

Depth and Surface Normal Estimation from a Single Image

Mian Wei
University of Toronto

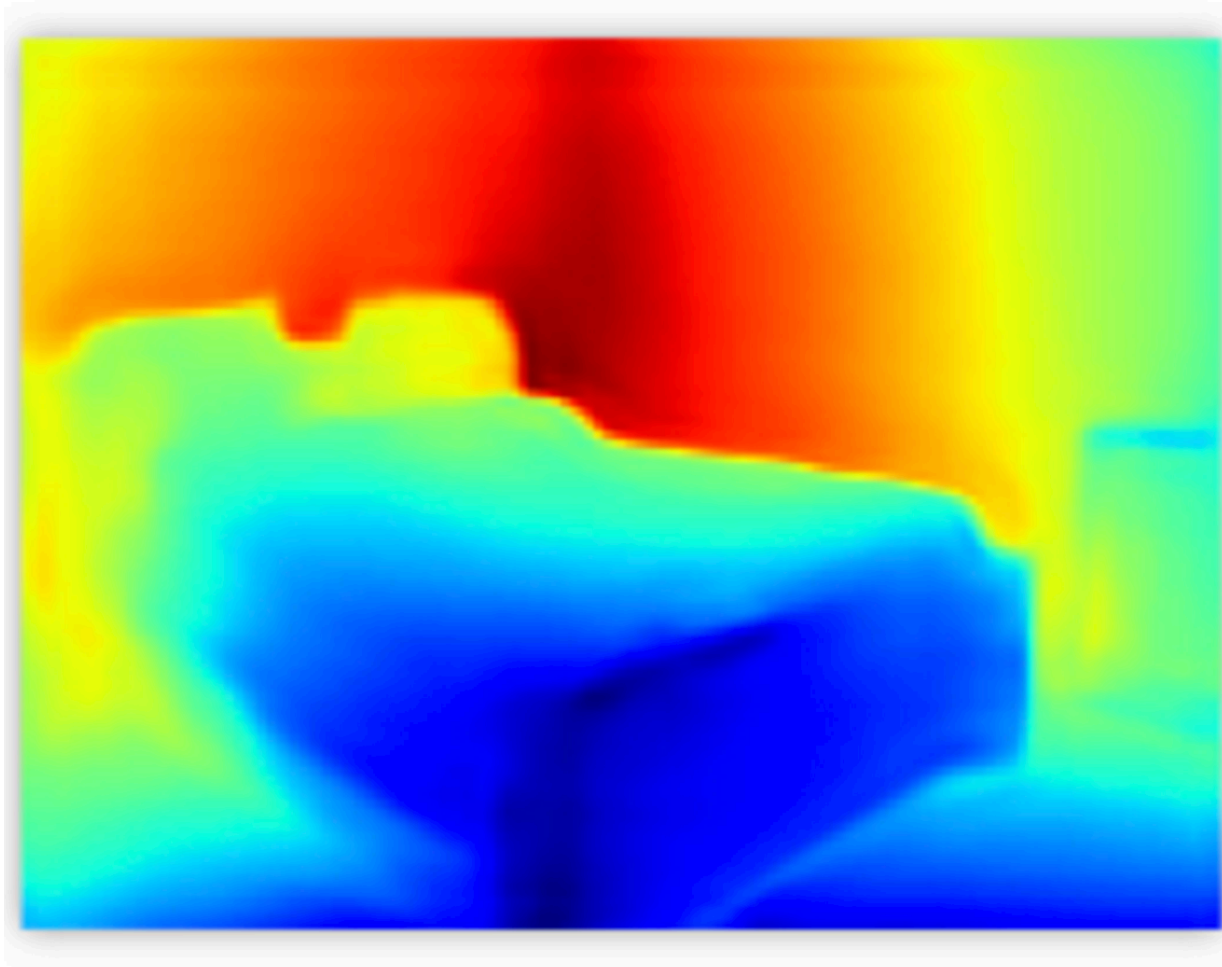
What is the problem?

Given one image

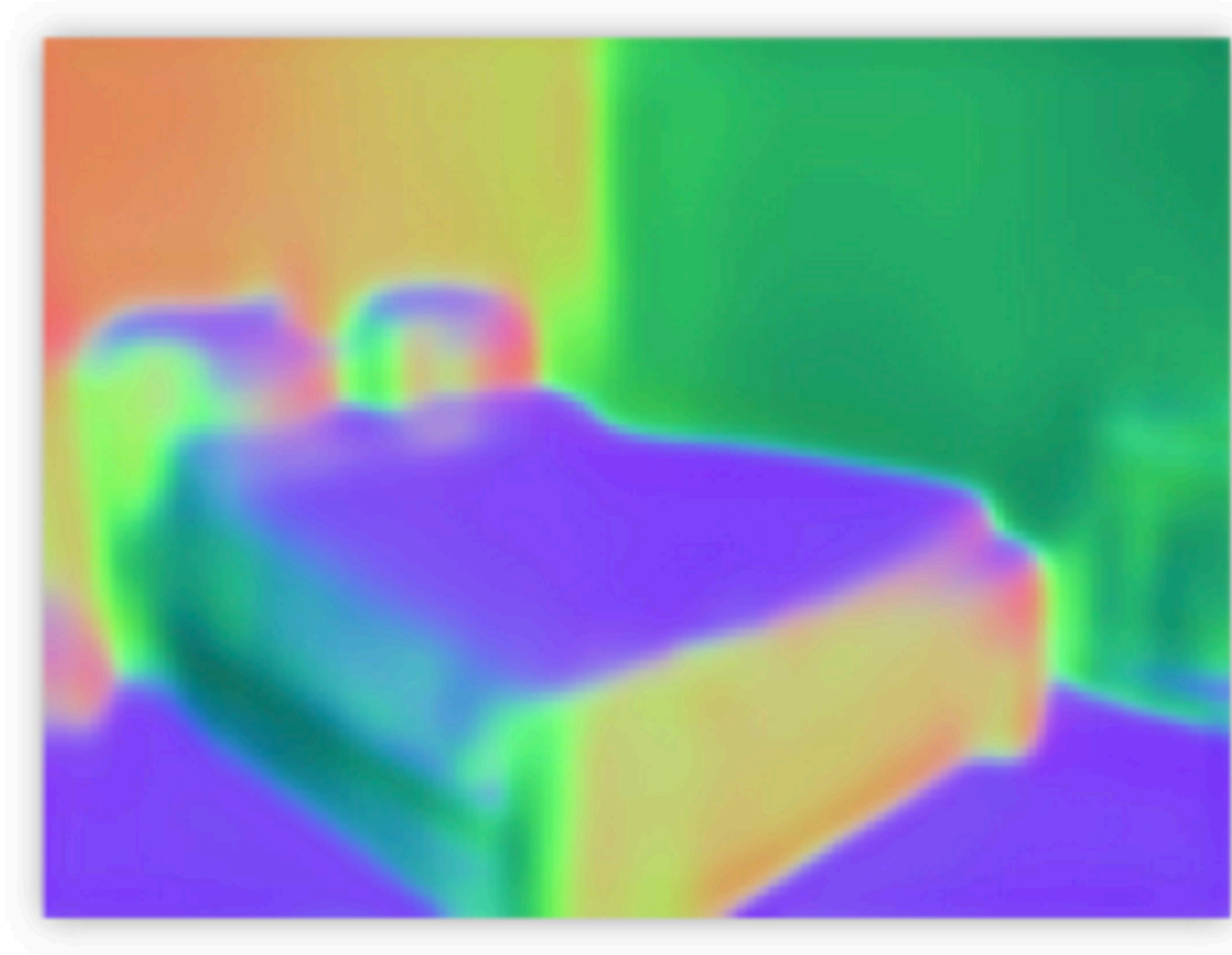


N. Silberman, D. Hoiem, P. Kohli, and R. Fergus, "Indoor segmentation and support inference from RGBD images," in *Proc. Eur. Conf. Comput. Vision*, 2012, pp. 746–760.

