

Machine Learning in Geometric Computer Vision

Part #1

Machine Learning in Geometric Computer Vision

Part #1

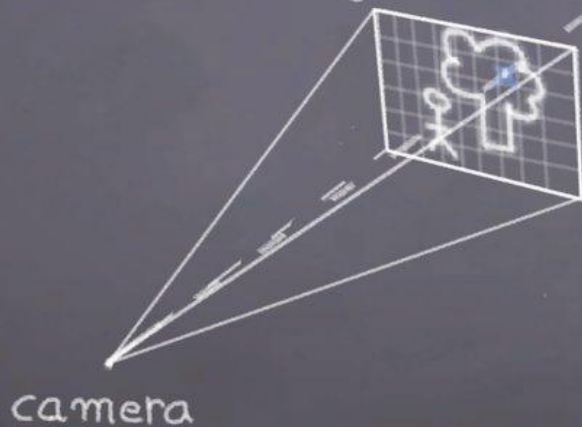
Machine Learning in Geometric Computer Vision

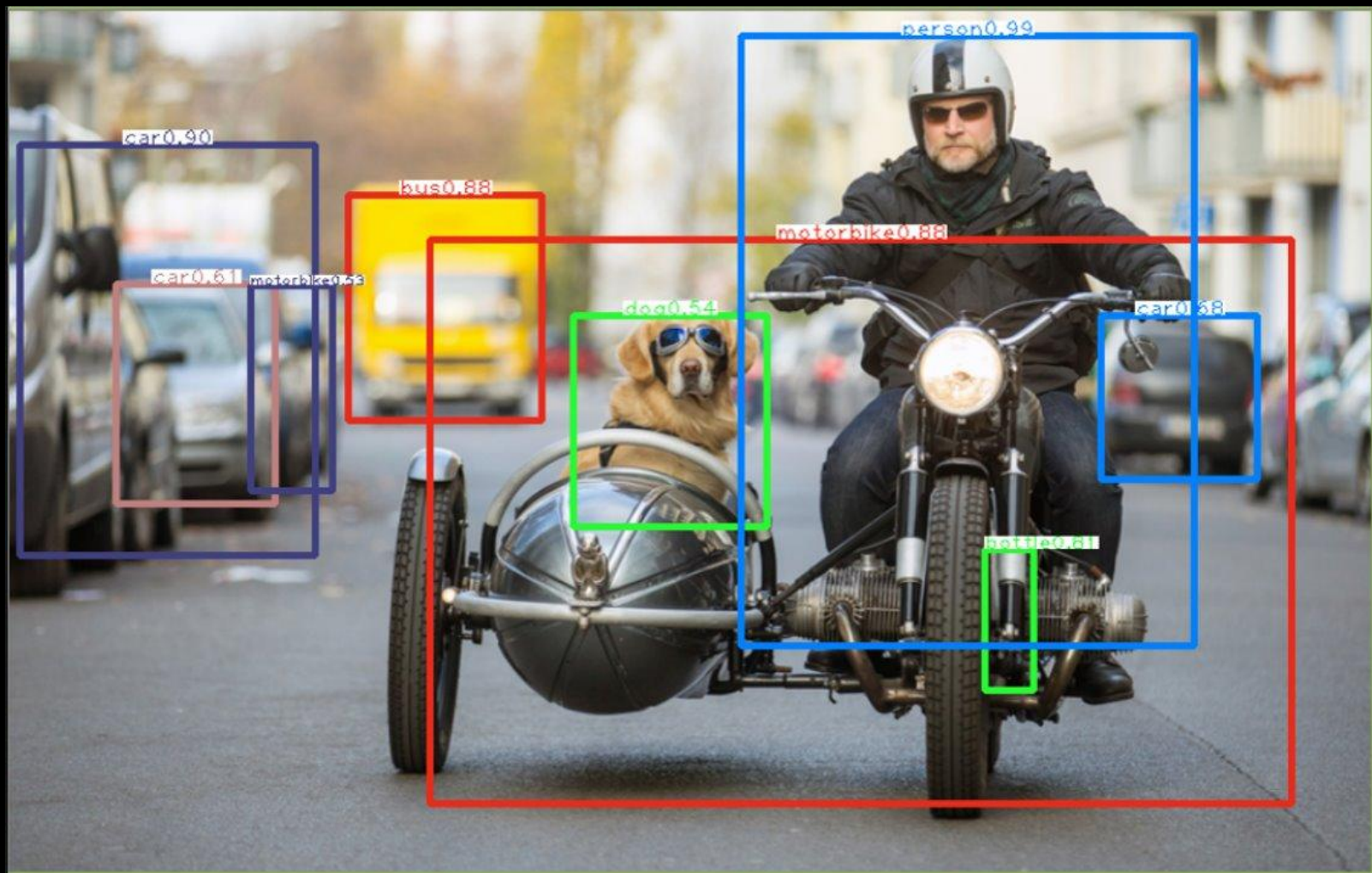
Part #1

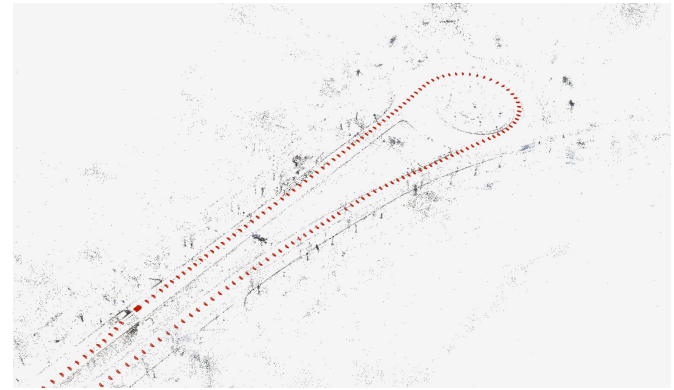
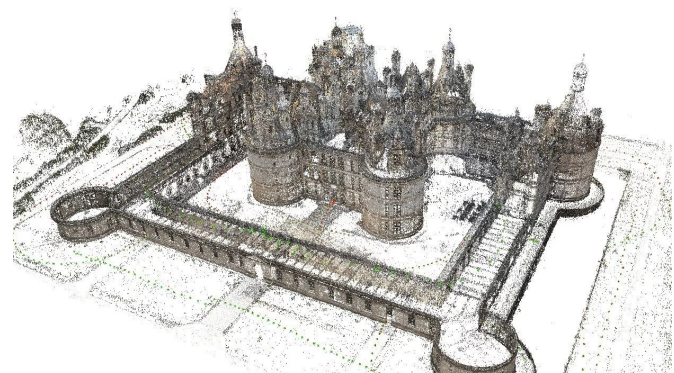
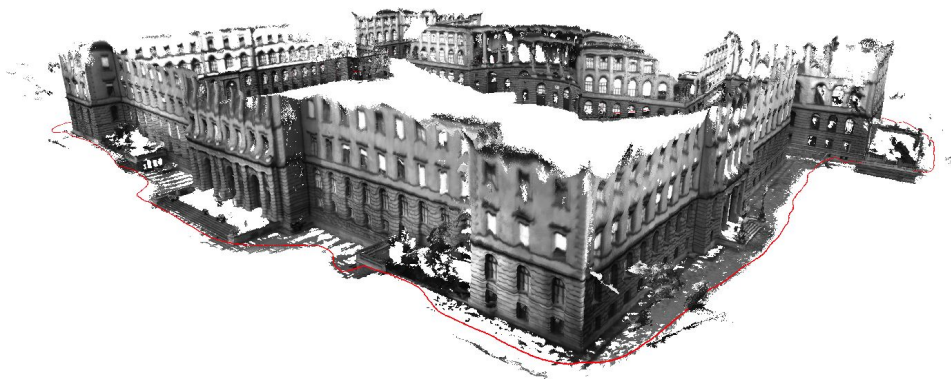
Computer Vision

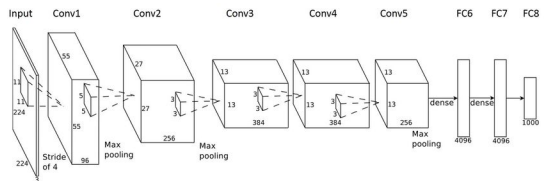
3D World

2D Image



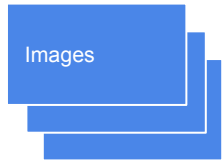




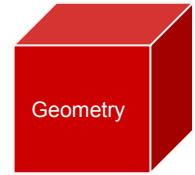


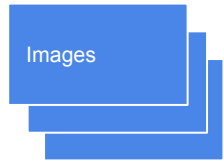


* Many great minds have been left out of this slide. Images have been resized for appearance purposes and their sizes don't correspond to any known metric of the individual's contribution to society.

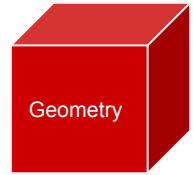


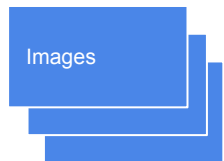
Geometric Algorithm
crafted by Humans





Machine Learned
Algorithm

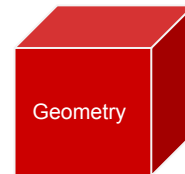


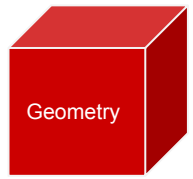
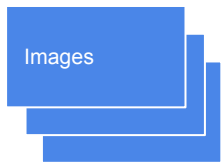


