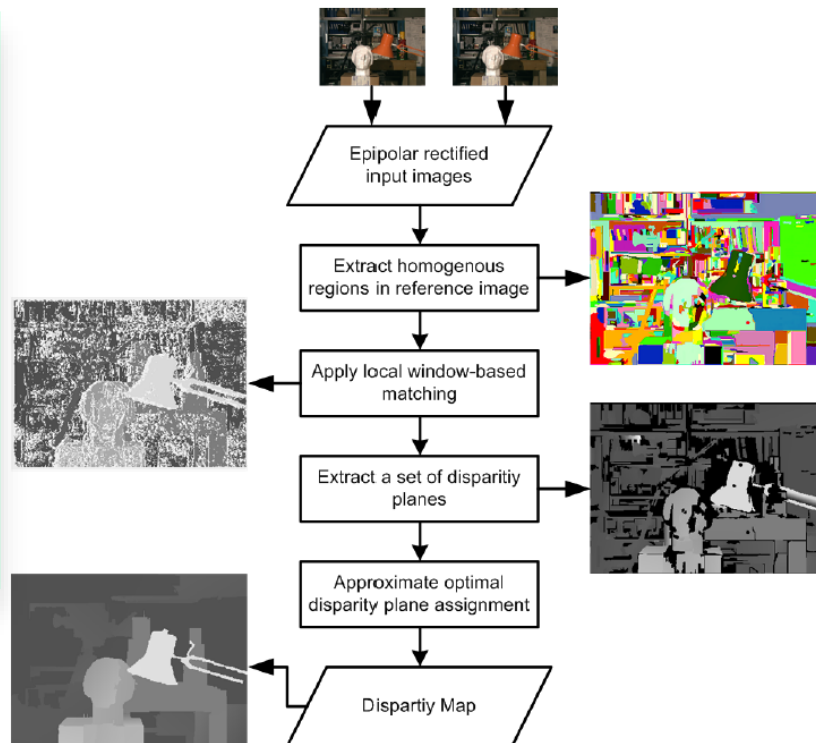
- SAD
- Census transform
- Feature based :SIFT

Cost aggregation

- Average over local region
- Cross-arm
- Guided/Bilateral filter

Optimization

Post Processing



CV based Stereo Disparity – Simple Algorithm

Cost computation

- SAD
- Census transform
- Feature based :SIFT

Cost aggregation

- Average over local region
- Cross-arm
- Guided/Bilateral filter

Optimization

- Semi-global matching
- Graph cut
- Belief propagation

Post Processing

CV based Stereo Disparity – Simple Algorithm

Cost computation

- SAD
- Census transform
- Feature based :SIFT

Cost aggregation

- Average over local region
- Cross-arm
- Guided/Bilateral filter

Optimization

- Semi-global matching
- Graph cut
- Belief propagation

Post Processing



CV based Stereo Disparity – Simple Algorithm

Cost computation

- SAD
- Census transform
- Feature based :SIFT

Cost aggregation

- Average over local region
- Cross-arm
- Guided/Bilateral filter

Optimization

- Semi-global matching
- Graph cut
- Belief propagation

Post Processing

- L-R consistency
- Subpixel refinement
- Segmentation techniques

- ❑ Feature Selection : CENSUS vs SAD vs SIFT vs ...
- ❑ Local or Small Neighbourhood Information
- ❑ ONLY Pixel Level Properties → NO SEMANTIC

Evolution → Towards Learning



CNN End-To-End : 2016+

End-To-End system removing separate pre & post processing

- "A Large Dataset to Train Convolutional Networks for Disparity, Optical Flow and Scene Flow Estimation", CVPR 2016

CNN Cost Function : ~2015

From hand crafted features to learned features, learning similarity between patches

- "A deep visual correspondence embedding model for stereo matching costs", ICCV 2015

Graph Based Methods : ~2011

Better correspondence searching, enhanced smoothness and occlusion handling

- "Kolmogorov and Zabih's graph cuts stereo matching algorithm", IPOL 2014
- "Pmbp: Patch match belief propagation for correspondence field estimation", IJCV 2014

