

# Depth Estimation Using Deep Neural Networks

Mușat Bogdan-Adrian

December 2017

# Problem Formulation

- Given a pair of stereo images, how to estimate the depth of the scene?

# Disparity

- Disparity represents the difference of perspective created by the horizontal or vertical separation of two cameras
- The human brain processes disparity information from both eyes to estimate the depth of real world objects
- Disparity is inverse proportional with depth; given the baseline distance  $b$  between the cameras and the camera focal length  $f$ , the depth  $\hat{d}$  from the predicted disparity  $d$  is simply  $\hat{d} = bf/d$
- Objects closer to the camera have bigger disparities, whilst objects that are farther have smaller disparities

# Proposed Solution

- Given a pair of rectified<sup>1</sup> stereo images as input, build a model which can accurately predict the disparity per-pixel (i.e. learn the  $\Delta$  values with which every pixel from the left image is shifted from the right one on the x axis)
- Such a system can be modeled by a Convolutional Neural Network
- The dataset used for experiments is FlyingThings3D<sup>2</sup>

---

<sup>1</sup>Vertically aligned

<sup>2</sup>N.Mayer et al. "A Large Dataset to Train Convolutional Networks for Disparity, Optical Flow, and Scene Flow Estimation". In: *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*. arXiv:1512.02134. 2016. URL: <http://lmb.informatik.uni-freiburg.de/Publications/2016/MIFDB16>.

## Challenges - Textureless Areas

- Depth estimation is an ill-posed problem: a pixel in a textureless area from the left image can belong to multiple pixels in the right image

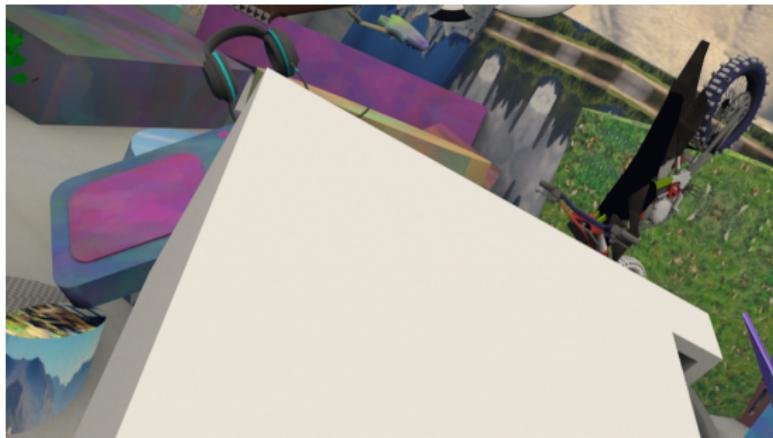


Figure 1: Example of objects with textureless areas

## Challenges - Object Occlusion

- Parts of an object may be visible in an image, while being absent in the other one because of the difference of perspective



Figure 2: Occlusion in stereo images

# Convolutional Neural Networks Architectures

- Supervised
  - Requires ground-truth disparity, which might be expensive to obtain
  - Based on disparity regression
  - Accurate predictions
- Unsupervised
  - Does not require any form of ground-truth
  - Based on an image reconstruction loss
  - Has artefacts around object boundaries, caused by occlusion

# DispNet<sup>3</sup> - Overview

- Image-to-image with contractive and expanding structure

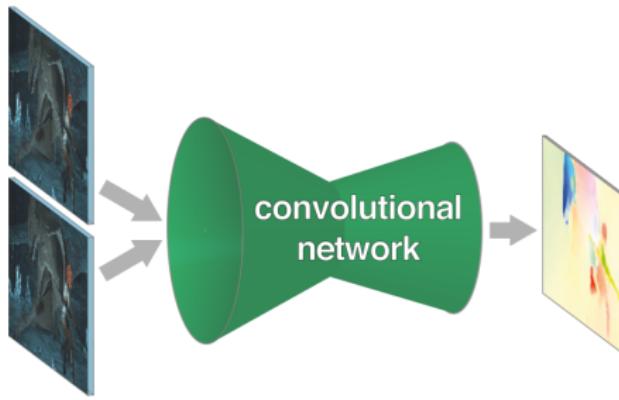


Figure 3: Hourglass structure of DispNet

---

<sup>3</sup>Nikolaus Mayer et al. “A Large Dataset to Train Convolutional Networks for Disparity, Optical Flow, and Scene Flow Estimation”. In: *CoRR* abs/1512.02134 (2015). arXiv: 1512.02134. URL: <http://arxiv.org/abs/1512.02134>.