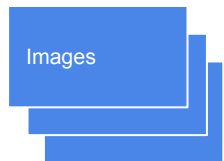
Machine Learned
Algorithm



Geometric Algorithm
crafted by Humans



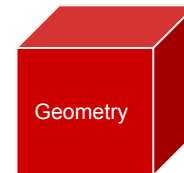


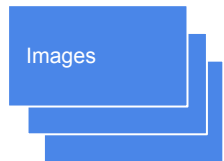


Geometric Algorithm
crafted by Humans



Machine Learned
Algorithm

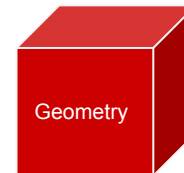




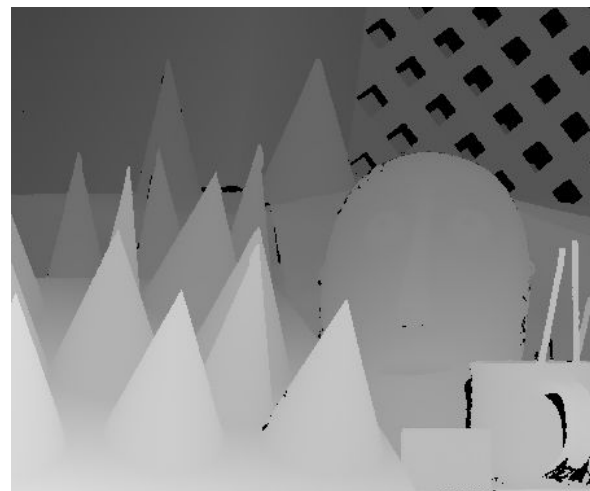
Geometric Algorithm
crafted by Humans



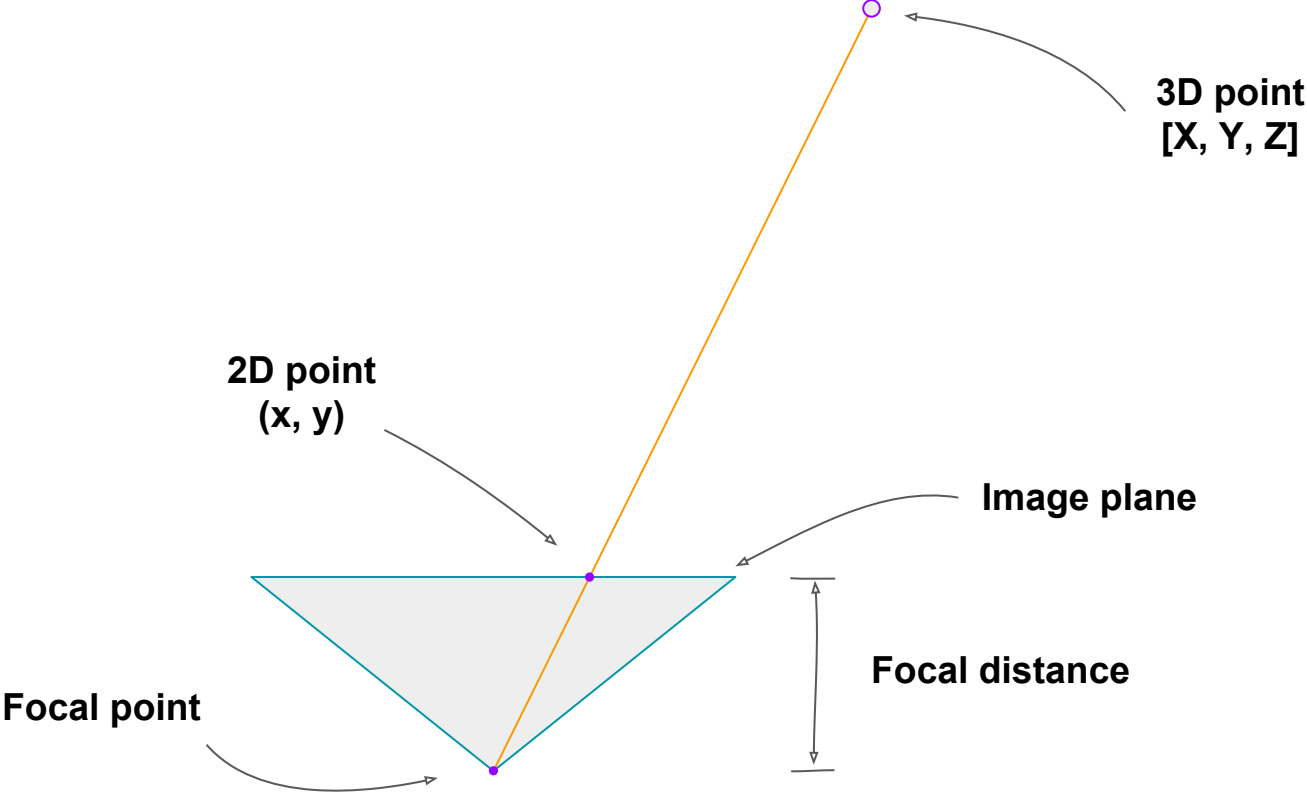
Machine Learned
Algorithm

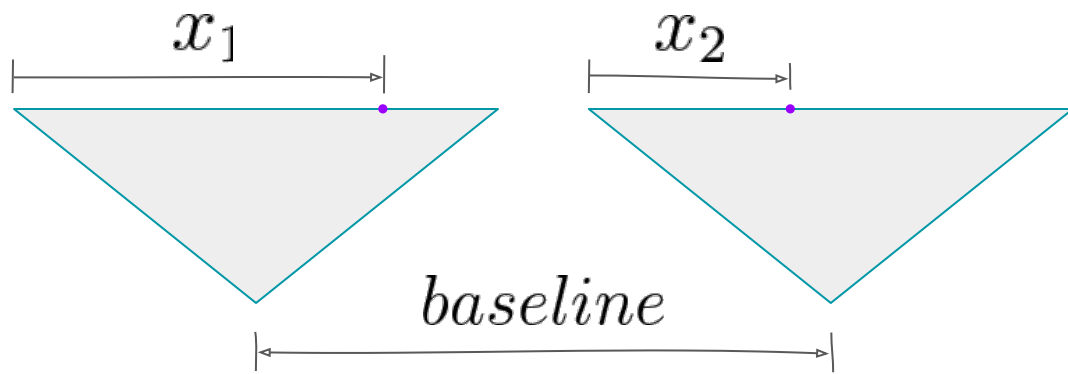


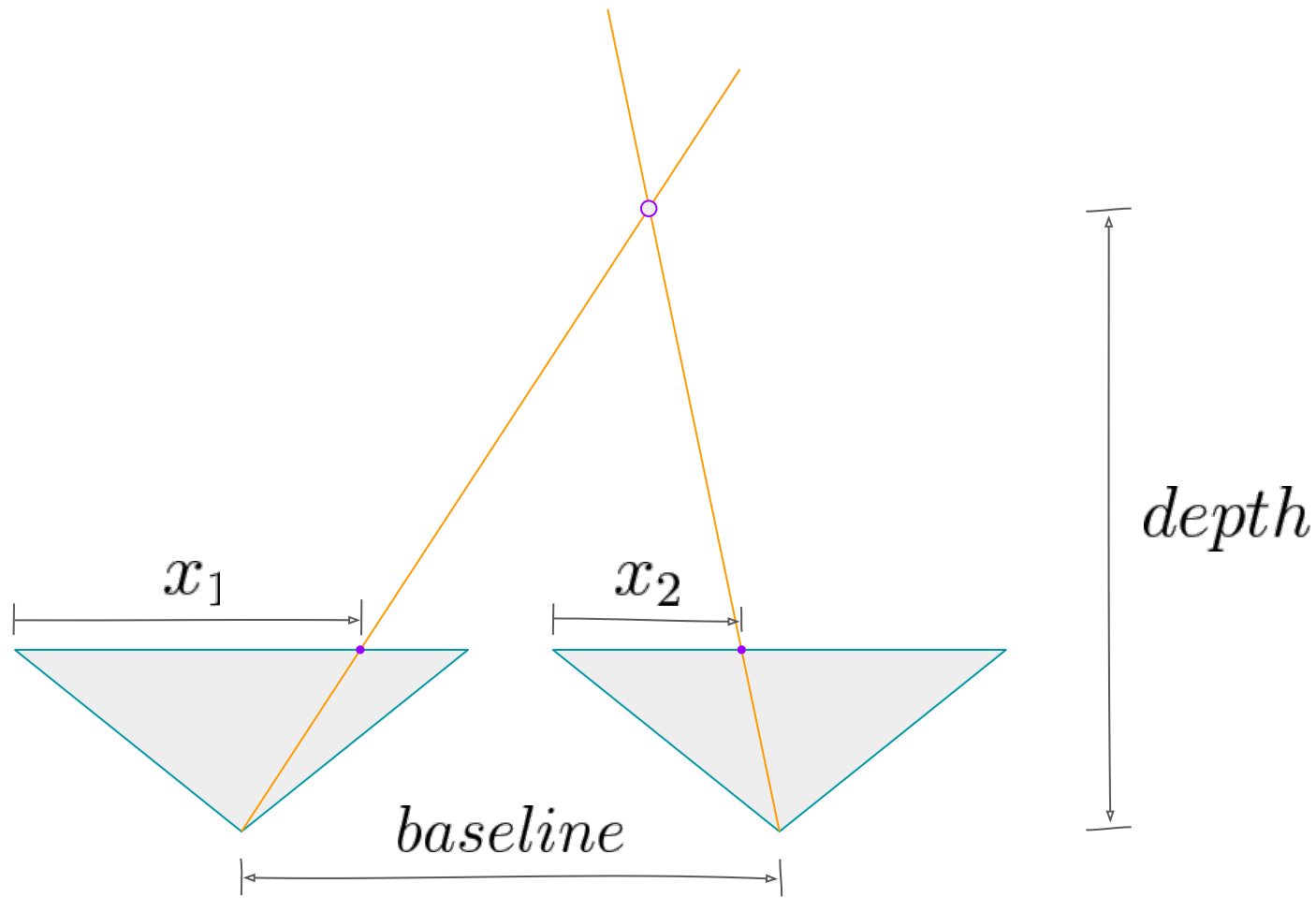
1. **Dense stereo correspondence**
2. **Visual odometry**
3. **Depth map fusion**
4. **SLAM**