Estimate the following:



Eigen, D. and Fergus, R. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. ICCV 2015

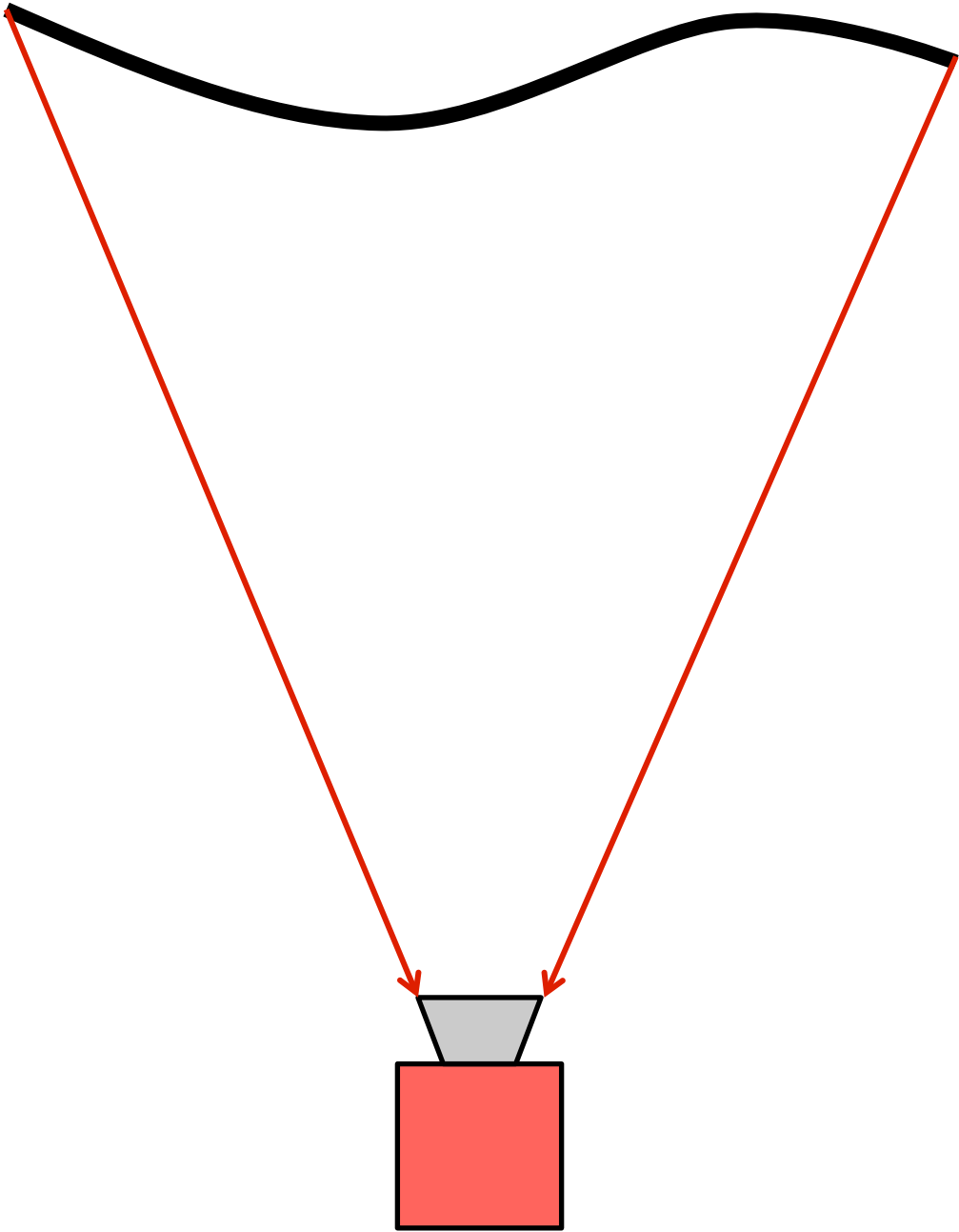


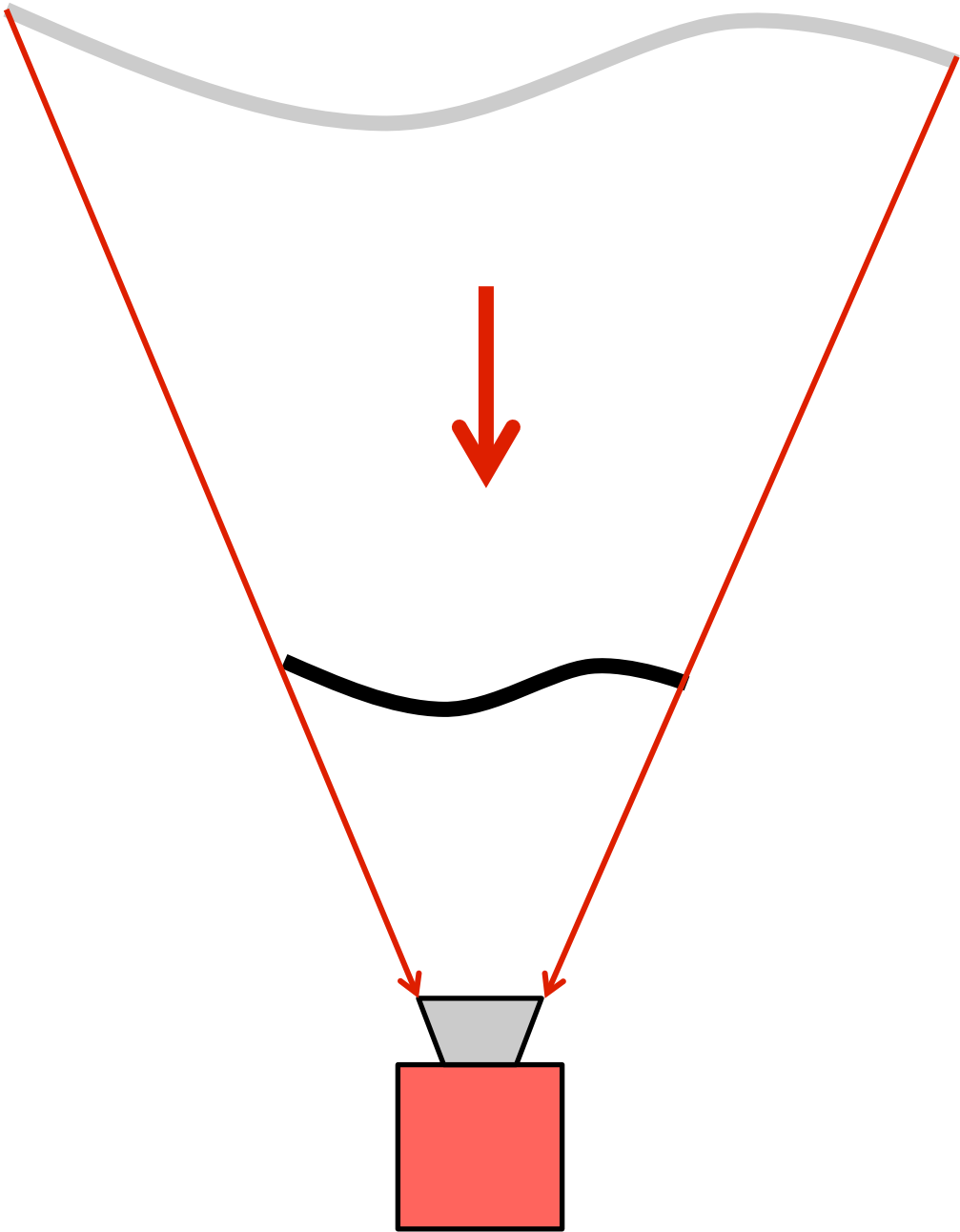
Eigen, D. and Fergus, R. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. ICCV 2015

Why is this hard?