Semi-global Matching : ~2008

Smoothness constraint

- "Stereo Processing by Semi-Global Matching and Mutual Information", TPAMI 2008

SAMSUNG

DEEP LEARNING APPROACHES

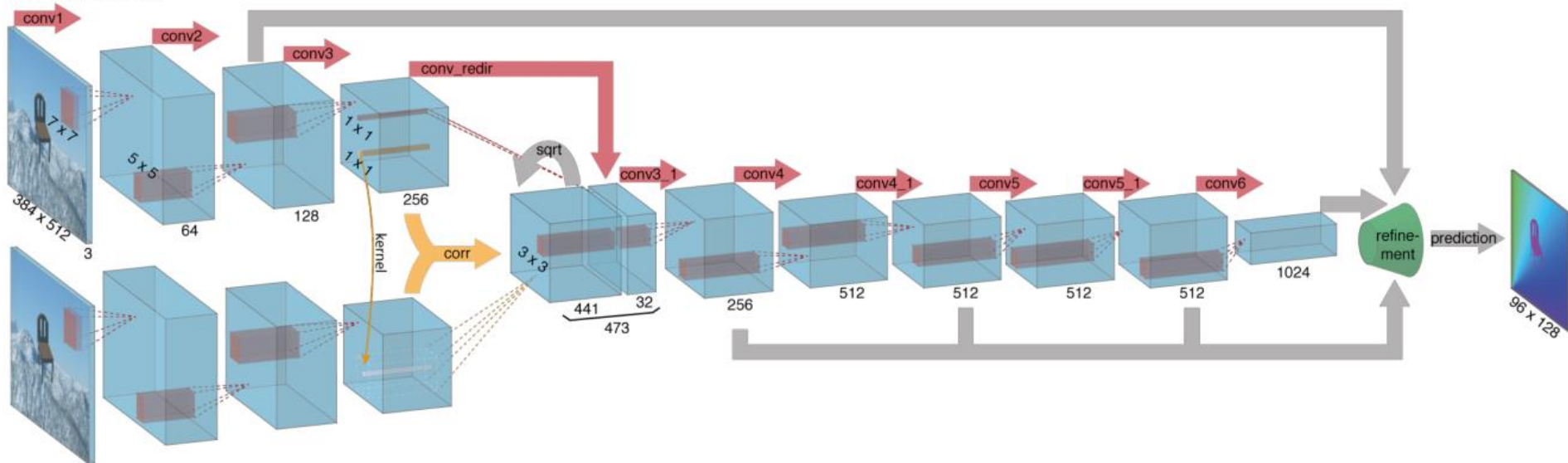
SAMSUNG

Deep Learning based Stereo Disparity

DispNet with Correlation Layer [DispNetC]

- First work: Mayer, N. et al. A Large Dataset to Train Convolutional Networks for Disparity, Optical Flow and Scene Flow Estimation. CVPR 2016