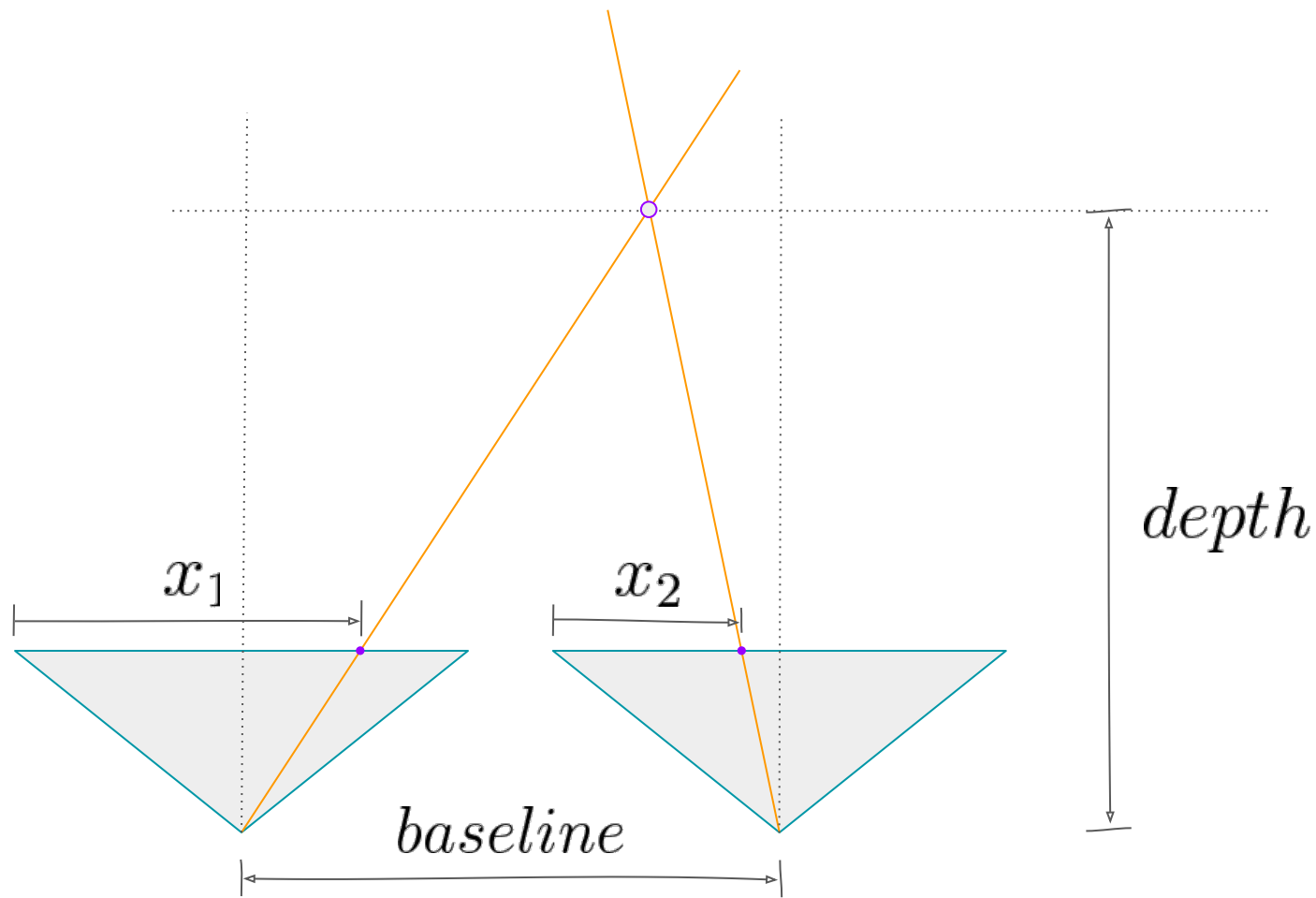


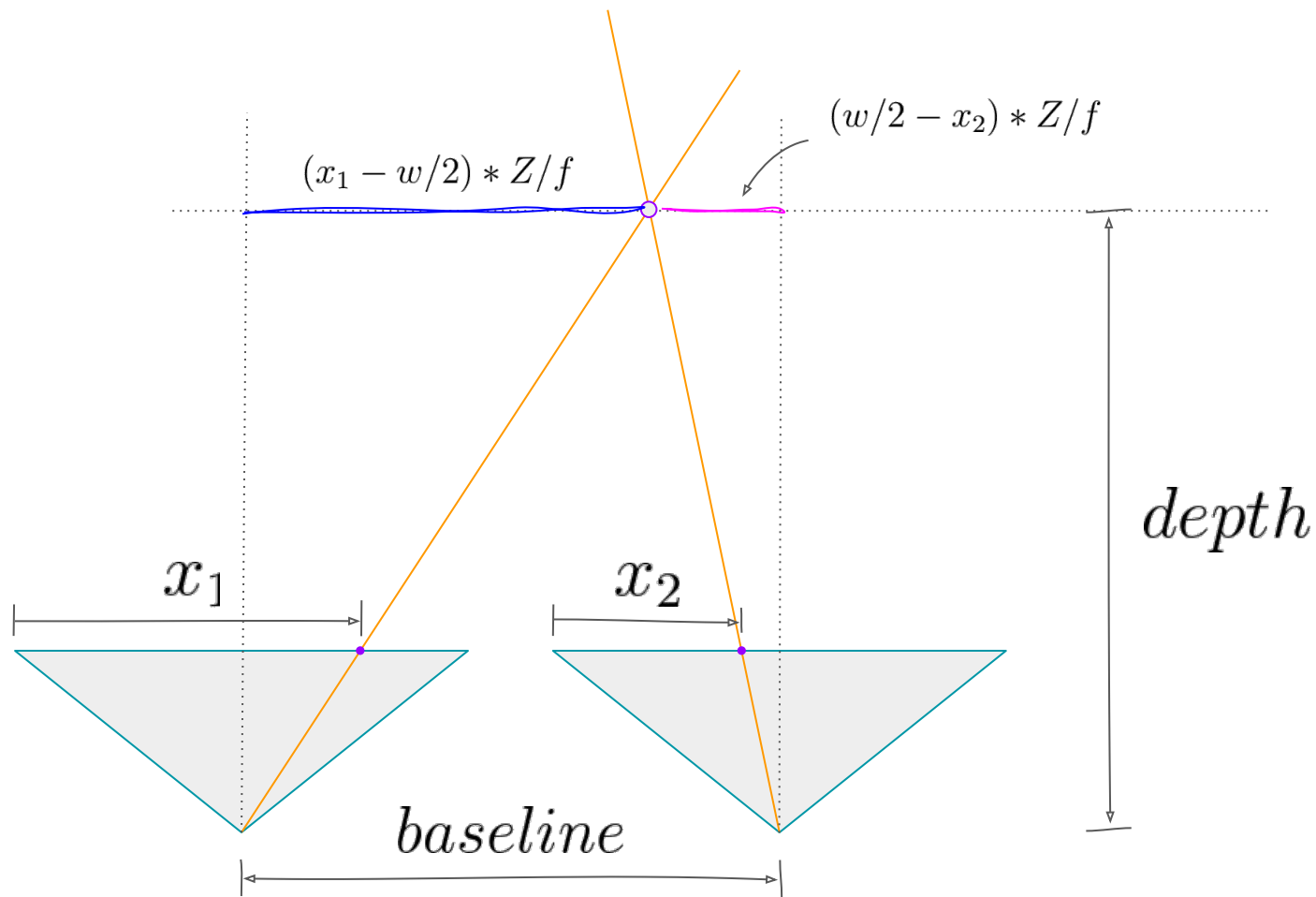
Pinhole camera model



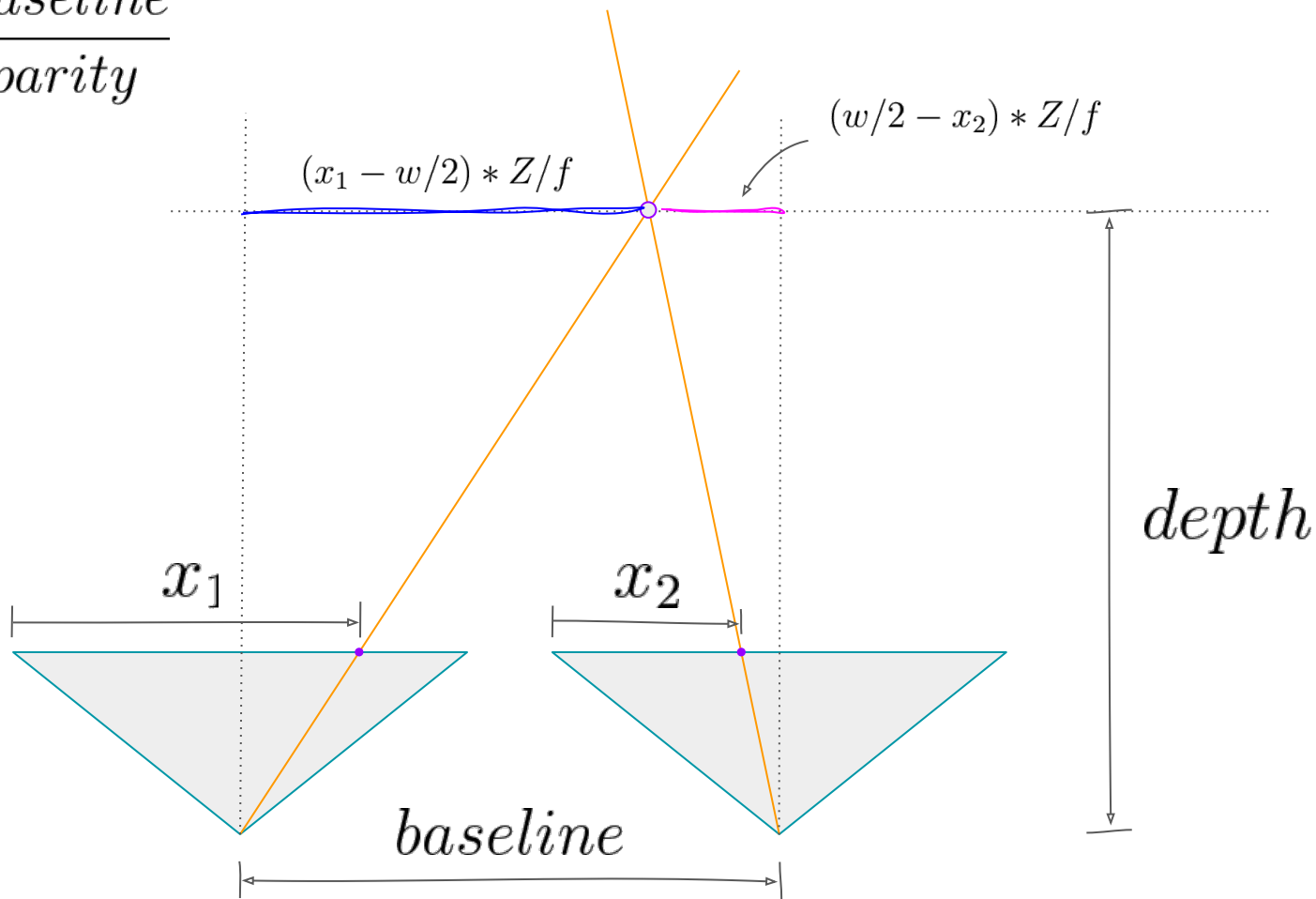






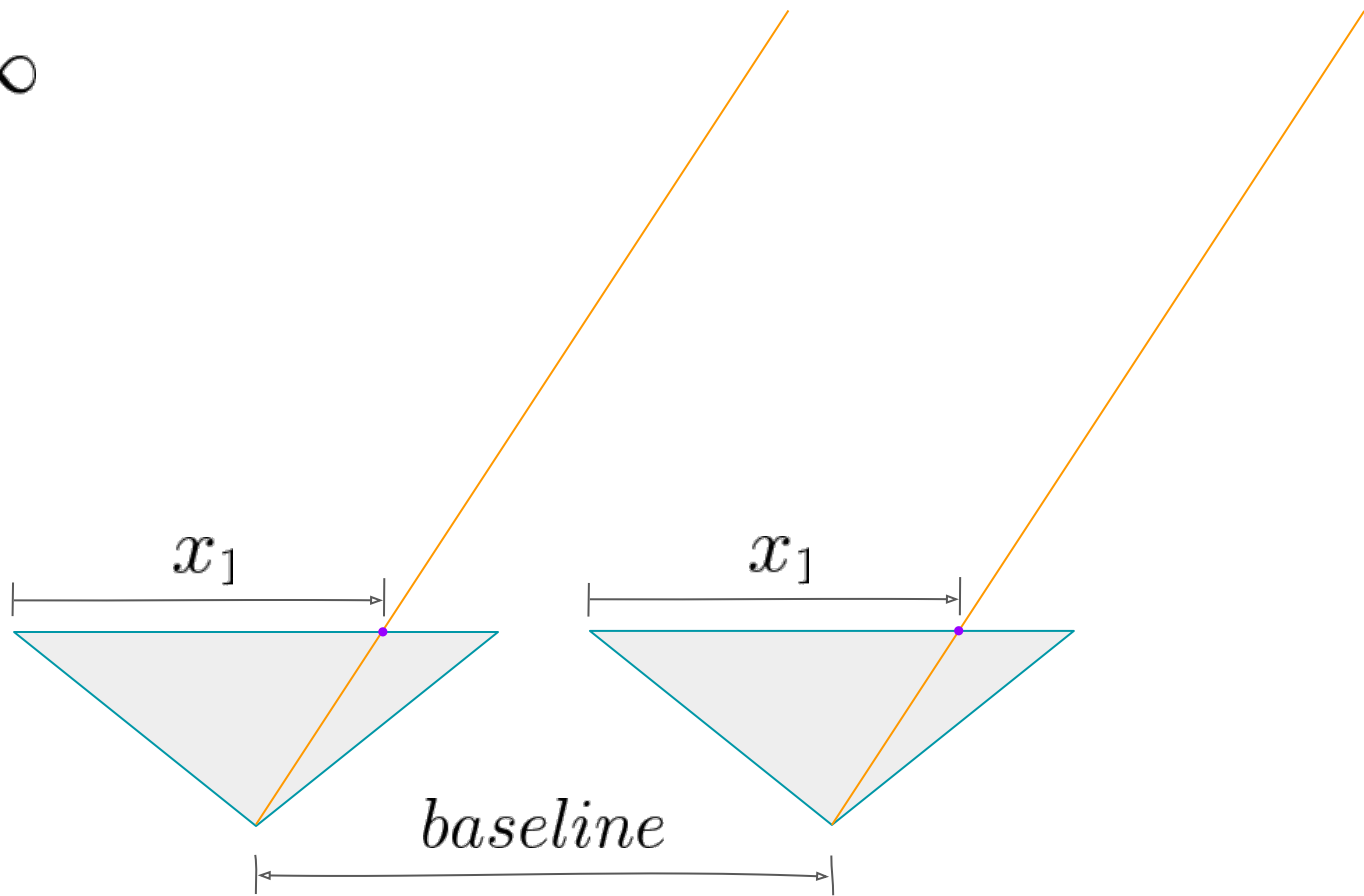


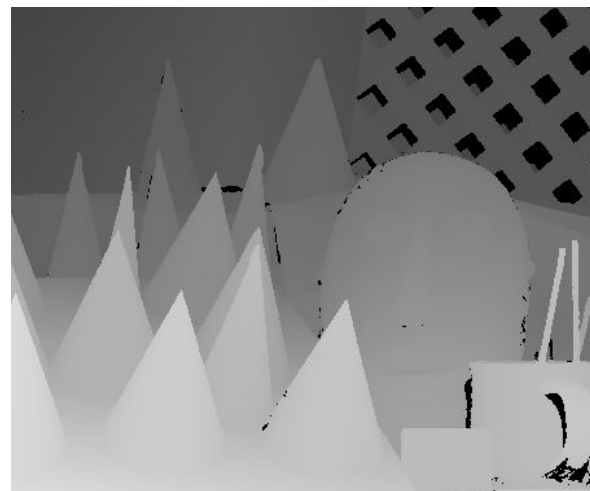
$$\text{depth} = \frac{f * \text{baseline}}{\text{disparity}}$$