Multiple ambiguities

Scale ambiguity

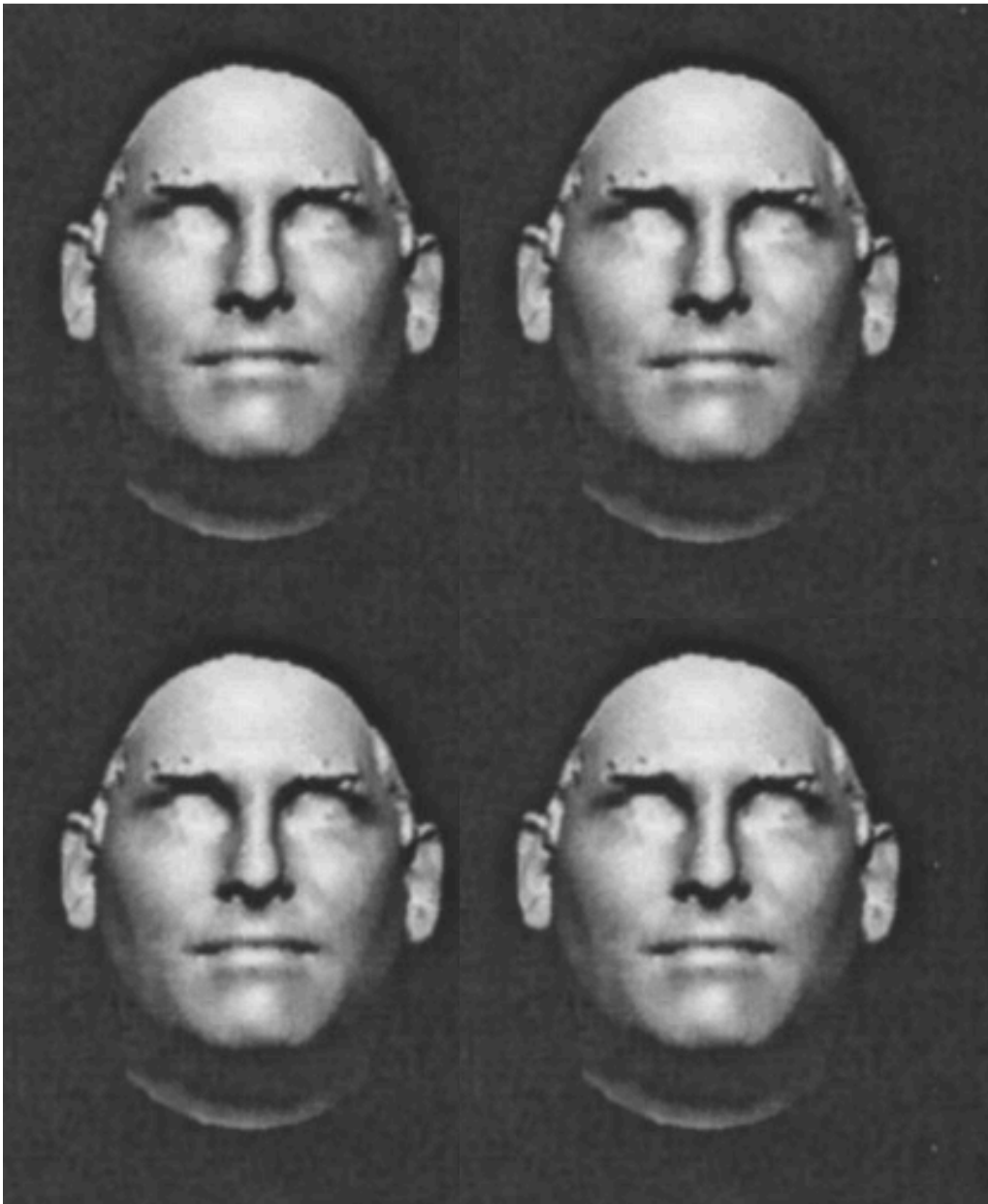




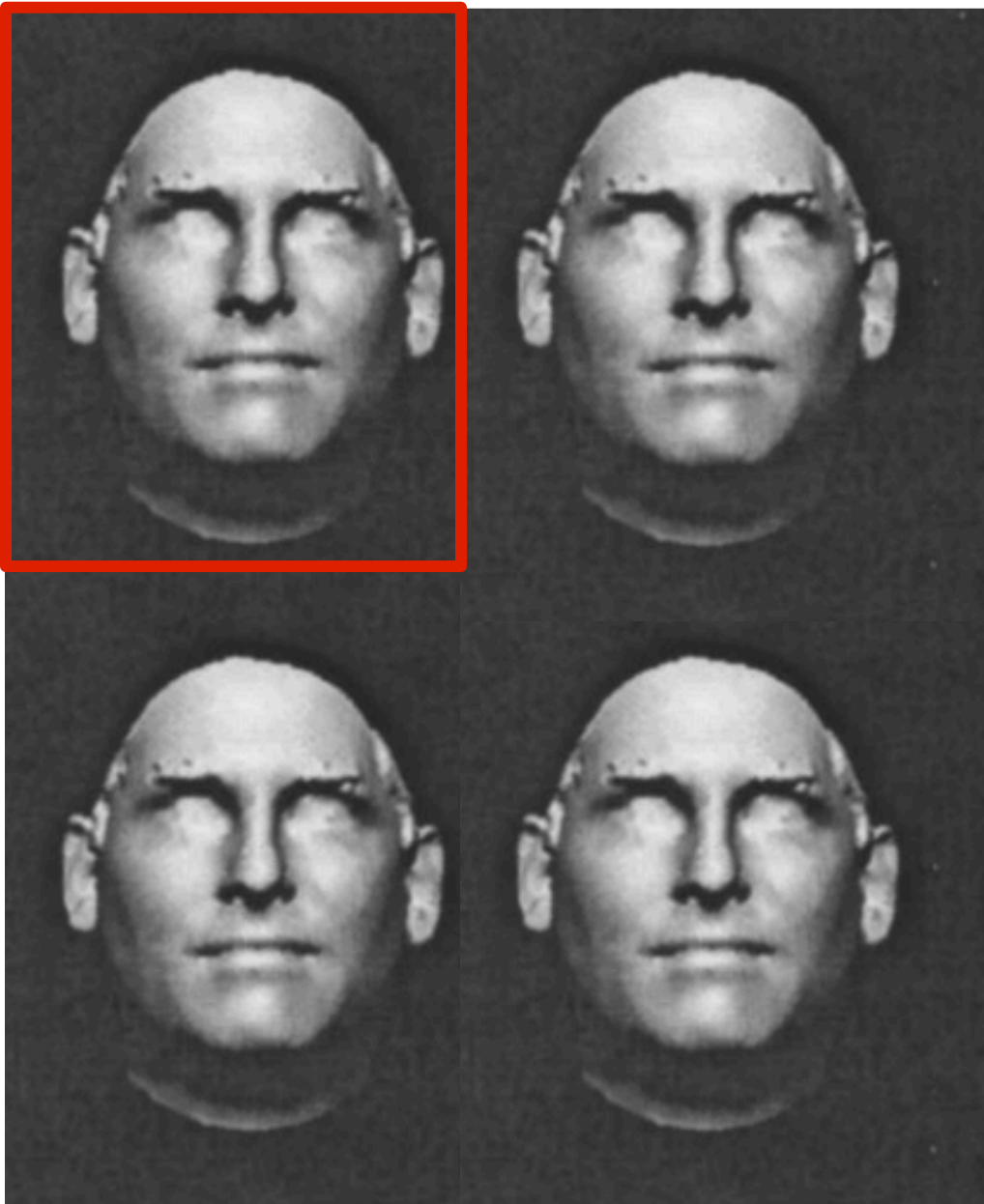
Bas-relief ambiguity

Let's play a game

Spot the Difference



P. Belhumeur, D. Kriegman, and A. Yuille, "The Bas-Relief Ambiguity," Proc. IEEE Conf. Computer Vision and Pattern Recognition, pp. 1040-1046, 1997.

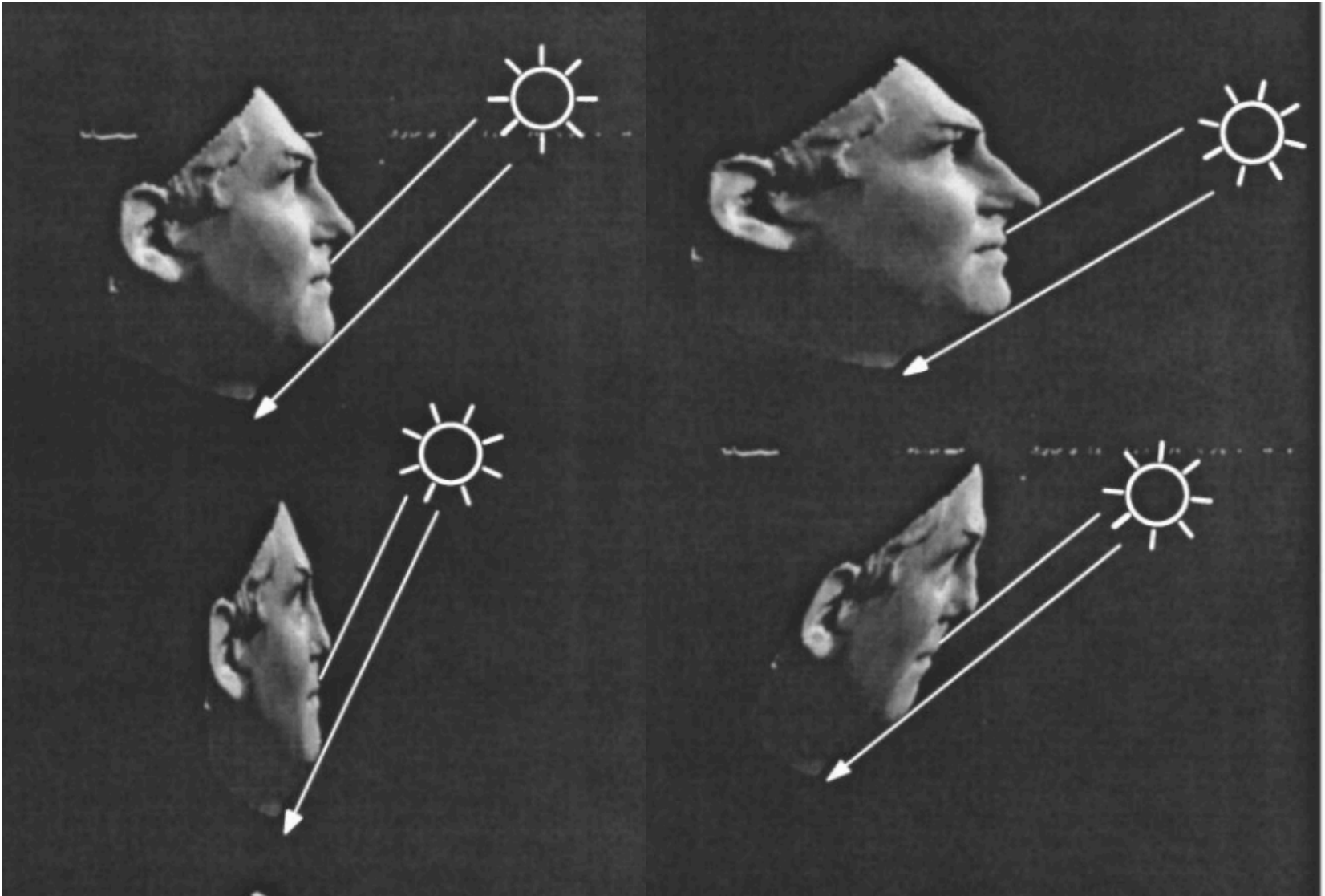


P. Belhumeur, D. Kriegman, and A. Yuille, "The Bas-Relief Ambiguity," Proc. IEEE Conf. Computer Vision and Pattern Recognition, pp. 1040-1046, 1997.