# DispNet - Architectural Details

- Fully Convolutional Network<sup>4</sup>
- Concatenates the “upconvolution” results with the features from the “contractive” part to recover spatial information lost by downsampling
- Uses smooth L1 loss between the resulted disparity computed at different pyramid levels of the network and the ground-truth

$$|x|_{smooth} = \begin{cases} 0.5x^2, & \text{if } |x| \leq 1 \\ |x| - 0.5, & \text{otherwise} \end{cases}$$

---

<sup>4</sup> Jonathan Long, Evan Shelhamer, and Trevor Darrell. “Fully Convolutional Networks for Semantic Segmentation”. In: *CoRR* abs/1411.4038 (2014). arXiv: 1411.4038. URL: <http://arxiv.org/abs/1411.4038>.

# DispNet - Architectural Details

- The upsampling part of the network can be constructed by using transposed convolutions<sup>5</sup> or by using interpolation methods like nearest neighbor or bilinear followed by a convolutional layer to smoothen the results

---

<sup>5</sup> Matthew D. Zeiler, Graham W. Taylor, and Rob Fergus. "Adaptive Deconvolutional Networks for Mid and High Level Feature Learning". In: *Proceedings of the 2011 International Conference on Computer Vision*. ICCV '11. Washington, DC, USA: IEEE Computer Society, 2011, pp. 2018–2025. ISBN: 978-1-4577-1101-5. DOI: 10.1109/ICCV.2011.6126474. URL: <http://dx.doi.org/10.1109/ICCV.2011.6126474>.

# DispNet - Results



Figure 4: Supervised depth estimation

# DispNet - Limitations

- Requires ground-truth data which might be expensive to obtain
- Current depth estimation hardware for collecting ground-truth data like LIDAR, Time of flight, Kinect may provide inaccurate predictions caused by environmental factors, which directly affects the network's performance

# Unsupervised Monocular Network<sup>6</sup> - Overview

- Similar architecture as DispNet

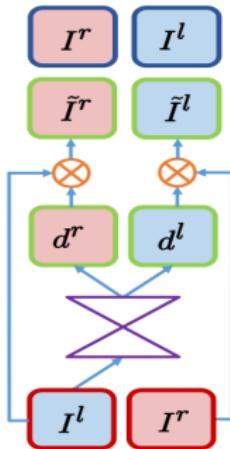


Figure 5: Unsupervised Monocular Network architecture

<sup>6</sup>Clément Godard, Oisin Mac Aodha, and Gabriel J. Brostow. “Unsupervised Monocular Depth Estimation with Left-Right Consistency”. In: *CoRR abs/1609.03677* (2016). arXiv: 1609.03677. URL: <http://arxiv.org/abs/1609.03677>.

# Unsupervised Monocular Network - Architectural Details

- Shares the same hourglass structure as DispNet and works the same, up until a point
- Instead of using ground-truth disparity, the network tries to reconstruct the right image from the left one, and vice-versa
- The network tries to achieve the needed disparities as to warp the left and right image using some sampling method
- The sampling method chosen is binilinear sampling as used in Spatial Transformer Networks<sup>7</sup>, which is fully differentiable
- During inference, only the left view image is used (monocular estimation)

---

<sup>7</sup>Max Jaderberg et al. “Spatial Transformer Networks”. In: *CoRR* abs/1506.02025 (2015). arXiv: 1506.02025. URL: <http://arxiv.org/abs/1506.02025>.