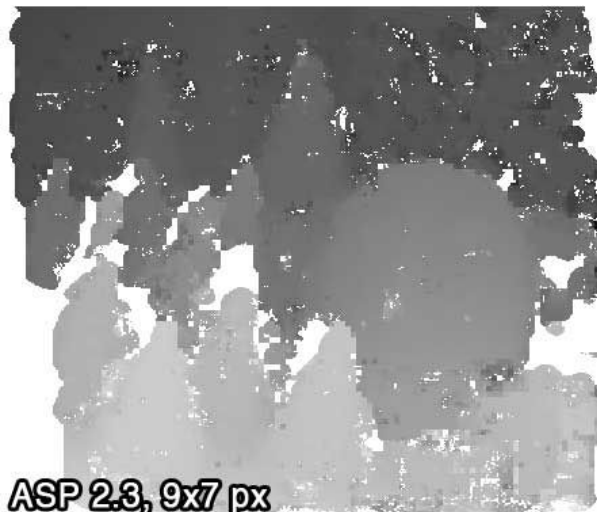


disparity = 0

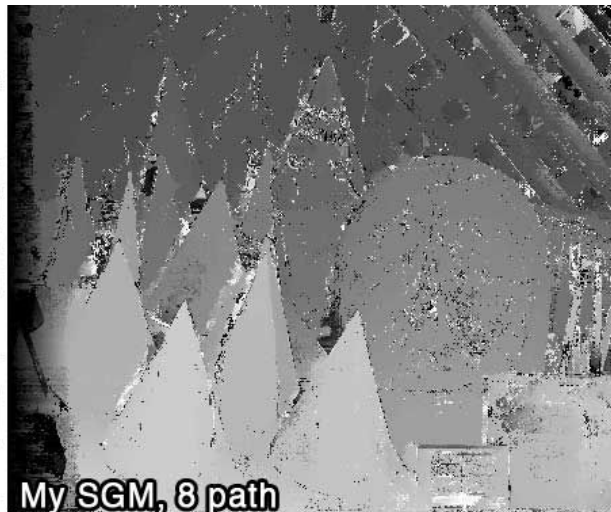
depth = ∞



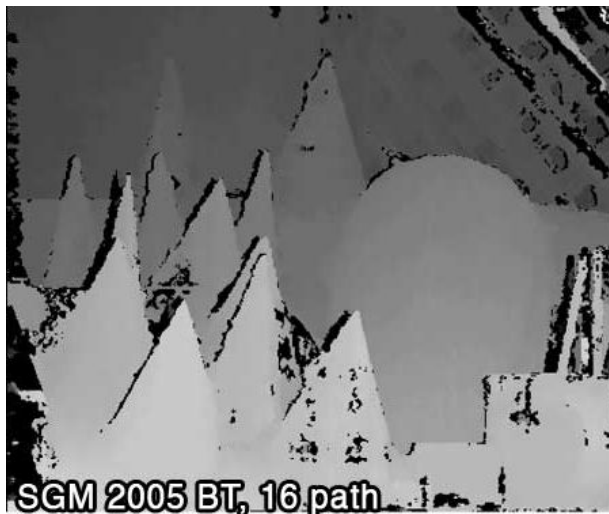




ASP 2.3, 9x7 px



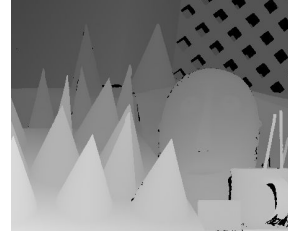
My SGM, 8 path

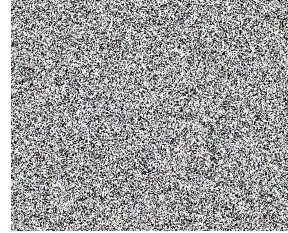


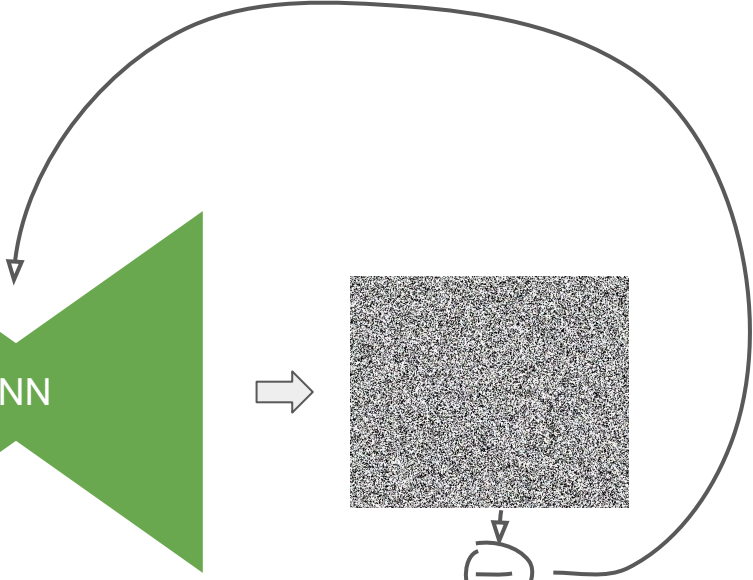
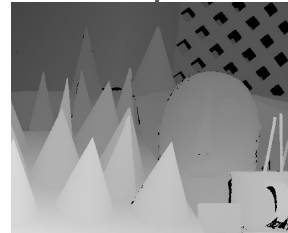
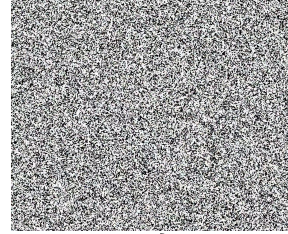
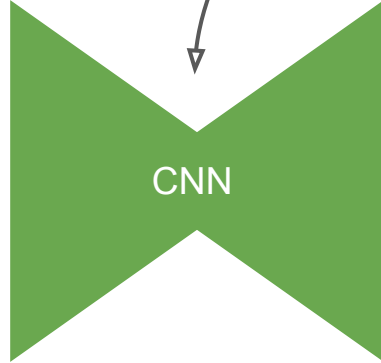
SGM 2005 BT, 16 path

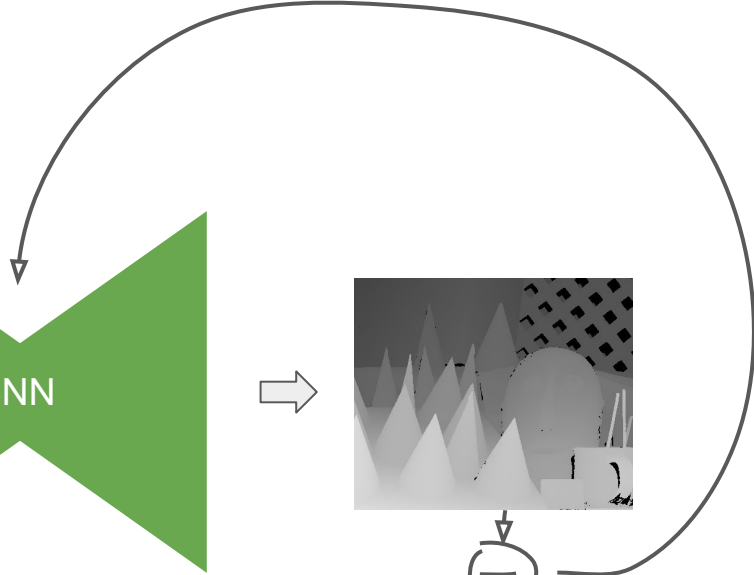
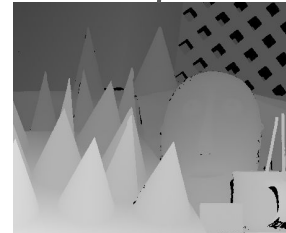
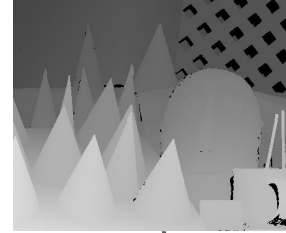
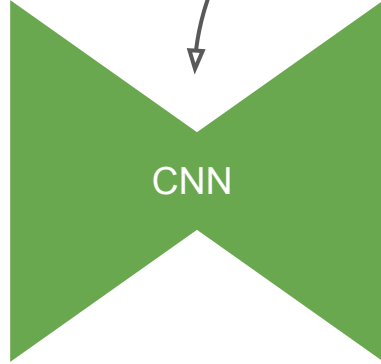


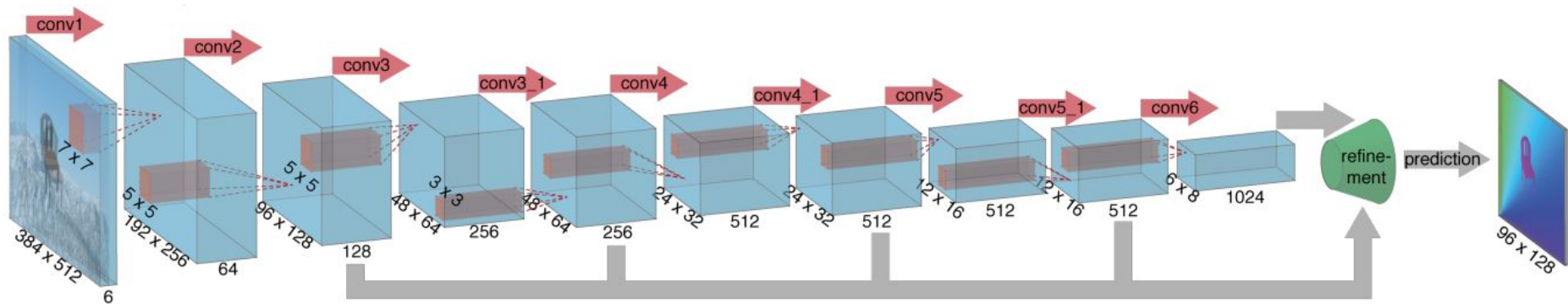
Ground Truth



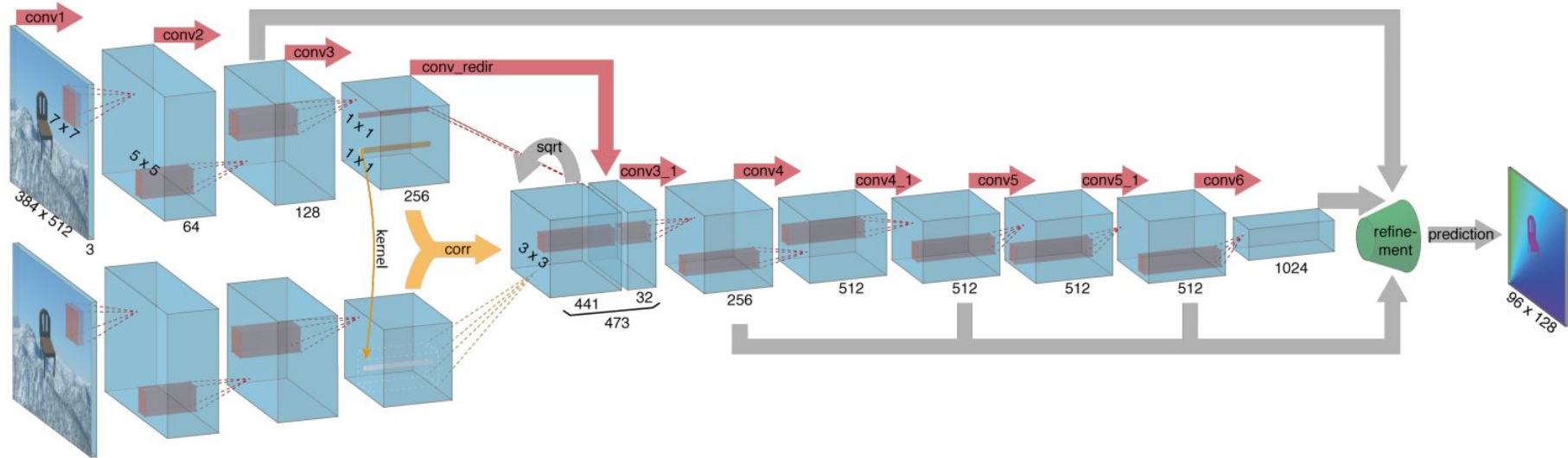








FlowNetCorr



RGB image (L)



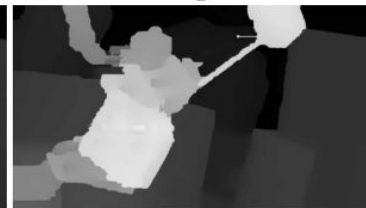
Disparity GT



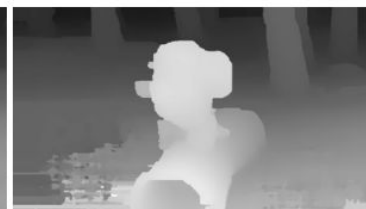
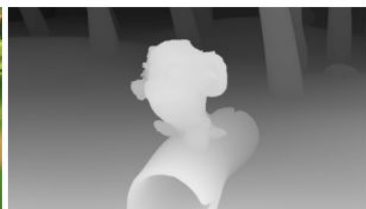
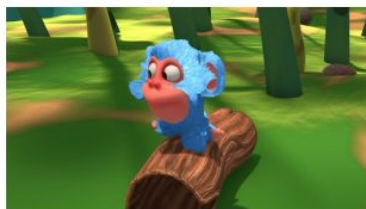
DispNetCorr1D



MC-CNN prediction



SGM prediction



A Large Dataset to Train Convolutional Networks for Disparity, Optical Flow, and Scene Flow Estimation

Nikolaus Mayer*¹, Eddy Ilg*¹, Philip Häusser*², Philipp Fischer*^{1†}

¹University of Freiburg ²Technical University of Munich

¹{mayern, ilg, fischer}@cs.uni-freiburg.de ²haeusser@cs.tum.edu

Daniel Cremers
Technical University of Munich
cremers@tum.de

Alexey Dosovitskiy, Thomas Brox
University of Freiburg
{dosovits, brox}@cs.uni-freiburg.de

Abstract

Recent work has shown that optical flow estimation can be formulated as a supervised learning task and can be successfully solved with convolutional networks. Training of the so-called FlowNet was enabled by a large synthetically generated dataset. The present paper extends the concept of optical flow estimation via convolutional networks to disparity and scene flow estimation. To this end, we propose three synthetic stereo video datasets with sufficient realism, variation, and size to successfully train large networks. Our datasets are the first large-scale datasets to enable training and evaluating scene flow methods. Besides the datasets, we present a convolutional network for real-time disparity estimation that provides state-of-the-art results. By combining a flow and disparity estimation network and training it jointly, we demonstrate the first scene flow estimation with a convolutional network.



Figure 1. Our datasets provide over 35 000 stereo frames with dense ground truth for optical flow, disparity and disparity change, as well as other data such as object segmentation.