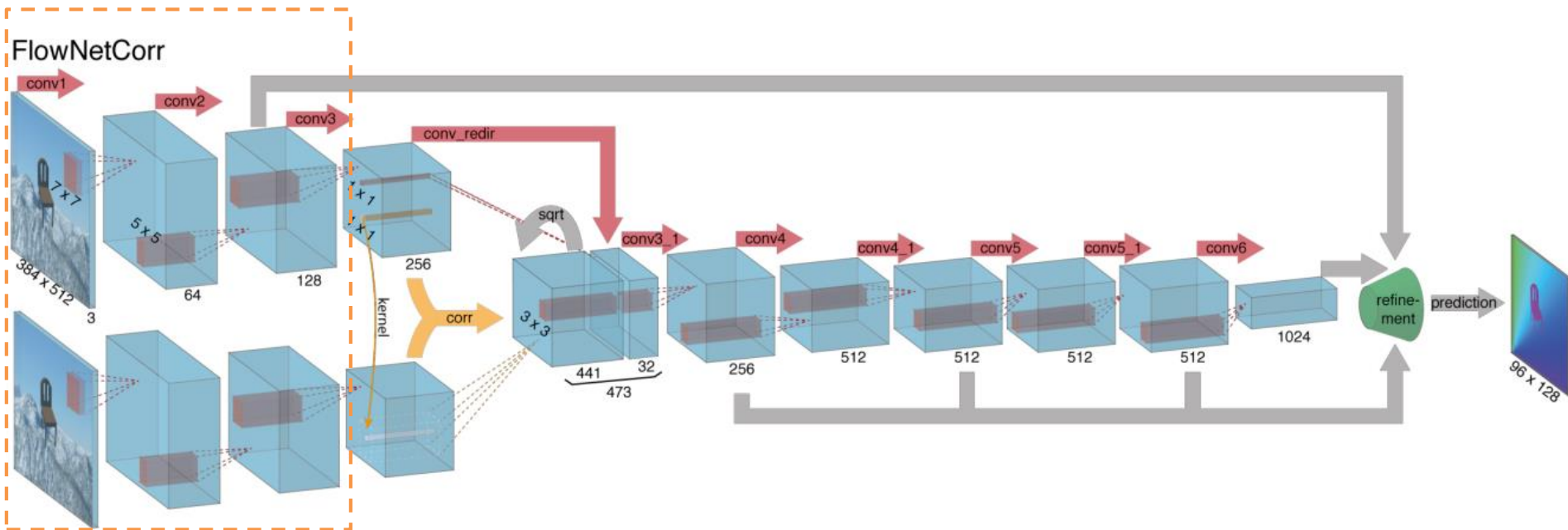
FlowNetCorr



Deep Learning based Stereo Disparity

DispNet with Correlation Layer [DispNetC]

- First work: Mayer, N. et al. A Large Dataset to Train Convolutional Networks for Disparity, Optical Flow and Scene Flow Estimation. CVPR 2016

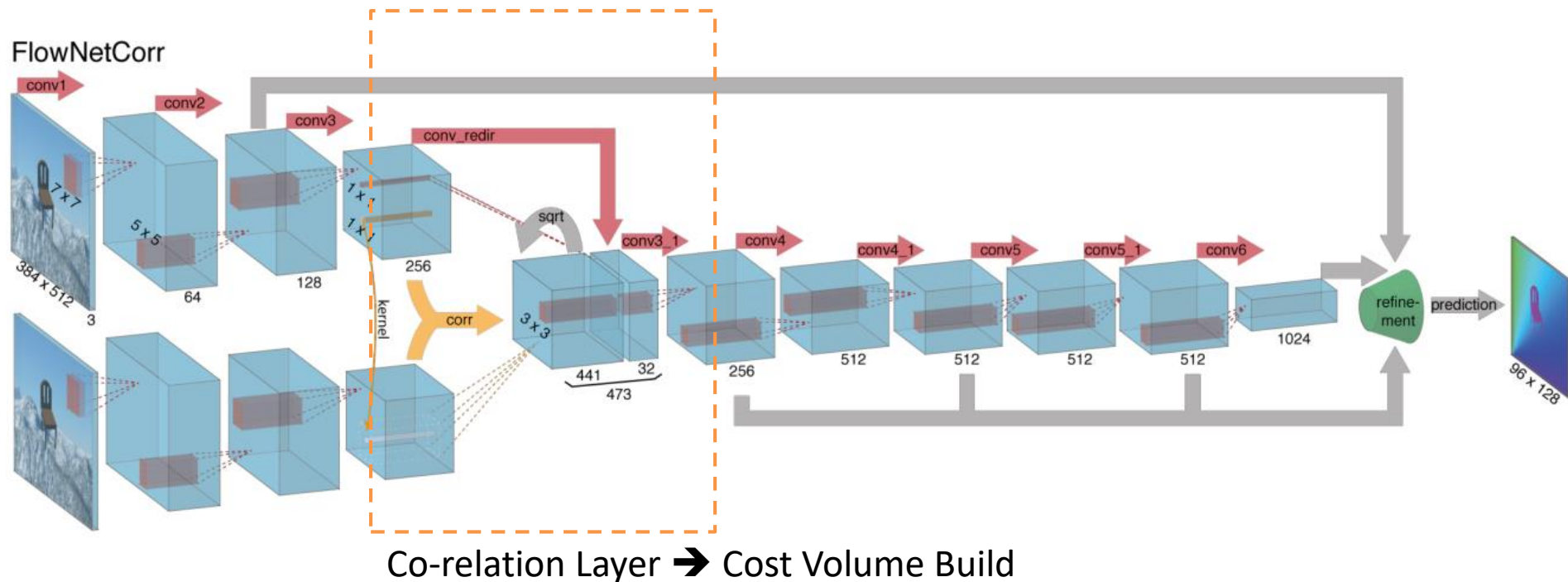


Feature Extraction → CENSUS

Deep Learning based Stereo Disparity

DispNet with Correlation Layer [DispNetC]