# Unsupervised Monocular Network - Architectural Details

- Uses a three term loss

$$C_s = \alpha_{ap}(C_{ap}^l + C_{ap}^r) + \alpha_{ds}(C_{ds}^l + C_{ds}^r) + \alpha_{lr}(C_{lr}^l + C_{lr}^r)$$

- Appearance Matching Loss

$$C_{ap}^l = \frac{1}{N} \sum_{i,j} \alpha \frac{1 - SSIM(I_{ij}^l, \tilde{I}_{ij}^l)}{2} + (1 - \alpha) \|I_{ij}^l - \tilde{I}_{ij}^l\|$$

- Disparity Smoothness Loss

$$C_{ds}^l = \frac{1}{N} \sum_{i,j} |\partial_x d_{ij}^l| e^{-|\partial_x I_{ij}^l|} + |\partial_y d_{ij}^l| e^{-|\partial_y I_{ij}^l|}$$

- Left-Right Disparity Consistency Loss

$$C_{lr}^l = \frac{1}{N} \sum_{i,j} |d_{ij}^l - d_{ij+d_{ij}^l}^r|$$

# Unsupervised Monocular Network - Results

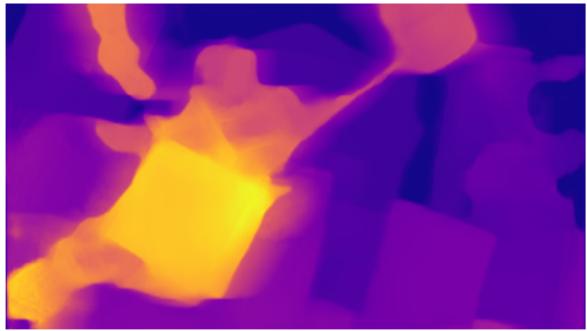


Figure 6: Unsupervised depth estimation

# Unsupervised Monocular Network - Limitations

- There are visible artefacts around the object boundaries due to the pixels in the occlusion region not being visible in both images
- Because this method relies solely on image reconstruction, textureless or transparent surfaces will produce inconsistent depths
- Complex objective loss

# Benchmark Evaluation

- One of the main benchmarks for depth estimation is KITTI 2015<sup>8</sup> and the evaluation metric used is called D1-all, which is the percentage of pixels for which the estimation error is larger than 3px and larger than 5% of the ground truth disparity at this pixel
- On this benchmark, DispNet achieves a D1-all error of 4.34%, while the unsupervised network gets a 23.81% error

---

<sup>8</sup>Moritz Menze and Andreas Geiger. "Object Scene Flow for Autonomous Vehicles". In: *Conference on Computer Vision and Pattern Recognition (CVPR)*. 2015.

# Public Implementations

- DispNet
  - <https://lmb.informatik.uni-freiburg.de/resources/software.php>
- Unsupervised Monocular Network
  - <https://github.com/mrharicot/monodepth>

# Conclusions

- Unsupervised methods work more alike the brain in the sense that depth is learnt by image reconstruction enforcement
- There should be a preference towards unsupervised models due to their lack of needing precise estimations
- For the moment, supervised networks perform much better than their unsupervised counterparts on popular benchmarks