A Large Dataset to Train Convolutional Networks for Disparity, Optical Flow, and Scene Flow Estimation

Nikolaus Mayer*¹, Eddy Ilg*¹, Philip Häusser*², Philipp Fischer*^{1†}

¹University of Freiburg ²Technical University of Munich

¹{mayern, ilg, fischer}@cs.uni-freiburg.de ²haeusser@cs.tum.edu

Daniel Cremers
Technical University of Munich
cremers@tum.de

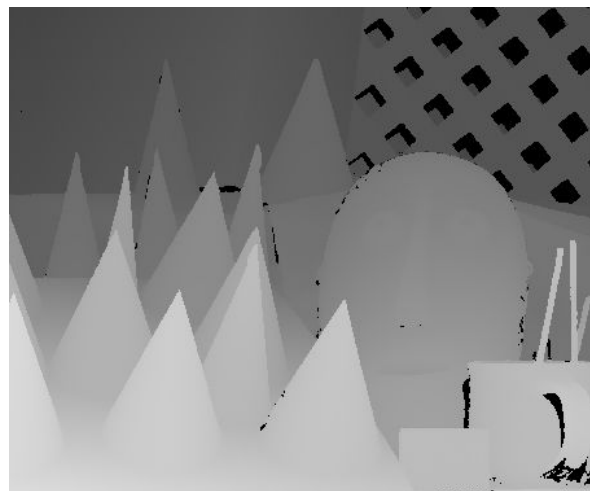
Alexey Dosovitskiy, Thomas Brox
University of Freiburg
{dosovits, brox}@cs.uni-freiburg.de

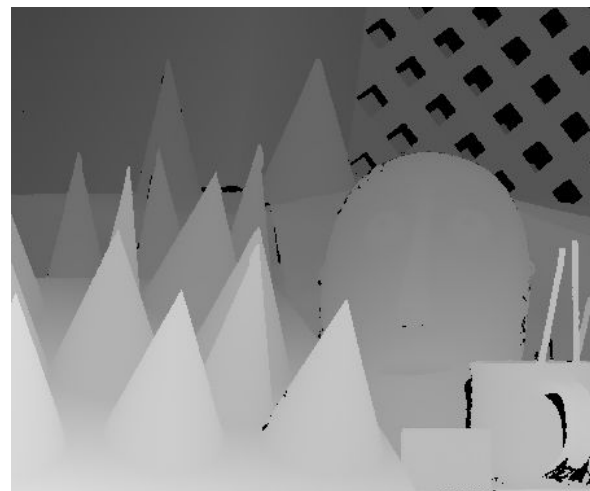
Abstract

Recent work has shown that optical flow estimation can be formulated as a supervised learning task and can be successfully solved with convolutional networks. Training of the so-called FlowNet was enabled by a large synthetically generated dataset. The present paper extends the concept of optical flow estimation via convolutional networks to disparity and scene flow estimation. To this end, we propose three synthetic stereo video datasets with sufficient realism, variation, and size to successfully train large networks. Our datasets are the first large-scale datasets to enable training and evaluating scene flow methods. Besides the datasets, we present a convolutional network for real-time disparity estimation that provides state-of-the-art results. By combining a flow and disparity estimation network and training it jointly, we demonstrate the first scene flow estimation with a convolutional network.



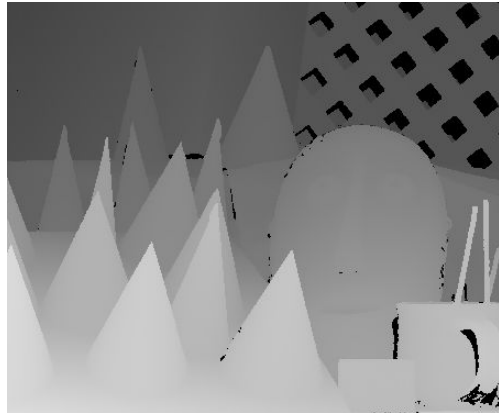
Figure 1. Our datasets provide over 35 000 stereo frames with dense ground truth for optical flow, disparity and disparity change, as well as other data such as object segmentation.





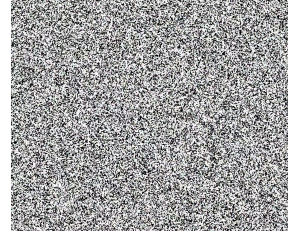


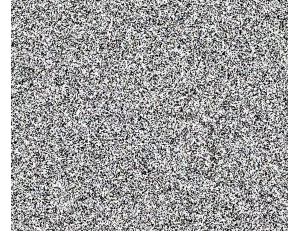
+



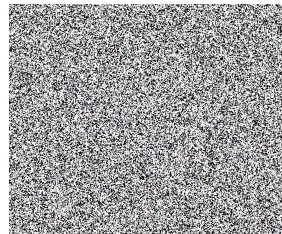
=







**Right
disparity
estimate**

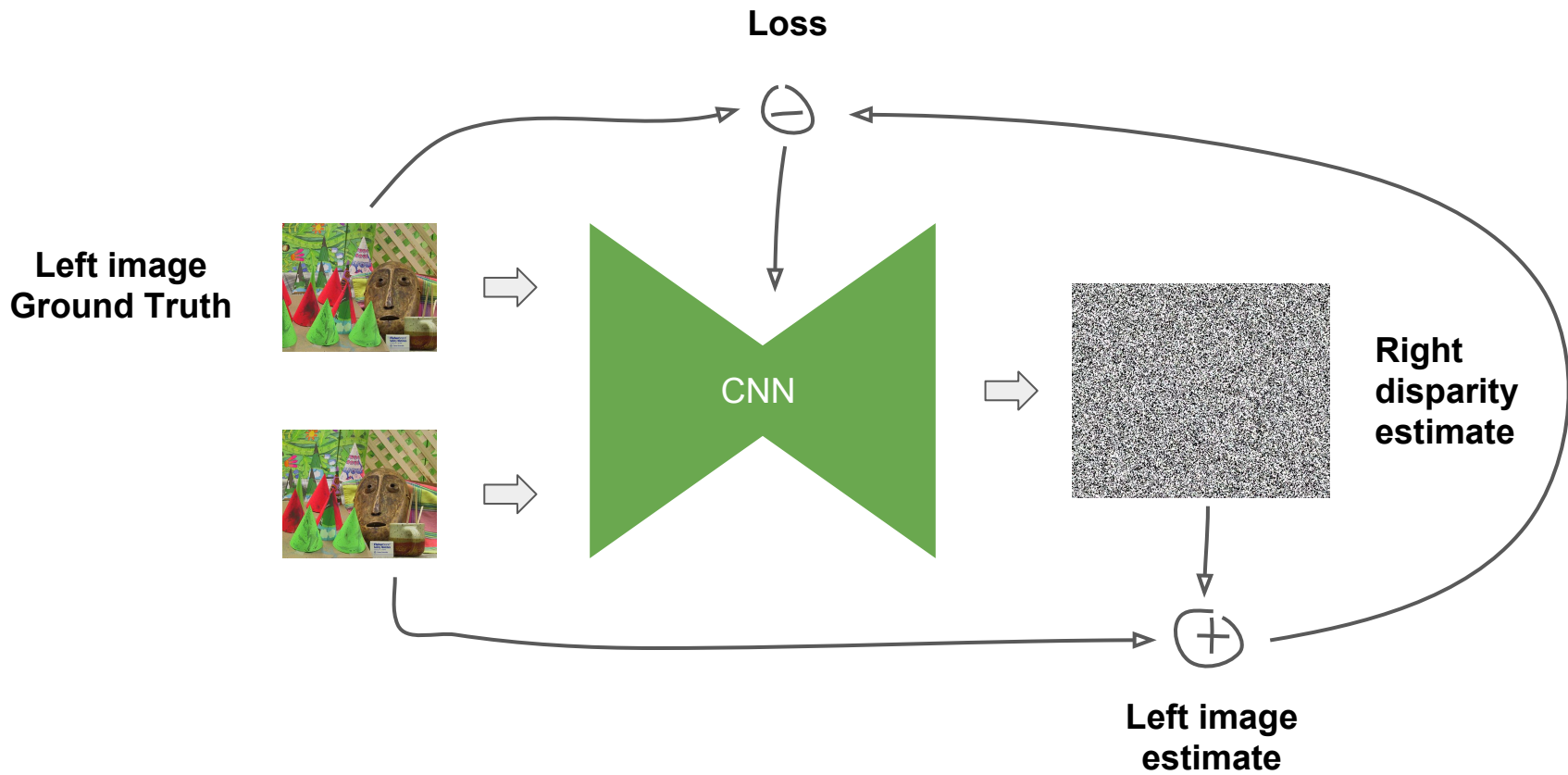


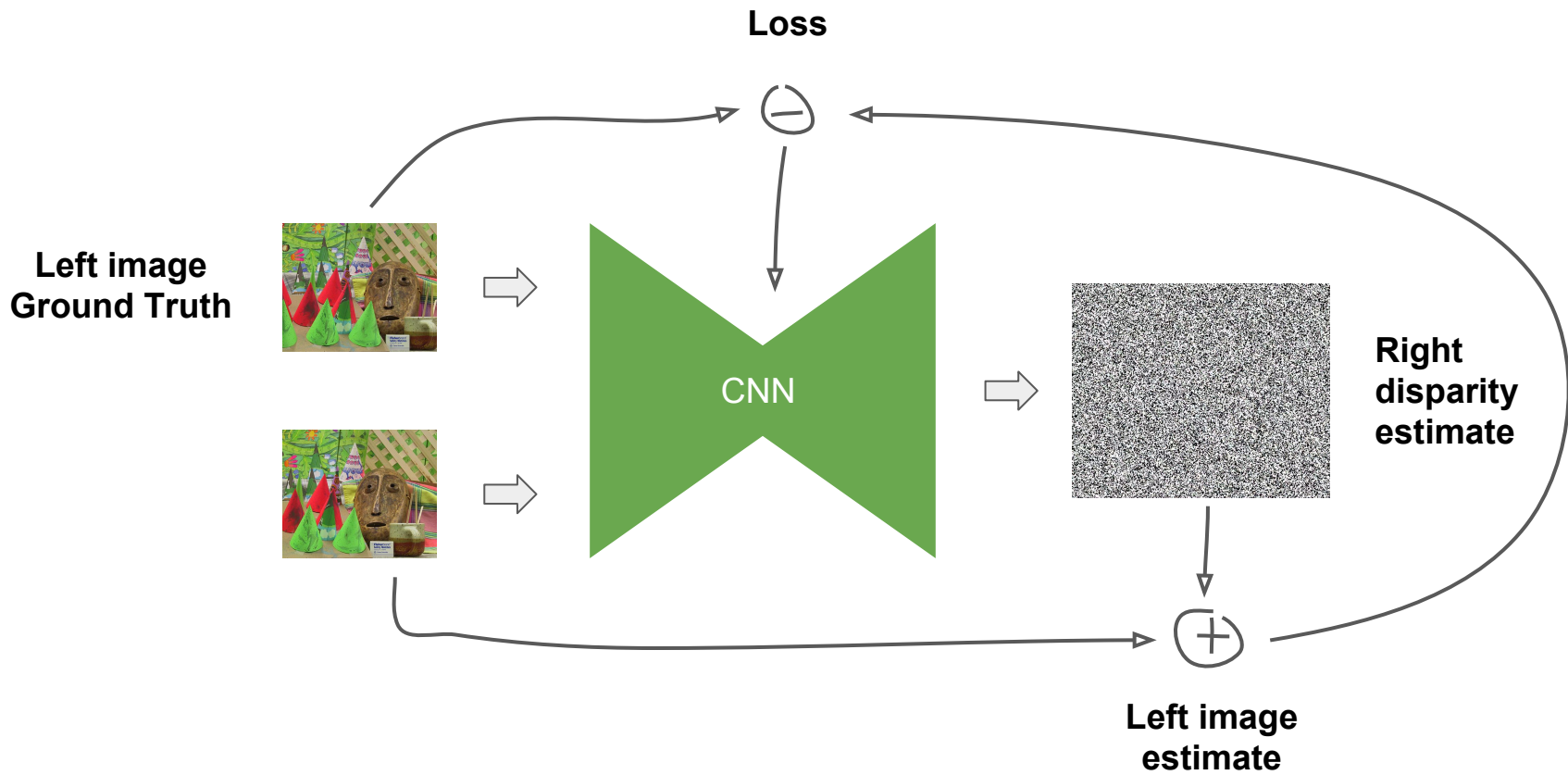
**Right
disparity
estimate**

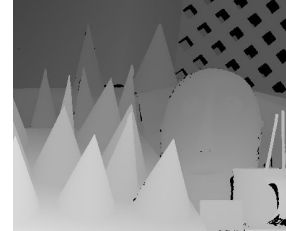


**Left image
estimate**









Unsupervised XXXXXXXXXX Depth Estimation with Left-Right Consistency

Clément Godard

Oisín Mac Aodha

Gabriel J. Brostow

University College London

<http://visual.cs.ucl.ac.uk/pubs/XXXXXXXXXX/>

Abstract

Learning based methods have shown very promising results for the task of XXXXXXXXXX. However, most existing approaches treat depth prediction as a supervised regression problem and as a result, require vast quantities of corresponding ground truth depth data for training. Just recording quality depth data in a range of environments is a challenging problem. In this paper, we innovate beyond existing approaches, replacing the use of explicit depth data during training with easier-to-obtain binocular stereo footage.