All the same



P. Belhumeur, D. Kriegman, and A. Yuille, "The Bas-Relief Ambiguity," Proc. IEEE Conf. Computer Vision and Pattern Recognition, pp. 1040-1046, 1997.



P. Belhumeur, D. Kriegman, and A. Yuille, "The Bas-Relief Ambiguity," Proc. IEEE Conf. Computer Vision and Pattern Recognition, pp. 1040-1046, 1997.

Family of transformation

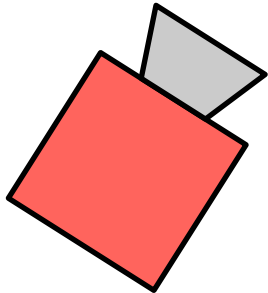
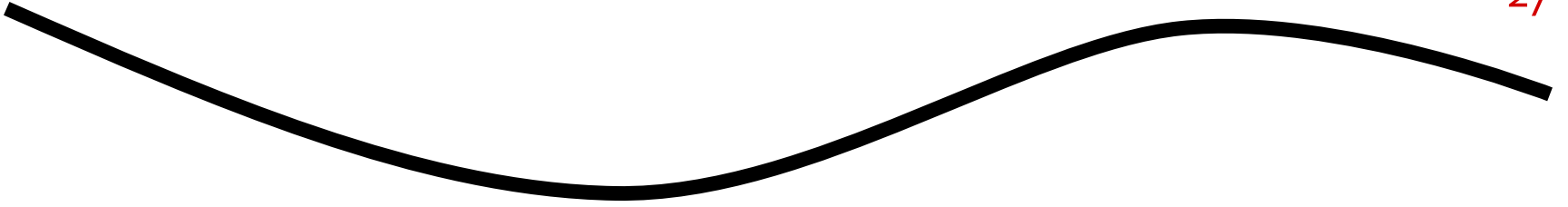
Generalized Bas-Relief