# Technical Slides - DispNet Architecture

Name	Kernel	Str.	Ch I/O	InpRes	OutRes	Input
conv1	$7 \times 7$	2	6/64	$768 \times 384$	$384 \times 192$	Images
conv2	$5 \times 5$	2	64/128	$384 \times 192$	$192 \times 96$	conv1
conv3a	$5 \times 5$	2	128/256	$192 \times 96$	$96 \times 48$	conv2
conv3b	$3 \times 3$	1	256/256	$96 \times 48$	$96 \times 48$	conv3a
conv4a	$3 \times 3$	2	256/512	$96 \times 48$	$48 \times 24$	conv3b
conv4b	$3 \times 3$	1	512/512	$48 \times 24$	$48 \times 24$	conv4a
conv5a	$3 \times 3$	2	512/512	$48 \times 24$	$24 \times 12$	conv4b
conv5b	$3 \times 3$	1	512/512	$24 \times 12$	$24 \times 12$	conv5a
conv6a	$3 \times 3$	2	512/1024	$24 \times 12$	$12 \times 6$	conv5b
conv6b	$3 \times 3$	1	1024/1024	$12 \times 6$	$12 \times 6$	conv6a
pr6+loss6	$3 \times 3$	1	1024/1	$12 \times 6$	$12 \times 6$	conv6b
upconv5	$4 \times 4$	2	1024/512	$12 \times 6$	$24 \times 12$	conv6b
iconv5	$3 \times 3$	1	1025/512	$24 \times 12$	$24 \times 12$	upconv5+pr6+conv5b
pr5+loss5	$3 \times 3$	1	512/1	$24 \times 12$	$24 \times 12$	iconv5
upconv4	$4 \times 4$	2	512/256	$24 \times 12$	$48 \times 24$	iconv5
iconv4	$3 \times 3$	1	769/256	$48 \times 24$	$48 \times 24$	upconv4+pr5+conv4b
pr4+loss4	$3 \times 3$	1	256/1	$48 \times 24$	$48 \times 24$	iconv4
upconv3	$4 \times 4$	2	256/128	$48 \times 24$	$96 \times 48$	iconv4
iconv3	$3 \times 3$	1	385/128	$96 \times 48$	$96 \times 48$	upconv3+pr4+conv3b
pr3+loss3	$3 \times 3$	1	128/1	$96 \times 48$	$96 \times 48$	iconv3
upconv2	$4 \times 4$	2	128/64	$96 \times 48$	$192 \times 96$	iconv3
iconv2	$3 \times 3$	1	193/64	$192 \times 96$	$192 \times 96$	upconv2+pr3+conv2
pr2+loss2	$3 \times 3$	1	64/1	$192 \times 96$	$192 \times 96$	iconv2
upconv1	$4 \times 4$	2	64/32	$192 \times 96$	$384 \times 192$	iconv2
iconv1	$3 \times 3$	1	97/32	$384 \times 192$	$384 \times 192$	upconv1+pr2+conv1
pr1+loss1	$3 \times 3$	1	32/1	$384 \times 192$	$384 \times 192$	iconv1

Figure 7: DispNet architecture

# Technical Slides - Unsupervised Monocular Network Architecture

“Encoder”							“Decoder”						
layer	k	s	chns	in	out	input	layer	k	s	chns	in	out	input
conv1	7	2	3/32	1	2	left	upconv7	3	2	512/512	128	64	conv7b
conv1b	7	1	32/32	2	2	conv1	iconv7	3	1	1024/512	64	64	upconv7+conv6b
conv2	5	2	32/64	2	4	conv1b	upconv6	3	2	512/512	64	32	iconv7
conv2b	5	1	64/64	4	4	conv2	iconv6	3	1	1024/512	32	32	upconv6+conv5b
conv3	3	2	64/128	4	8	conv2b	upconv5	3	2	512/256	32	16	iconv6
conv3b	3	1	128/128	8	8	conv3	iconv5	3	1	512/256	16	16	upconv5+conv4b
conv4	3	2	128/256	8	16	conv3b	upconv4	3	2	256/128	16	8	iconv5
conv4b	3	1	256/256	16	16	conv4	iconv4	3	1	128/128	8	8	upconv4+conv3b
conv5	3	2	256/512	16	32	conv4b	disp4	3	1	128/2	8	8	iconv4
conv5b	3	1	512/512	32	32	conv5	upconv3	3	2	128/64	8	4	iconv4
conv6	3	2	512/512	32	64	conv5b	iconv3	3	1	130/64	4	4	upconv3+conv2b+disp4*
conv6b	3	1	512/512	64	64	conv6	disp3	3	1	64/2	4	4	iconv3
conv7	3	2	512/512	64	128	conv6b	upconv2	3	2	64/32	4	2	iconv3
conv7b	3	1	512/512	128	128	conv7	iconv2	3	1	66/32	2	2	upconv2+conv1b+disp3*
							disp2	3	1	32/2	2	2	iconv2
							upconv1	3	2	32/16	2	1	iconv2
							iconv1	3	1	18/16	1	1	upconv1+disp2*
							disp1	3	1	16/2	1	1	iconv1

Figure 8: Unsupervised Monocular Network architecture

# Technical Slides - Transposed Convolutions

- The input (bottom) is specially padded and then a traditional convolutional layer is applied