We propose a novel training objective that enables our convolutional neural network to learn to perform single image depth estimation, despite the absence of ground truth depth data. Exploiting epipolar geometry constraints, we generate disparity images by training our network with an image reconstruction loss. We show that solving for image reconstruction alone results in poor quality depth images. To overcome this problem, we propose a novel training loss that enforces consistency between the disparities produced relative to both the left and right

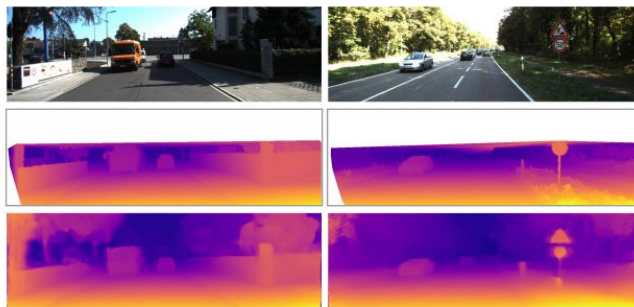
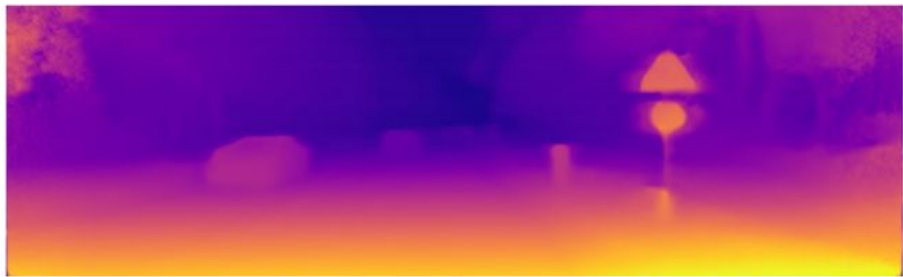
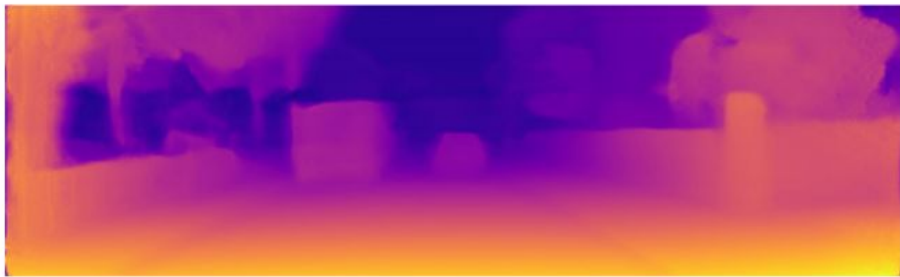
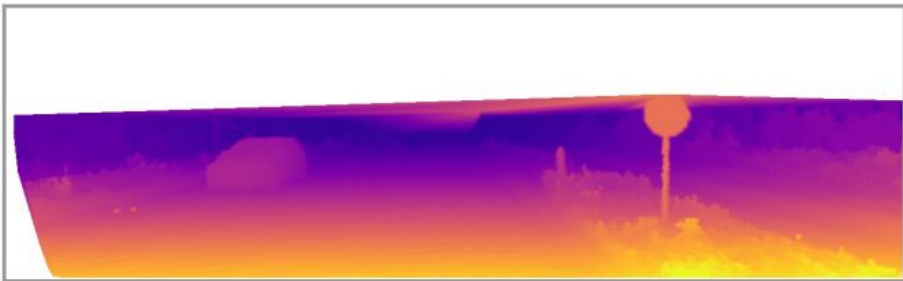
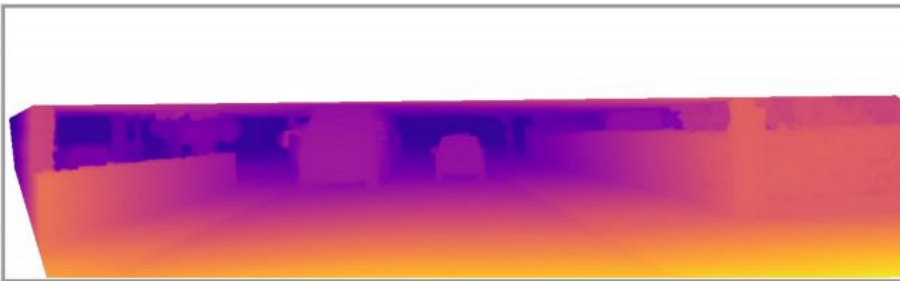
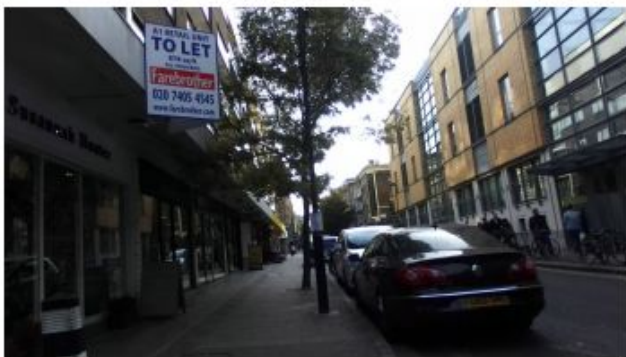
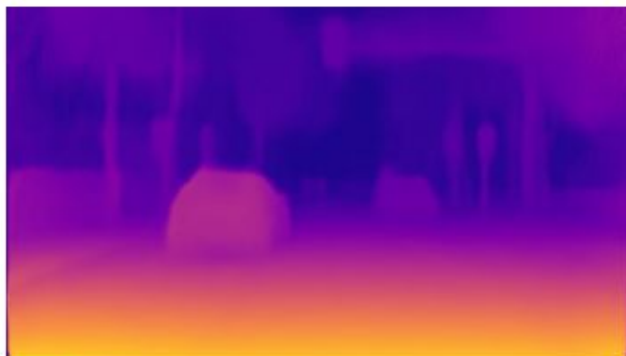
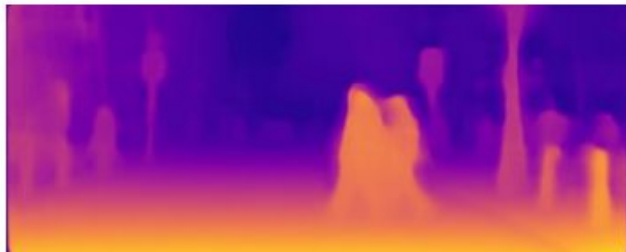


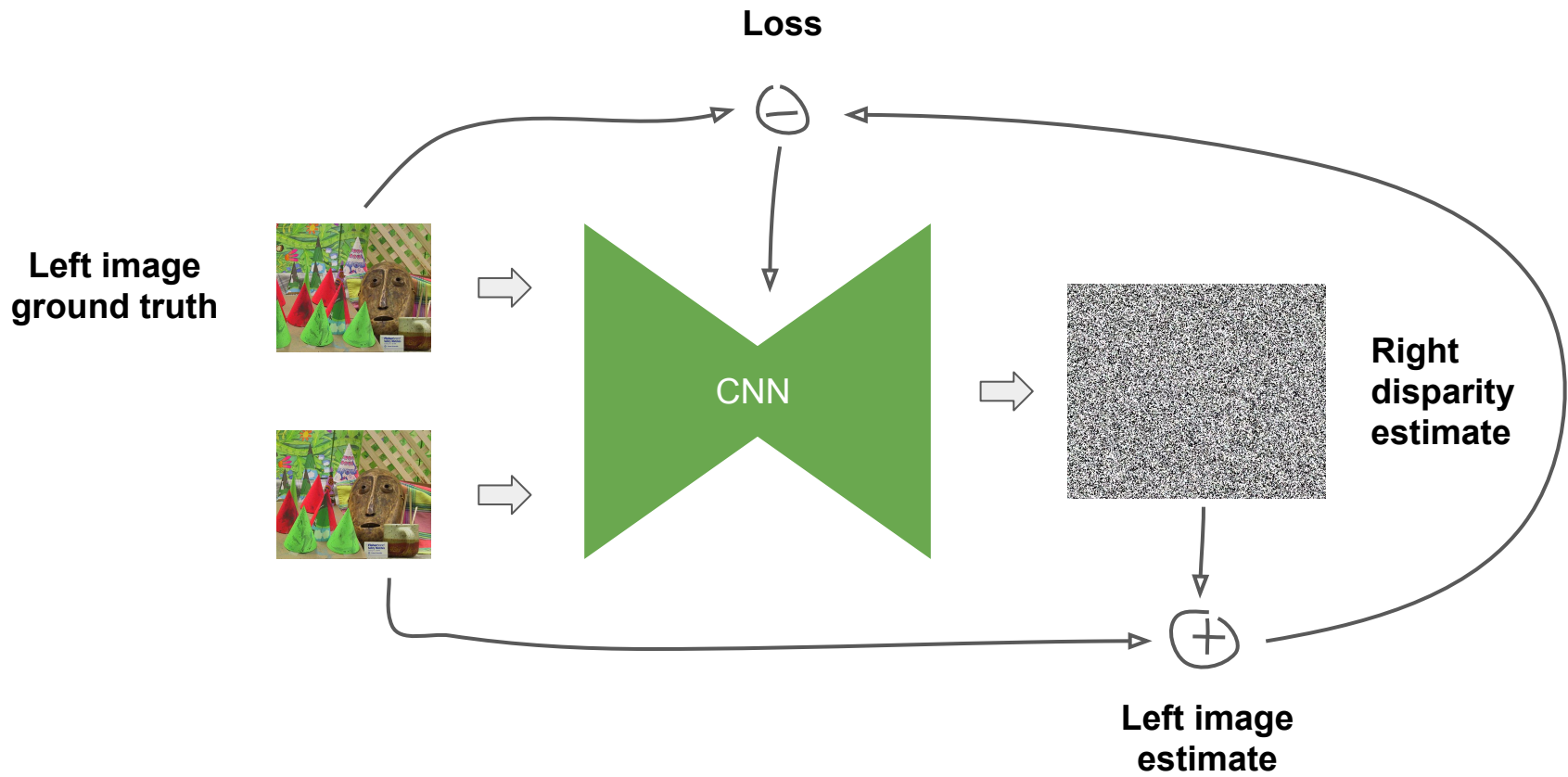
Figure 1. Our depth prediction results on KITTI 2015. Top to bottom: input image, ground truth disparities, and our result. Our method is able to estimate depth for thin structures such as street signs and poles.