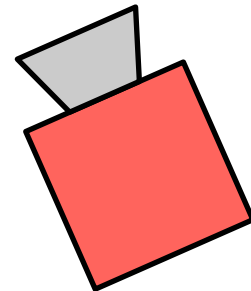
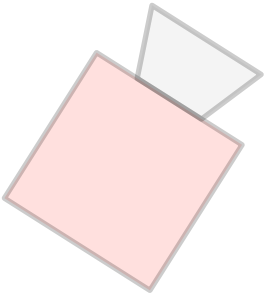
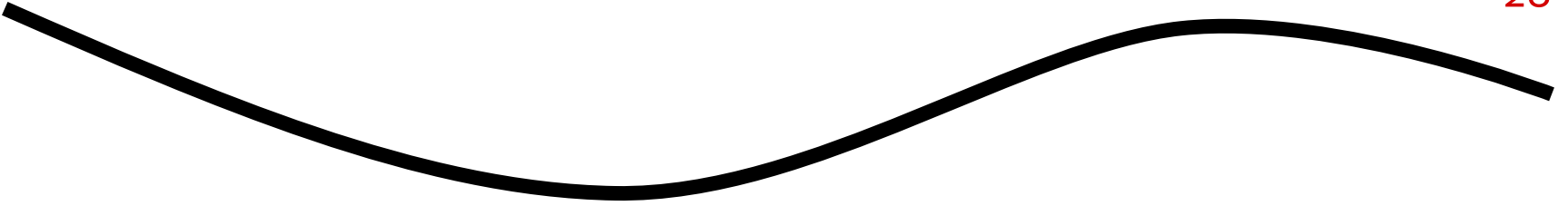
Change shape and illumination

Yield same image

Existing works

Multi-view Stereo

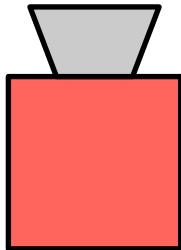
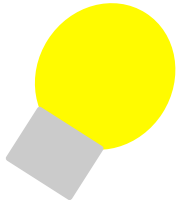
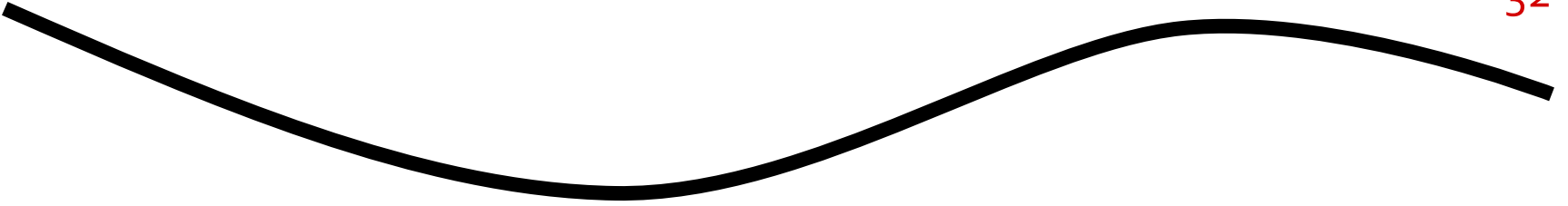


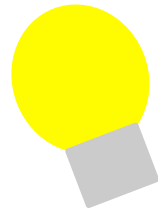
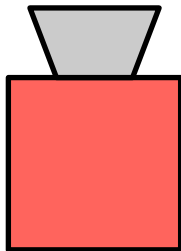
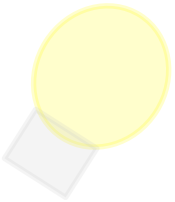
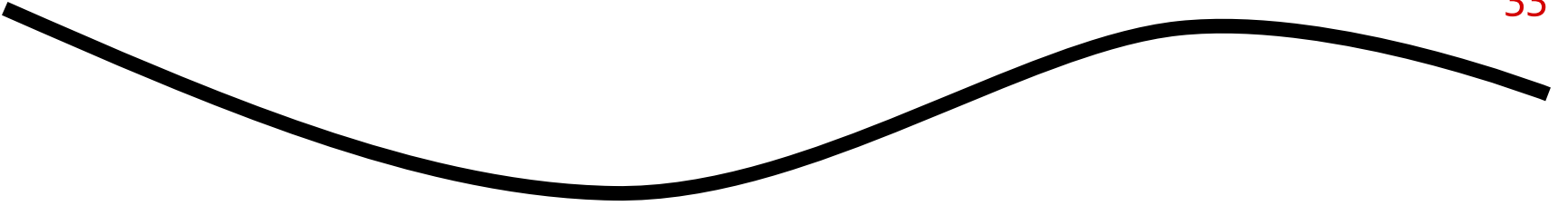


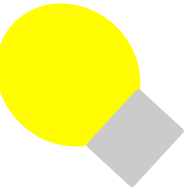
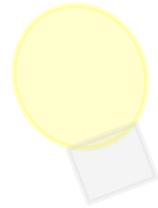
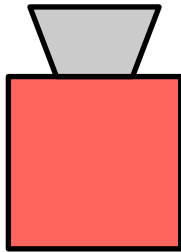
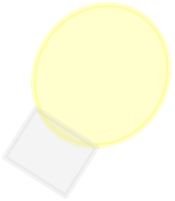
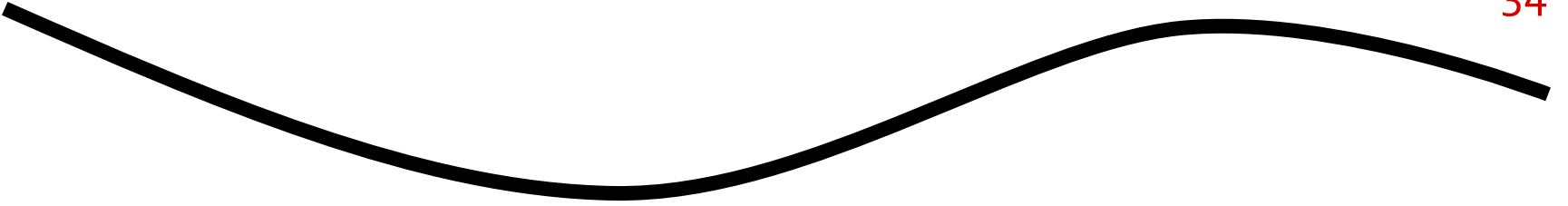
Photometric Stereo