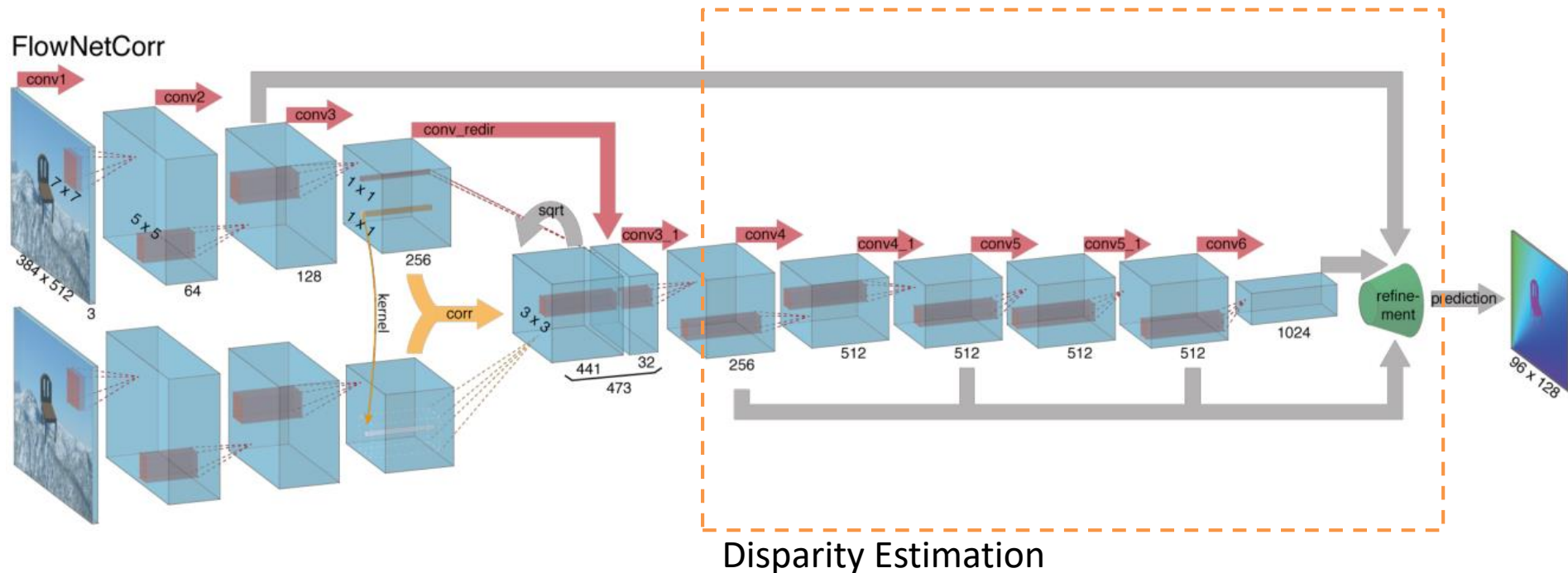
- First work: Mayer, N. et al. A Large Dataset to Train Convolutional Networks for Disparity, Optical Flow and Scene Flow Estimation. CVPR 2016



Deep Learning based Stereo Disparity

DispNet with Correlation Layer [DispNetC]

- First work: Mayer, N. et al. A Large Dataset to Train Convolutional Networks for Disparity, Optical Flow and Scene Flow Estimation. CVPR 2016



Deep Learning : DispNetC Results

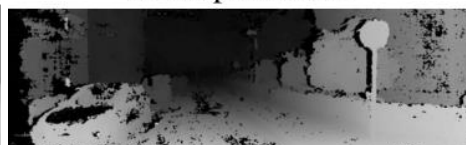
RGB image (L)



DispNetCorr1D-K



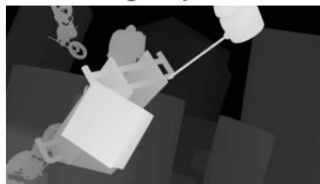
SGM prediction



RGB image (L)