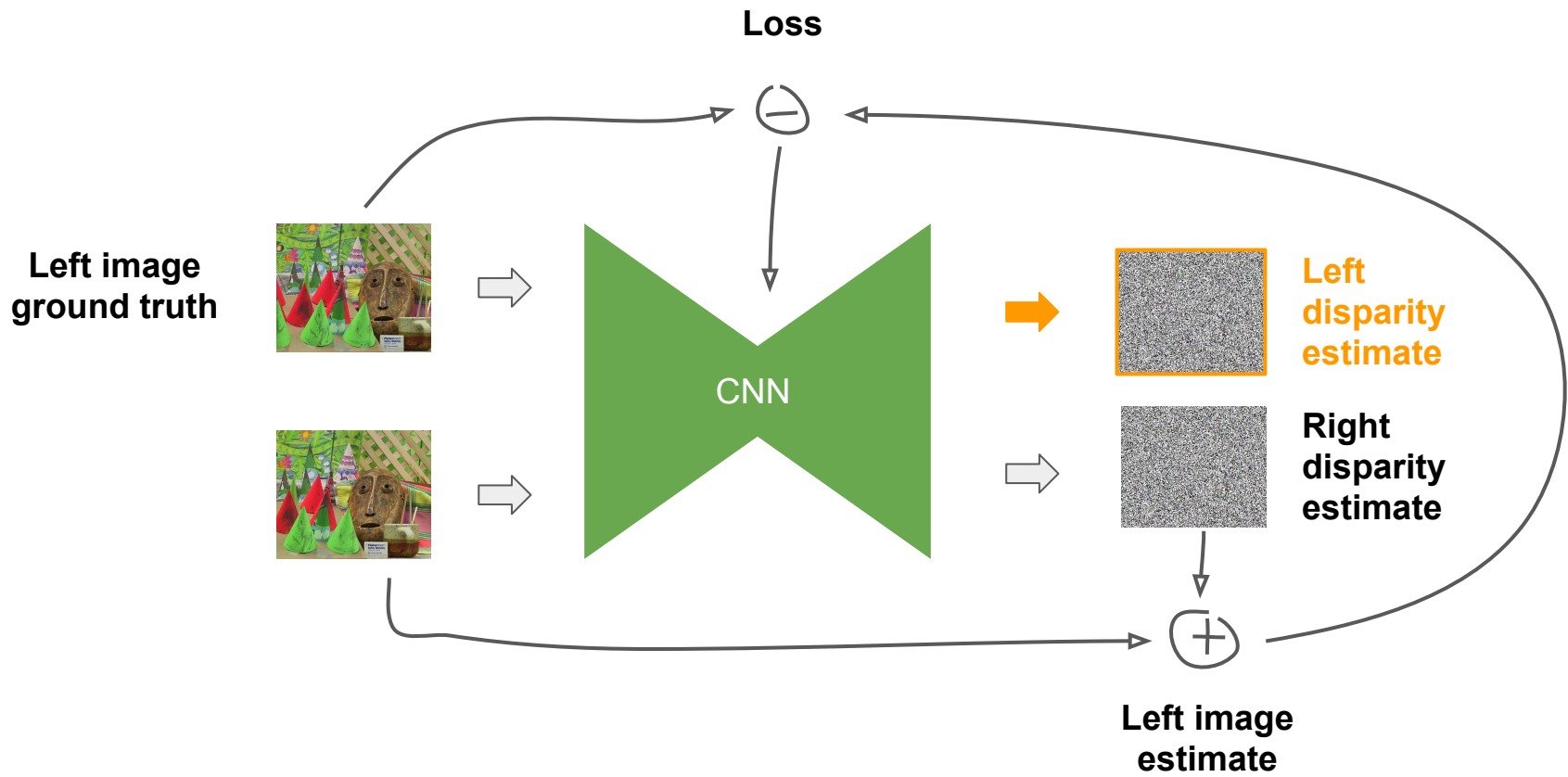
have been restricted to scenes where large image collections and their corresponding pixel depths are available.

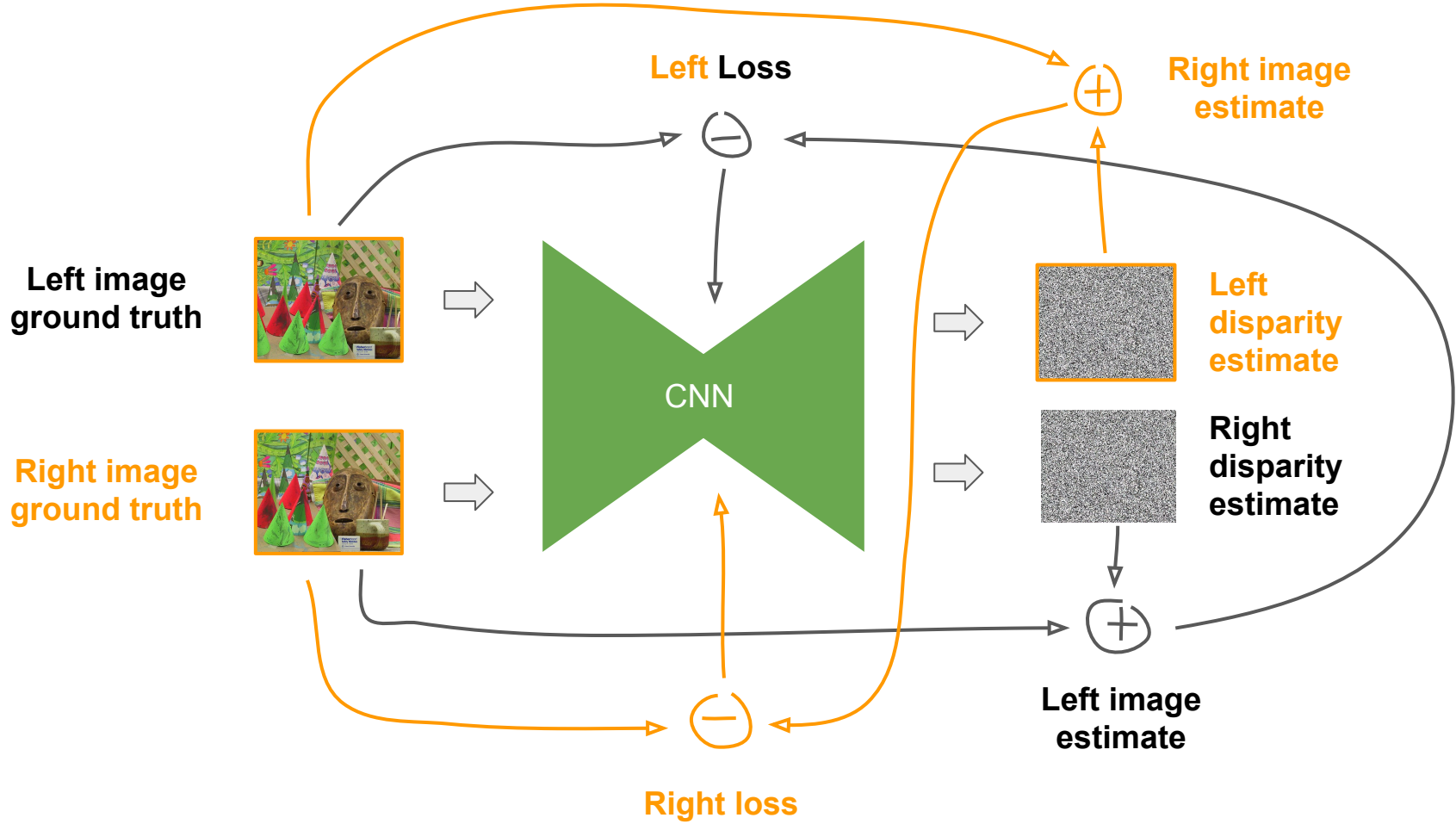
Understanding the shape of a scene from a single image, independent of its appearance, is a fundamental problem in machine perception. There are many applications such as synthetic object insertion in computer graphics [29], synthetic