Collimated Light Sources

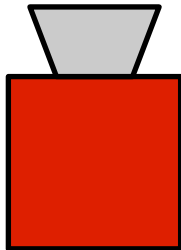
Light rays parallel







Shape from Focus

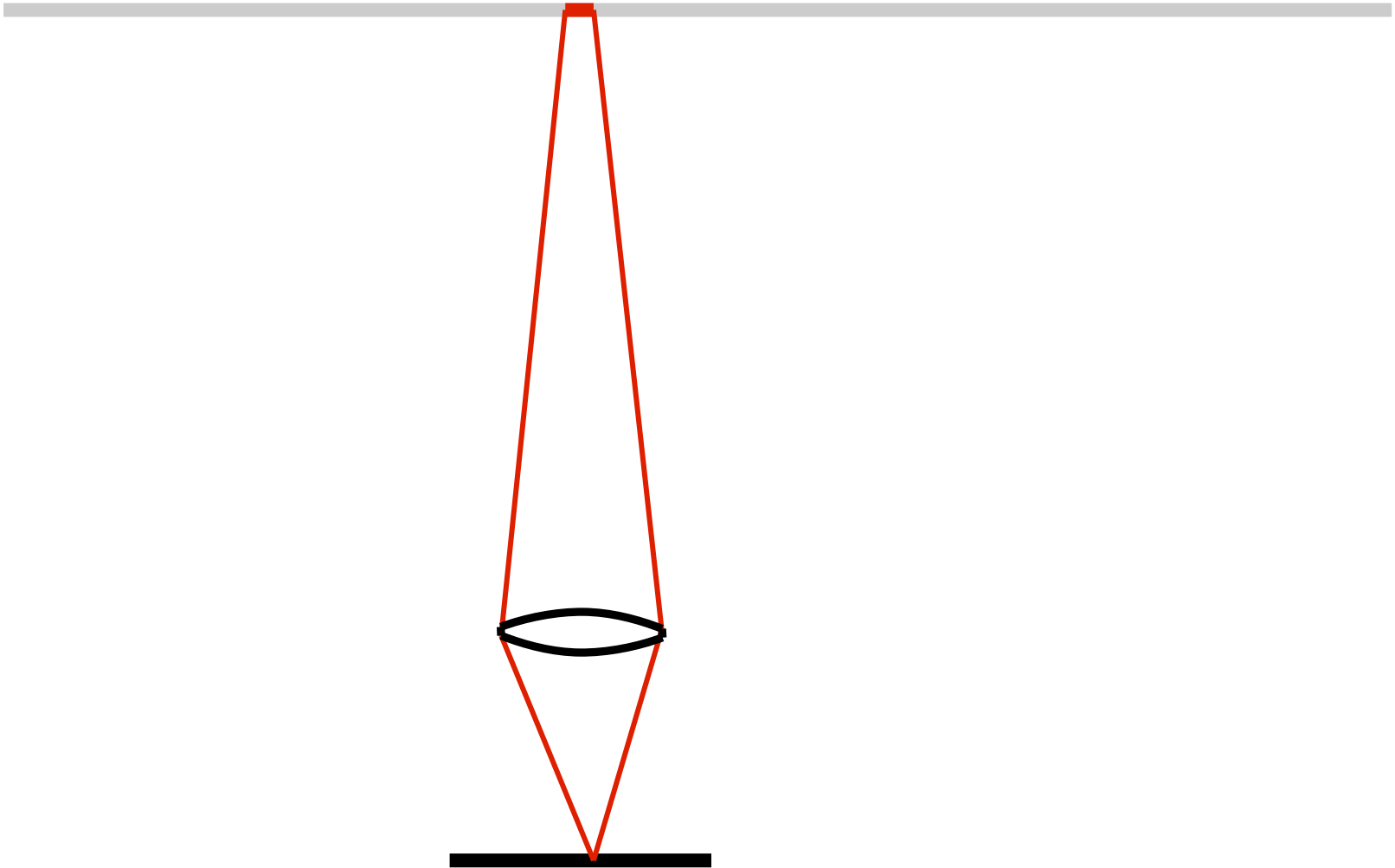


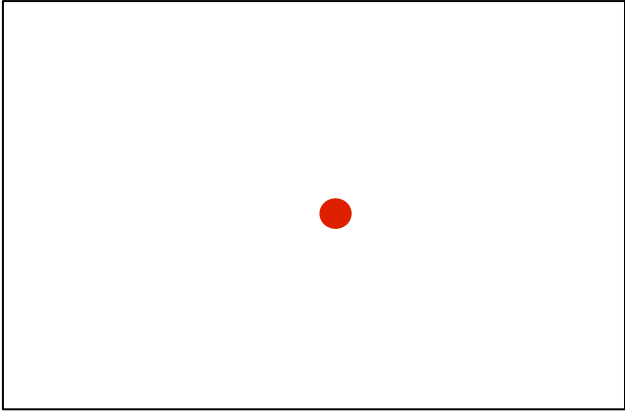
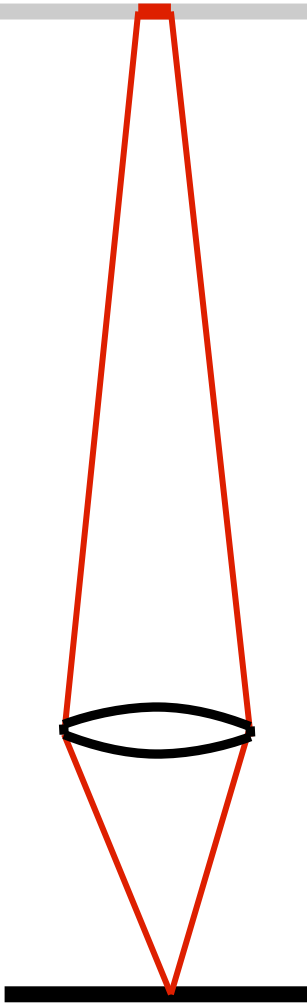


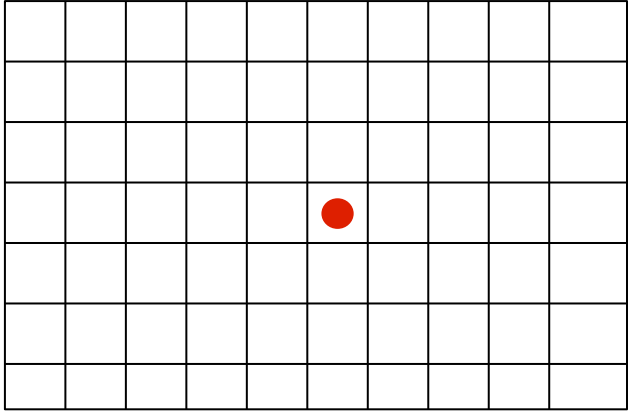
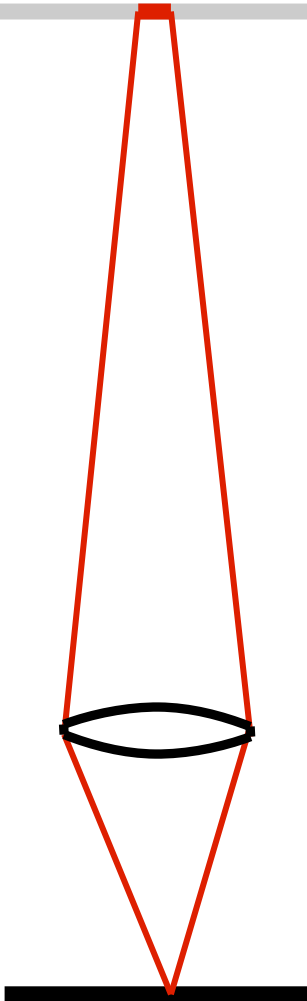