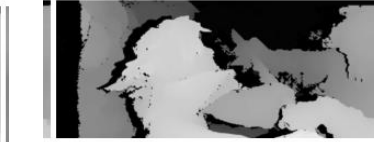
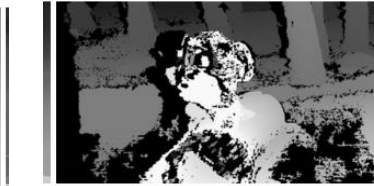
Disparity GT



DispNetCorr1D

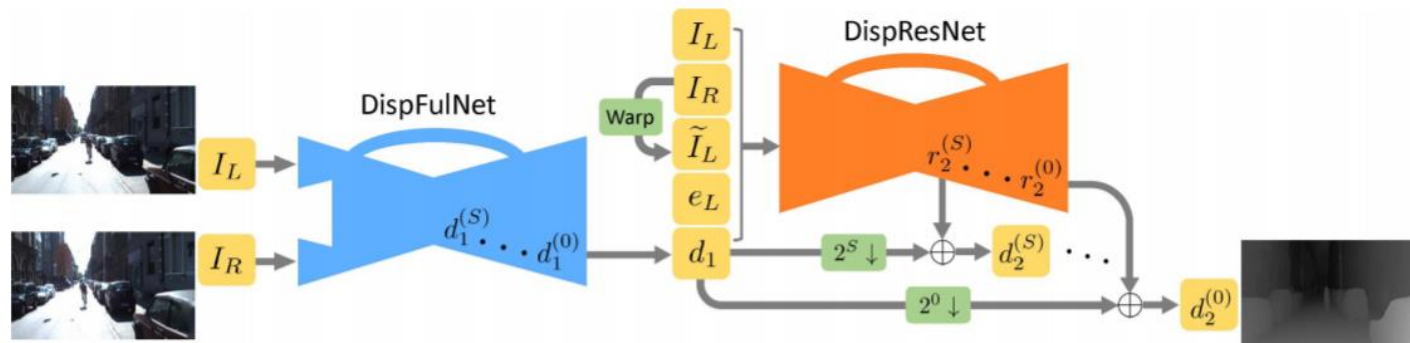


SGM prediction



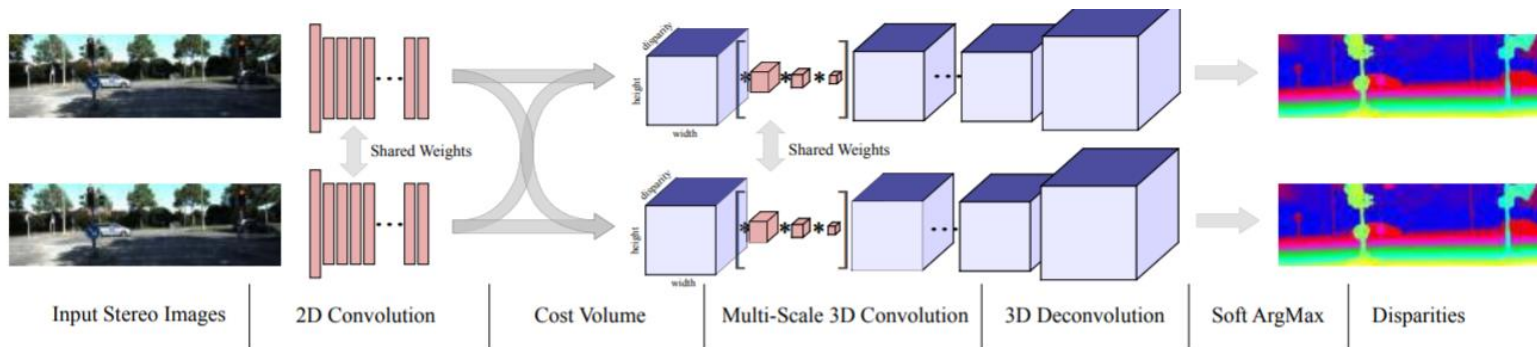
CRL : Cascaded Residual Learning (ICCV 2017)

<https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=8265317>



GC-Net : Geometry and Context Network (ICCV 2017)