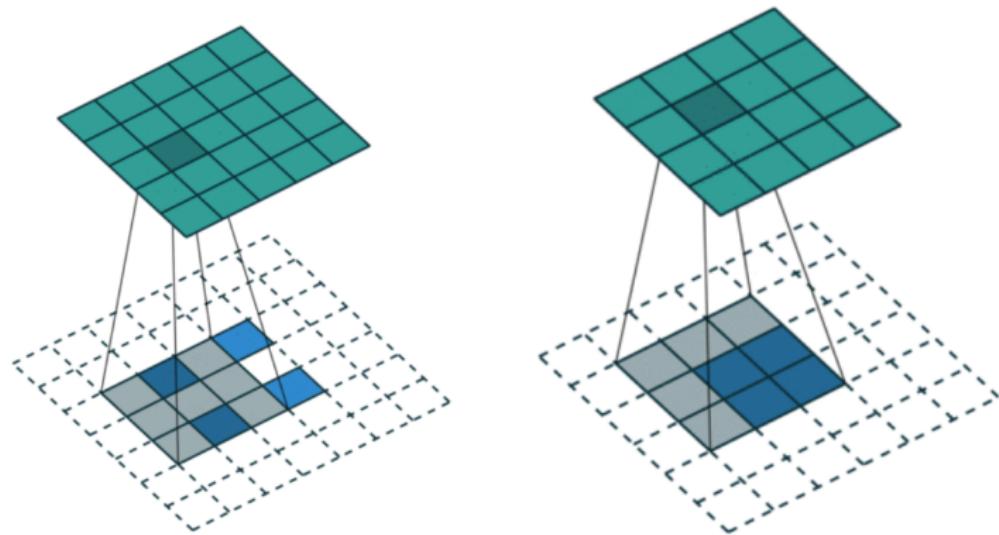


Figure 9: Transposed Convolutions

# Technical Slides - Structural Similarity (SSIM)<sup>9</sup>

$$SSIM(x, y) = \frac{(2\mu_x\mu_y + c_1)(2\sigma_{xy} + c_2)}{(\mu_x^2 + \mu_y^2 + c_1)(\sigma_x^2 + \sigma_y^2 + c_2)}$$

- $\mu_x$  the average of  $x$ ,  $\mu_y$  the average of  $y$ ;
- $\sigma_x^2$  the variance of  $x$ ,  $\sigma_y^2$  the variance of  $y$ ;
- $\sigma_{xy}$  the covariance of  $x$  and  $y$ ;
- $c1 = (k_1 L)^2$ ,  $c2 = (k_2 L)^2$  two variables to stabilize the division with weak denominator;
- $L$  the dynamic range of the pixel-values (usually  $2^{\# \text{bits-per-pixel}} - 1$ );
- $k_1 = 0.01$  and  $k_2 = 0.03$  by default.

---

<sup>9</sup>Zhou Wang et al. "Image Quality Assessment: From Error Visibility to Structural Similarity". In: *Trans. Img. Proc. (2004)*.

# Technical Slides - Smooth L1 Loss

$$|x|_{smooth} = \begin{cases} 0.5x^2, & \text{if } |x| \leq 1 \\ |x| - 0.5, & \text{otherwise} \end{cases}$$

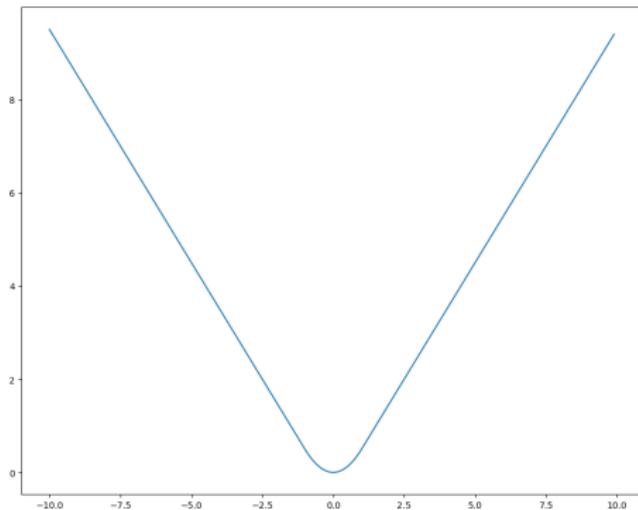


Figure 10: Plot for smooth L1