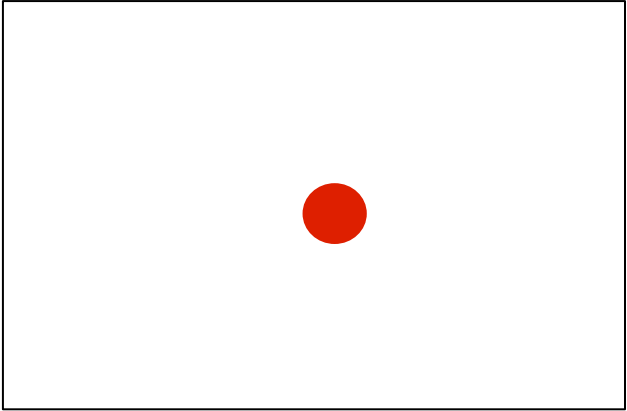
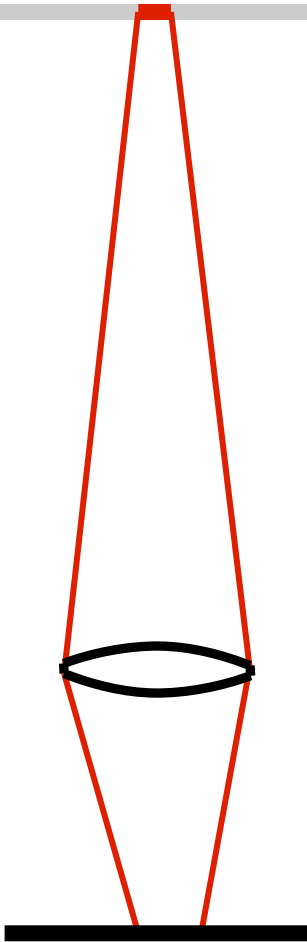


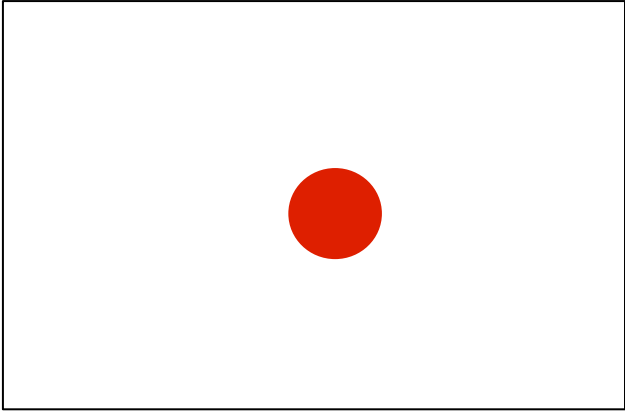
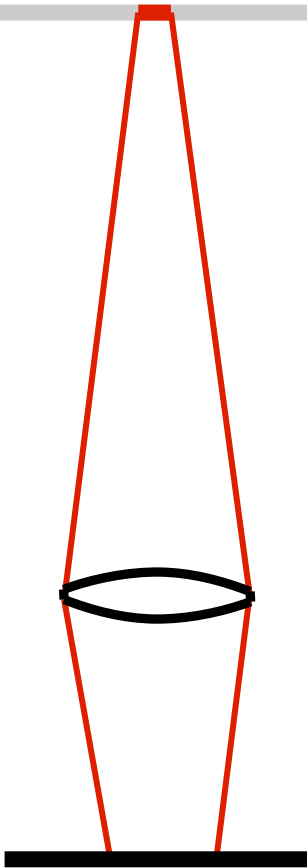