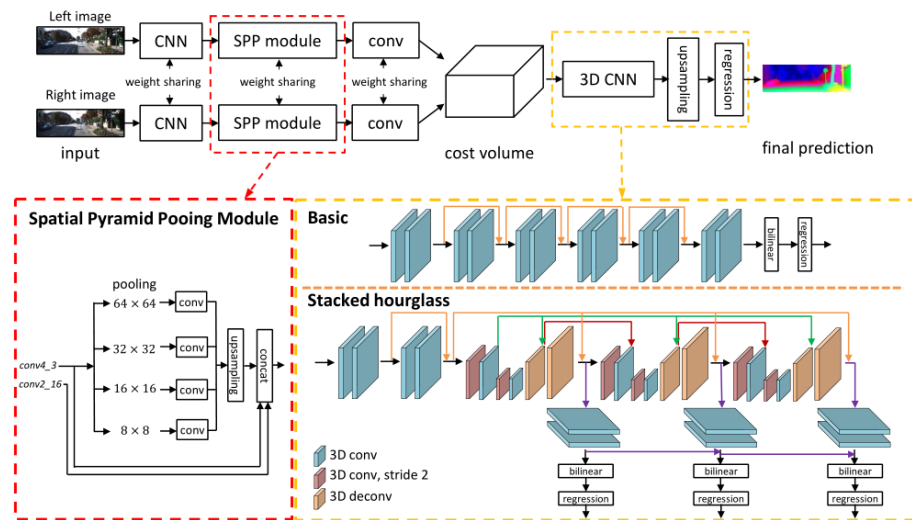
<https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=8237279>



Recent CNN Architectures

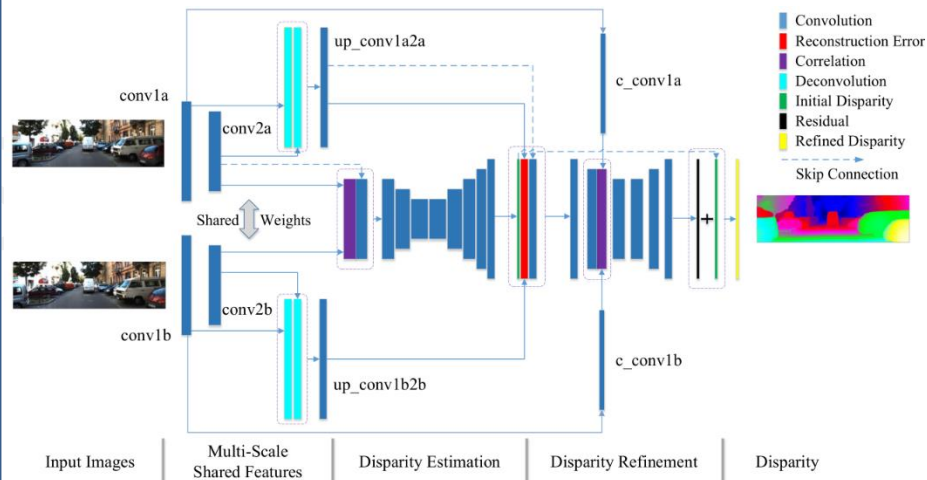
PSM-Net : Pyramid Stereo Matching Network (CVPR 2018)

<https://arxiv.org/pdf/1803.08669.pdf>



iResNet : Learning for Disparity Estimation through Feature Constancy (CVPR 2018)

<https://arxiv.org/pdf/1712.01039.pdf>



Open Challenges

Robust Vision Challenge (CVPR 2018) : “ foster the development of **vision systems** that are **robust** and consequently perform **well** on a **variety of datasets** with different characteristics”
<http://www.robustvision.net/>