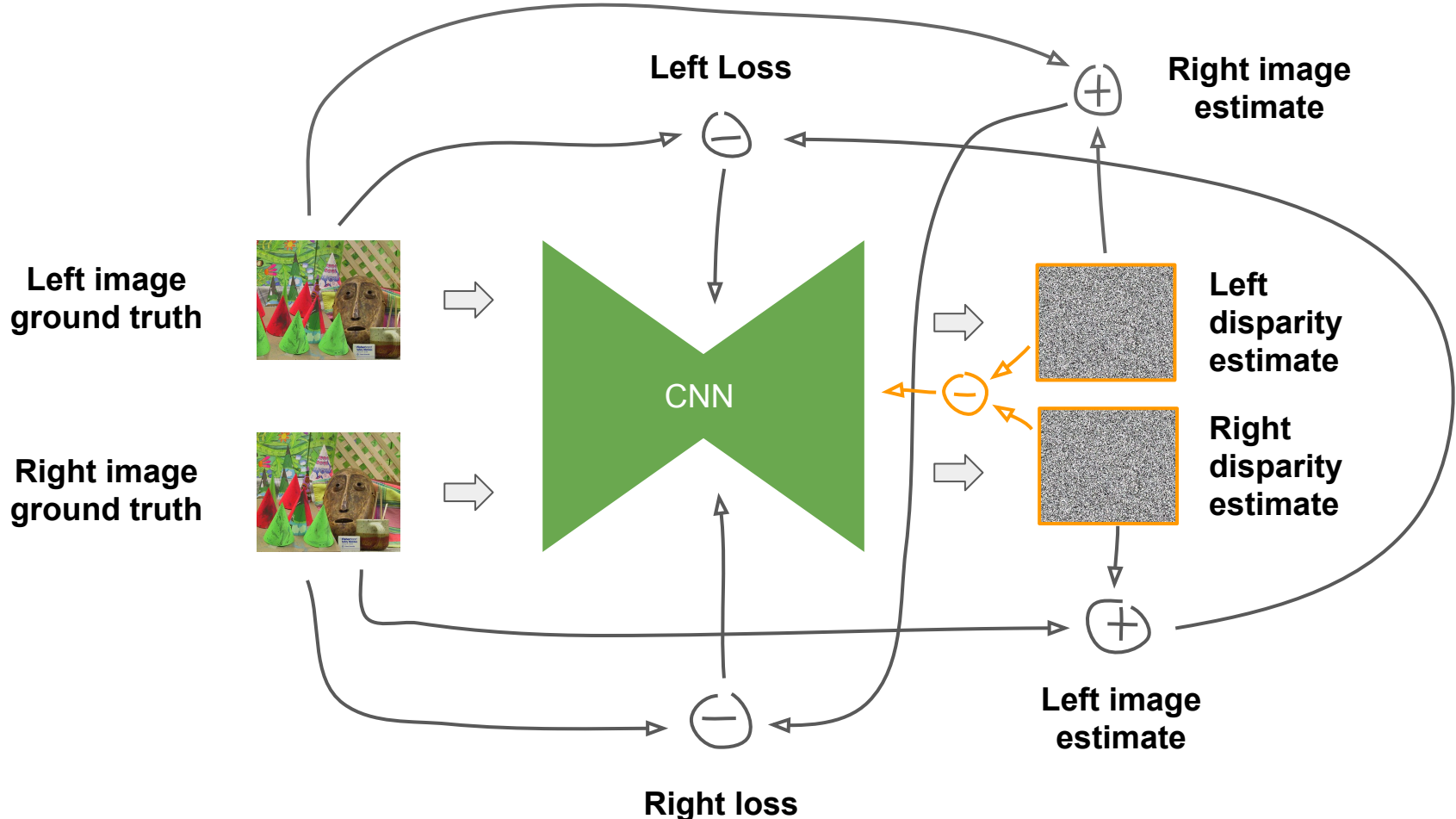


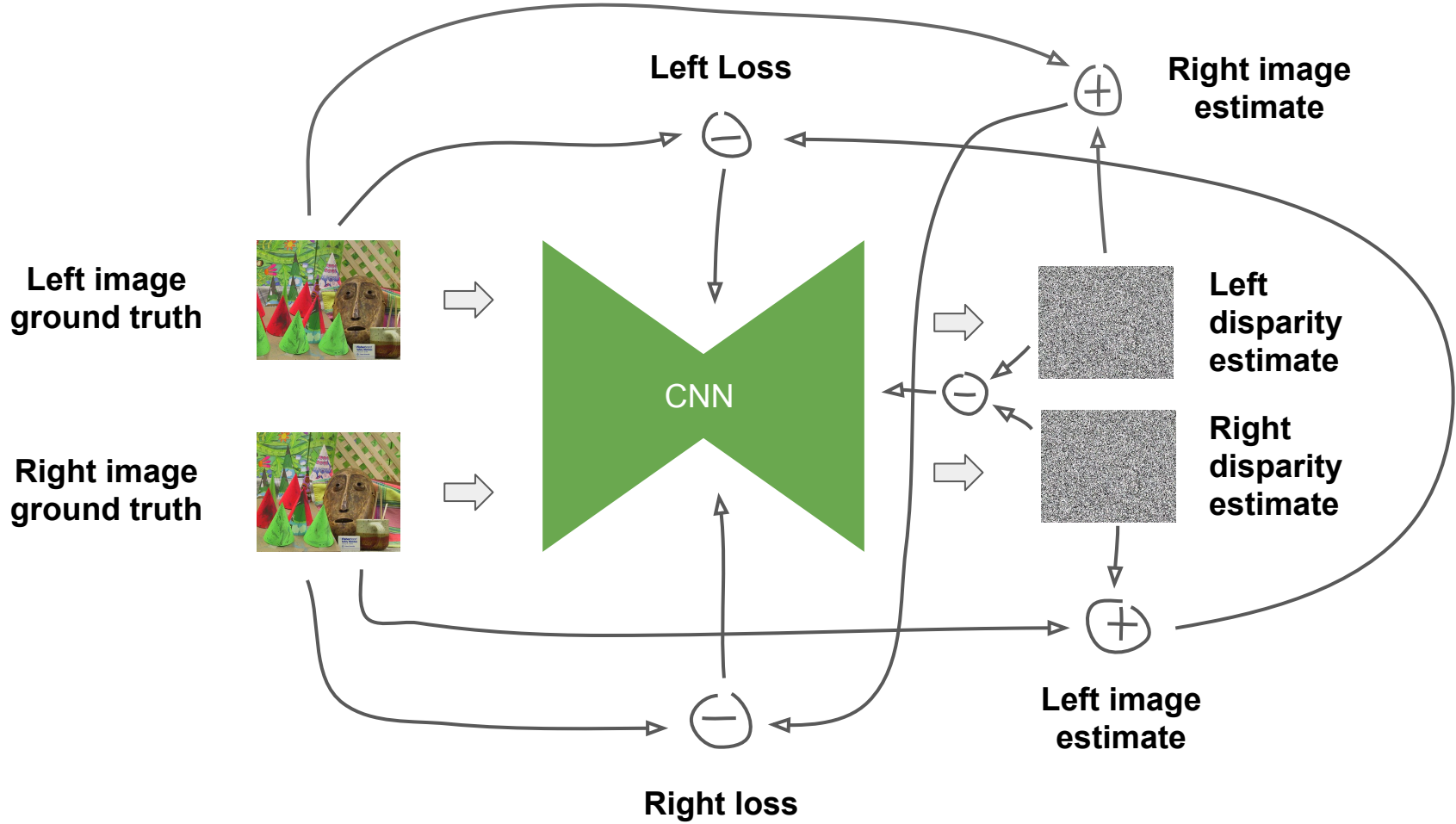


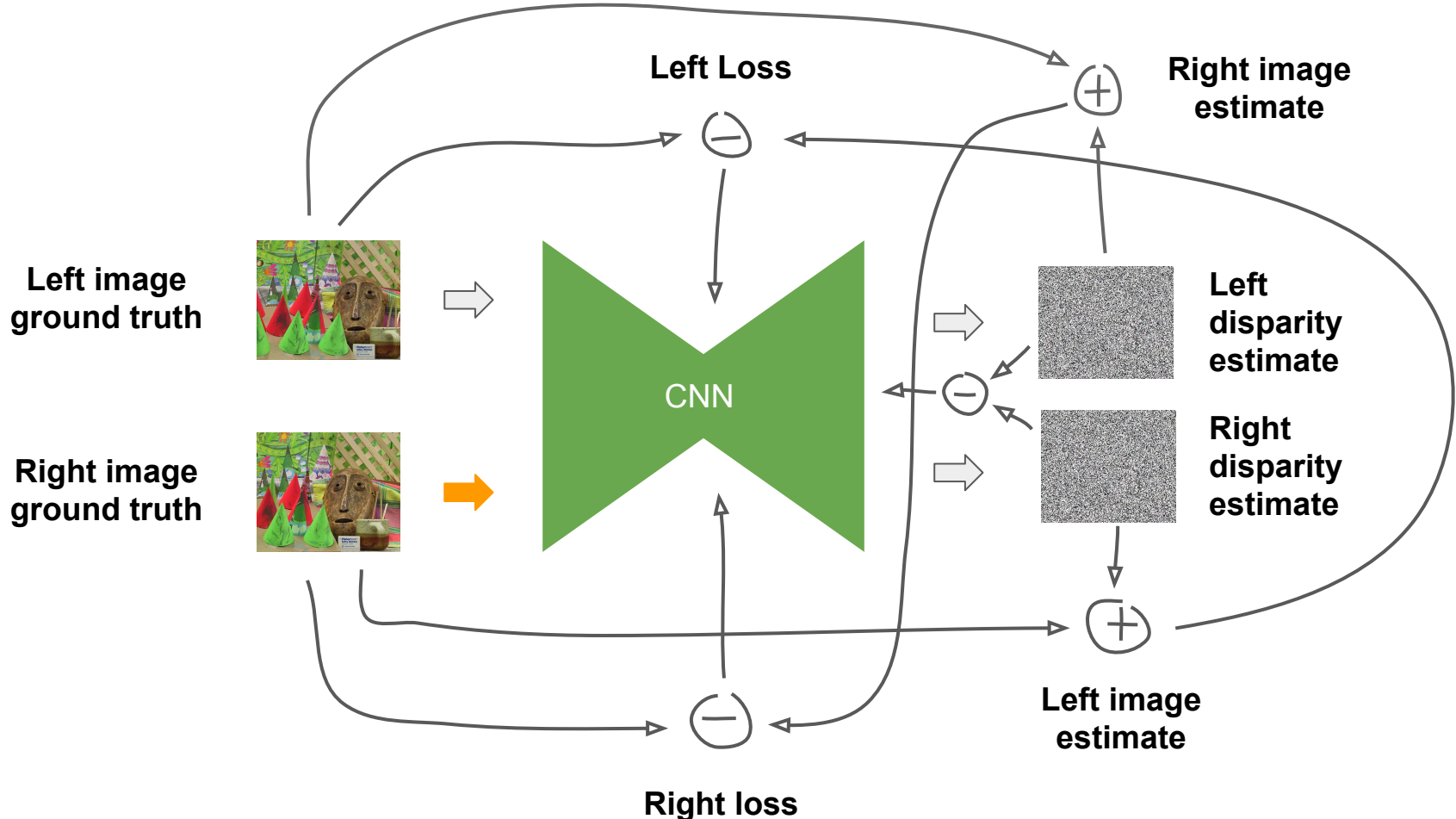












Left image ground truth



Right image ground truth



Left Loss



CNN

Right image estimate



Left disparity estimate



Right disparity estimate

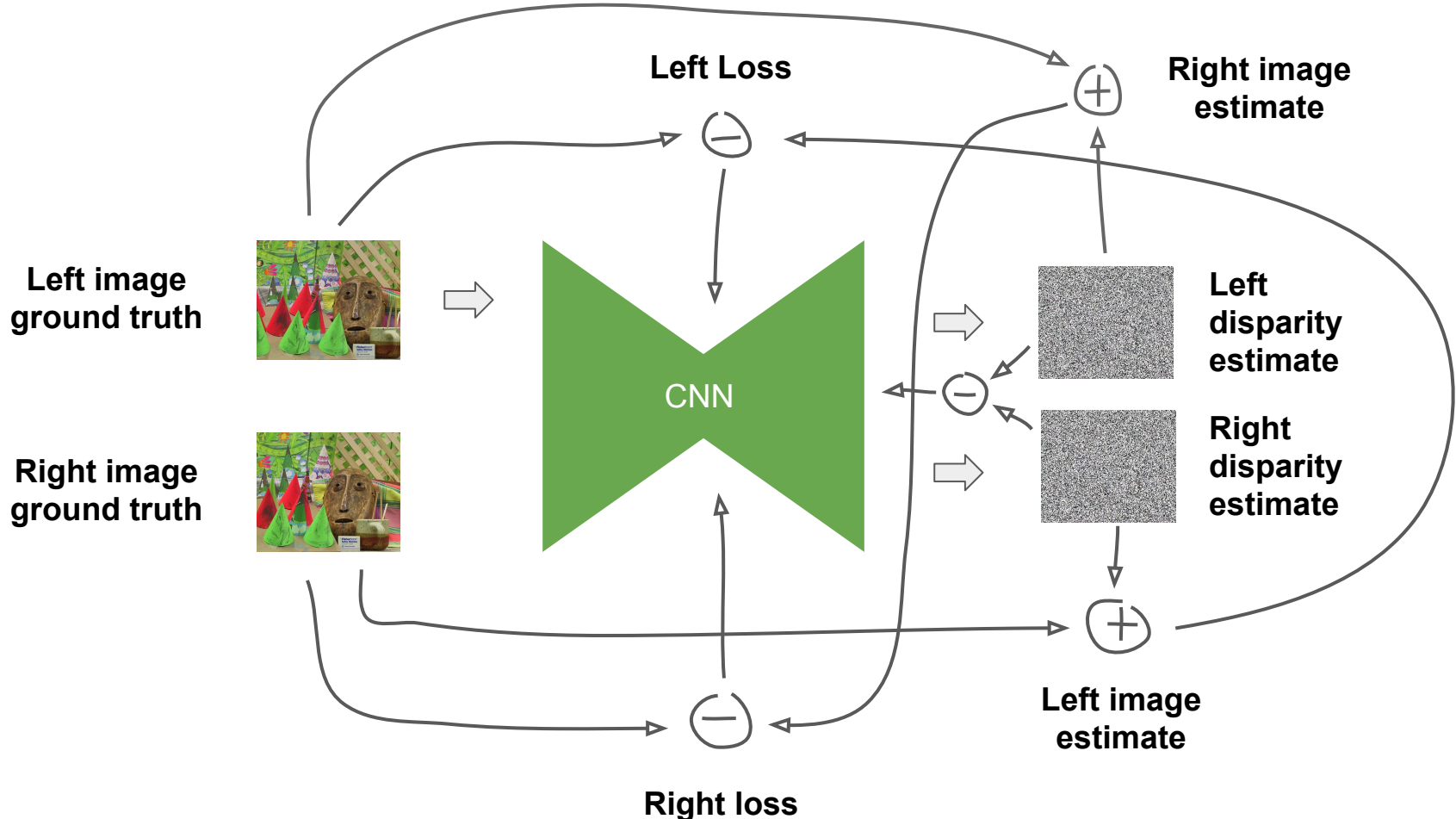


Left image estimate

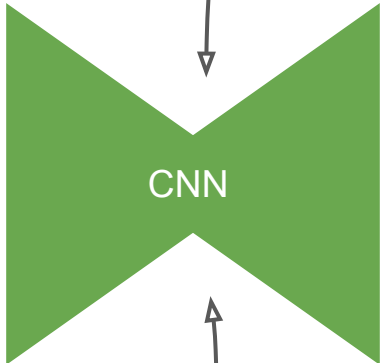


Right loss





Left image ground truth



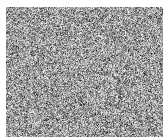
Left Loss



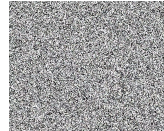
Right image estimate



Right image ground truth



Left disparity estimate



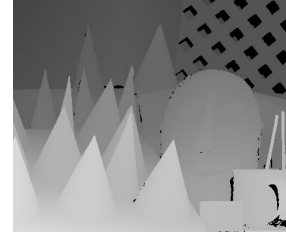
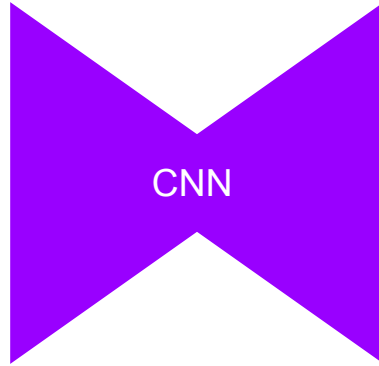
Right disparity estimate



Left image estimate

Right loss





Unsupervised **Monocular** Depth Estimation with Left-Right Consistency

Clément Godard

Oisín Mac Aodha

Gabriel J. Brostow

University College London

<http://visual.cs.ucl.ac.uk/pubs/monoDepth/>

Abstract

Learning based methods have shown very promising results for the task of depth estimation in single images. However, most existing approaches treat depth prediction as a supervised regression problem and as a result, require vast quantities of corresponding ground truth depth data for training. Just recording quality depth data in a range of environments is a challenging problem. In this paper, we innovate beyond existing approaches, replacing the use of explicit depth data during training with easier-to-obtain binocular stereo footage.

We propose a novel training objective that enables our convolutional neural network to learn to perform single image depth estimation, despite the absence of ground truth depth data. Exploiting epipolar geometry constraints, we generate disparity images by training our network with an image reconstruction loss. We show that solving for image reconstruction alone results in poor quality depth images. To overcome this problem, we propose a novel training loss that enforces consistency between the disparities produced relative to both the left and right

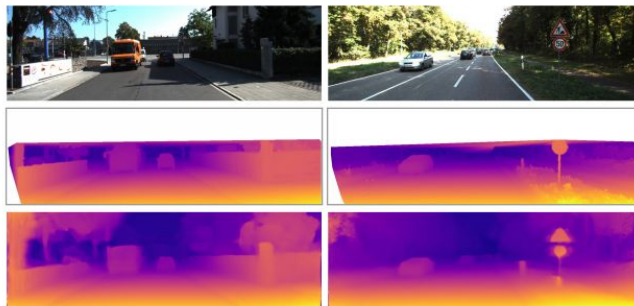


Figure 1. Our depth prediction results on KITTI 2015. Top to bottom: input image, ground truth disparities, and our result. Our method is able to estimate depth for thin structures such as street signs and poles.

have been restricted to scenes where large image collections and their corresponding pixel depths are available.

Understanding the shape of a scene from a single image, independent of its appearance, is a fundamental problem in machine perception. There are many applications such as synthetic object insertion in computer graphics [29], synthetic

Differentiable Bilinear Sampling

Next time