Geometric
Corrections

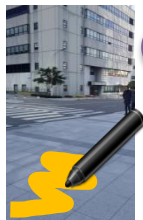
Occlusions



Target

Reference

Dense Data
Annotation



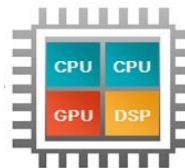
DNN

Specular ,
Transparent

Uniform
Regions



Real Time
Performance



Leader Board Position

Low-res two-view results - ETH3D

https://www.eth3d.net/low_res_two_view?mask=all&metric=bad-4-0

ETH3D Benchmark

Home Documentation Datasets Benchmarks About Submit

Coverage: dense Set: Test Metric: 99% quantile [px] Mask: all

Download this table as CSV

Method	Info	all	lakes. 1l	lakes. 1s	sand box 1l	sand box 1s	stora. room 1l	stora. room 1s	stora. room 2l	stora. room 2s	stora. room 2.1l	stora. room 2.1s	stora. room 2.2l	stora. room 2.2s	stora. room 3l	stora. room 3s	tunnel 1l	tunnel 1s	tunnel 2l
DISCO	00	2.40	0.99	11.99	1.81	1.12	3.47	1.24	3.34	2.57	2.73	1.96	4.52	1.90	4.18	3.62	0.37	0.27	0.31
DLCB_ROB	00	2.55	1.32	10.26	1.31	1.04	2.94	1.29	3.81	2.73	2.93	2.17	3.37	1.53	10.40	3.26	0.42	0.35	0.41
iResNet_ROB	00	2.72	10.38	14.41	1.46	0.99	3.74	0.66	4.01	2.84	2.04	2.05	1.98	1.38	4.73	1.91	0.35	0.20	0.28
LALA_ROB	00	2.88	2.13	19.08	1.10	0.99	3.91	1.81	4.50	3.04	3.32	1.87	2.75	1.66	4.06	4.02	0.47	0.44	0.55
DN-CSS_ROB	00	2.89	9.72	19.71	2.07	0.90	4.01	0.61	4.06	3.35	1.66	1.89	1.39	1.06	3.24	1.88	0.44	0.19	0.44
ETE_ROB	00	2.93	2.12	17.80	1.13	0.94	3.71	1.53	4.29	2.98	3.32	2.13	2.34	1.51	8.68	2.82	0.44	0.40	0.51
NCCL2	00	2.96	4.07	15.75	1.24	1.03	3.48	2.55	3.79	2.39	2.78	1.80	2.41	1.67	10.27	2.48	0.46	0.41	0.57

vision.middlebury.edu/stereo/eval3

Not secure | vision.middlebury.edu/stereo/eval3/

Stereo Evaluation • Datasets • Code • Submit

Middlebury Stereo Evaluation - Version 3

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 0.5 bad 1.0 bad 2.0 bad 4.0 avgerr rms A50 A90 A95 A99 time time/MP time/GD

Mask: nonocc all

☐ plot selected ☐ show invalid [Reset sort](#) [Reference list](#)

Date	A99 (pixels)	Name	Res	Avg	Weight	Austr	AustrP	Bicyc2	Class	ClassE	Compu	Crusa	CrusaP	DjemB	DjemBL	Hoops	Livgrm	Nkuba	Plants	Stairs
						MP: 5.5 nd: 290 im0: im1 GT	MP: 5.6 nd: 290 im0: im1 GT	MP: 5.5 nd: 250 im0: im1 GT	MP: 5.7 nd: 610 im0: im1 GT	MP: 5.7 nd: 610 im0: im1 GT	MP: 5.5 nd: 256 im0: im1 GT	MP: 5.5 nd: 800 im0: im1 GT	MP: 5.7 nd: 800 im0: im1 GT	MP: 5.7 nd: 320 im0: im1 GT	MP: 5.7 nd: 320 im0: im1 GT	MP: 5.7 nd: 410 im0: im1 GT	MP: 5.9 nd: 320 im0: im1 GT	MP: 5.5 nd: 570 im0: im1 GT	MP: 5.8 nd: 320 im0: im1 GT	MP: 5.5 nd: 450 im0: im1 GT
						nonocc	nonocc	nonocc	nonocc	nonocc	nonocc	nonocc	nonocc	nonocc	nonocc	nonocc	nonocc	nonocc	nonocc	nonocc
						↓↑	↓↑	↓↑	↓↑	↓↑	↓↑	↓↑	↓↑	↓↑	↓↑	↓↑	↓↑	↓↑	↓↑	↓↑
05/22/18		DN-CSS_ROB	H	82.0		118 2	107 1	40.5 1	71.4 4	146 9	29.0 1	99.9 2	105 3	21.6 2	66.4 18	171 30	57.3 8	65.7 2	113 1	127 7
10/10/18		DISCO	H	86.6		108 1	110 2	46.1 2	102 7	146 8	31.0 2	108 3	101 2	17.1 1	104 30	166 18	51.5 1	68.6 3	142 3	87.2 3
05/31/18		iResNet_ROB	H	87.5		121 4	118 3	55.7 3	86.6 5	114 2	95.0 20	70.9 1	73.7 1	32.1 10	47.5 10	172 31	59.2 12	62.1 1	142 2	145 8
05/26/18		NOSS_ROB	H	104		130 21	125 8	71.2 8	63.5 1	137 5	97.0 25	126 9	123 8	41.7 24	41.9 3	152 9	58.1 10	106 13	177 29	174 28
03/06/18		NOSS	H	104		130 21	125 8	71.2 8	63.5 1	137 5	97.0 25	126 9	123 8	41.7 24	41.9 3	152 9	58.1 10	106 13	177 29	174 27
05/01/18		PSMNet_ROB	Q	106		145 83	143 54	84.1 27	106 8	128 3	58.0 5	112 4	111 4	34.3 12	139 43	169 24	62.5 13	98.4 5	163 9	118 8
12/11/17		OVOD	H	108		120 3	118 4	69.4 4	122 10	183 10	99.0 35	129 15	130 14	34.8 14	55.3 14	149 3	64.5 19	80.5 4	168 14	153 9
06/22/17		LocalExp	H	100		128 16	126 12	70.9 7	89.5 6	175 13	97.0 26	124 10	125 5	27.6 10	43.9 6	158 16	63.0 18	108 16	175 21	163 16

** As on Oct-2018

More references on Stereo :CV method

- H. Hirschmüller. Stereo processing by semi-global matching and mutual information. PAMI 30(2):328-341, 2008
- S. Drouyer, et al. Sparse stereo disparity map densification using hierarchical image segmentation. 13th International Symposium on Mathematical Morphology.
- L. Li, X. Yu, S. Zhang, X. Zhao, and L. Zhang. 3D cost aggregation with multiple minimum spanning trees for stereo matching. Applied Optics 56(12):3411-3420, 2017.
- L. Li, S. Zhang, X. Yu, and L. Zhang. PMSC: PatchMatch-based superpixel cut for accurate stereo matching. IEEE Trans on Circuits and Systems for Video Technology, 2016.
- “Multiview Geometry in Computer Vision”, book by Hartley and Zisserman

More references on Stereo : End to end CNN method

- J. Chang and Y. Chen: Pyramid Stereo Matching Network. arXiv preprint arXiv:1803.08669 2018.
- Z. Liang, Y. Feng, Y. Guo and H. Liu: [Learning for Disparity Estimation through Feature Constancy](#). arXiv preprint arXiv:1712.01039 2017.
- J. Pang, et al: [Cascade residual learning: A two-stage convolutional neural network for stereo matching](#). ICCV Workshop on Geometry Meets Deep Learning 2017.

Tutorial Material on Stereo

- <http://www.cse.psu.edu/~rtc12/CSE486/lecture09.pdf>
- <http://www.inf.u-szeged.hu/~kato/teaching/computervision/02-CameraGeometry.pdf>
- <http://www.ics.uci.edu/~majumder/vispercep/chap8notes.pdf>
- <http://vision.deis.unibo.it/~smatt/Seminars/StereoVision.pdf>
- <https://courses.cs.washington.edu/courses/cse455/09wi/Lects/lect16.pdf>

Stereo Datasets:

- Middlebury: <http://vision.middlebury.edu/stereo/eval3/>
- Kitti: http://www.cvlibs.net/datasets/kitti/eval_scene_flow.php?benchmark=stereo
- SceneFlow: <https://lmb.informatik.uni-freiburg.de/resources/datasets/SceneFlowDatasets.en.html>
- ETH dataset: https://www.eth3d.net/low_res_two_view

- Extend sincere gratitude to following members for providing their valuable support and help
 - ☐ Kunal Swami
 - ☐ Dr. Rituparna Sarkar
 - ☐ Yash Harbhajanka
 - ☐ Bhushan Bhagwan Gawde
 - ☐ Dr. Lokesh Boregowda
- Images are borrowed from various sources and internet
 - ☐ *"Stereo Vision: Algorithms and Applications" by Stefano Mattocia*
 - ☐ *"On Building an Accurate Stereo Matching System on Graphics Hardware" by Xing Mei*
 - ☐ *Etc...*



SAMSUNG