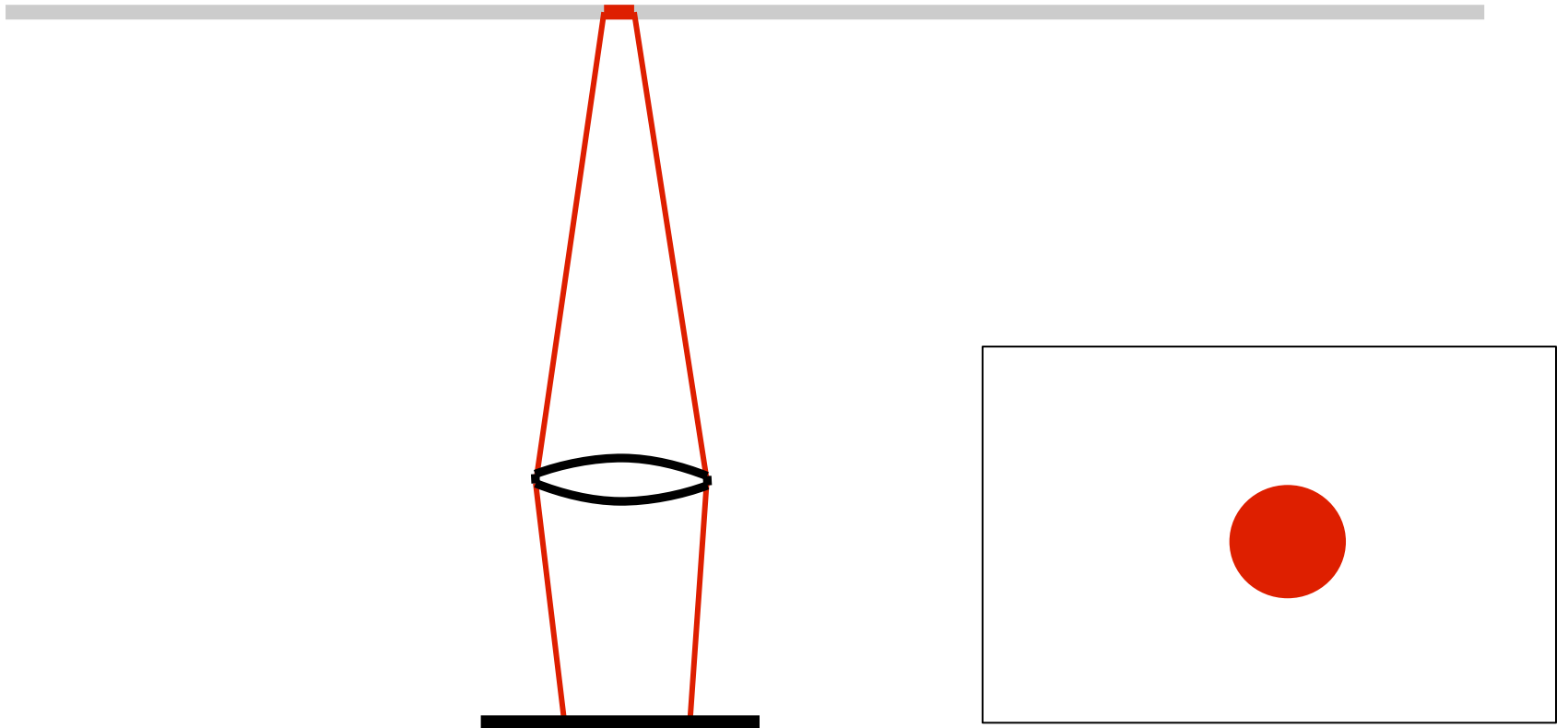


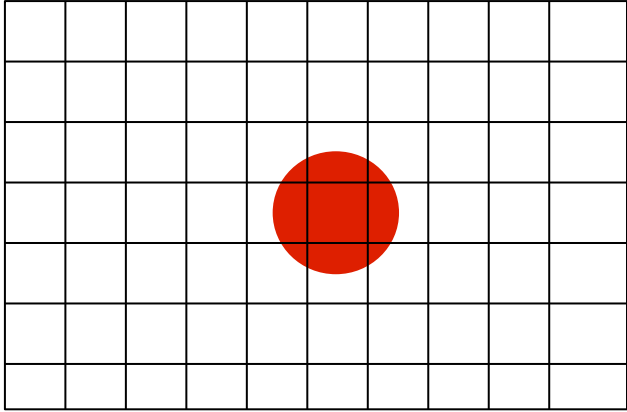
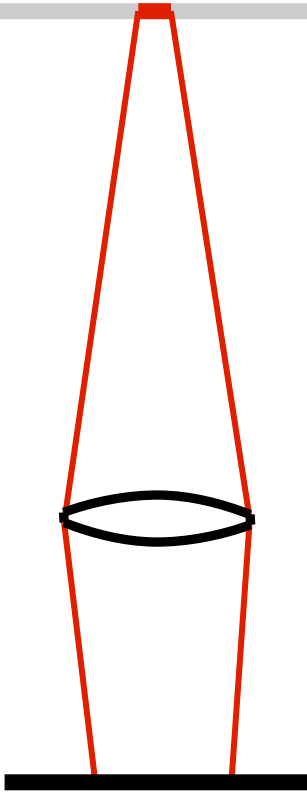




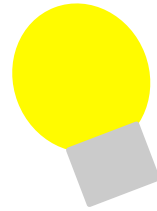
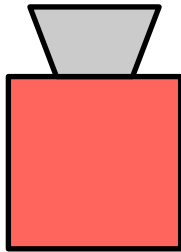
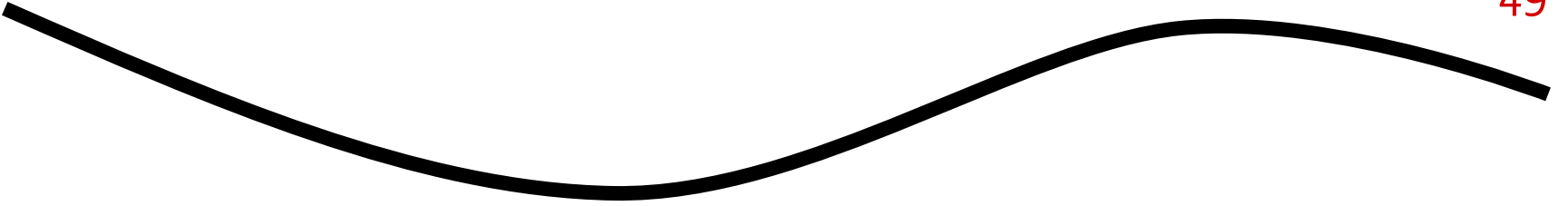


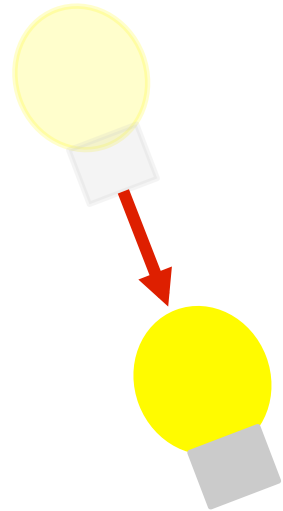
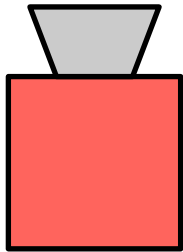
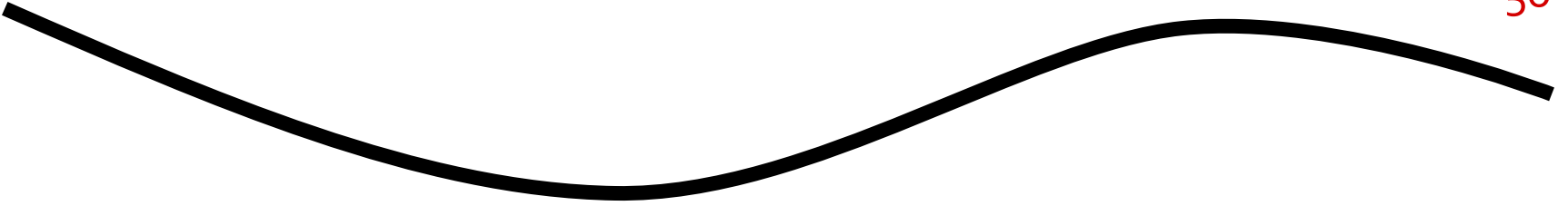






Light Fall-off Stereo





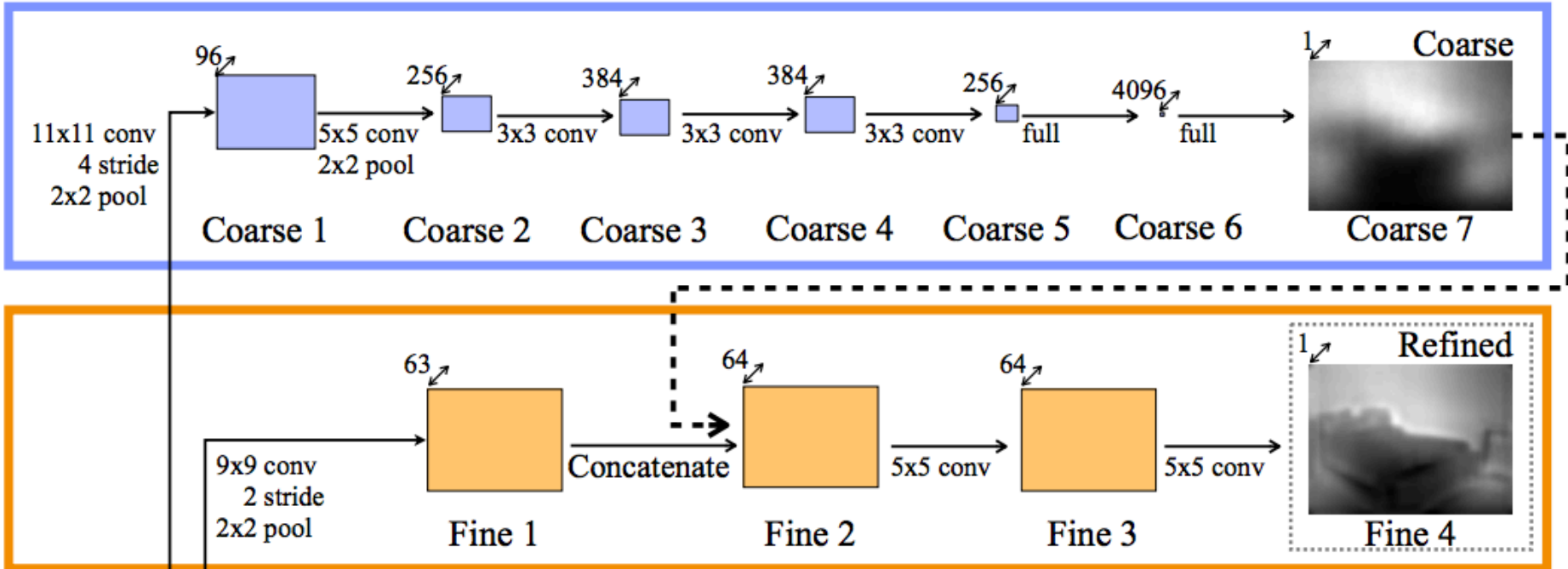
Specialized Hardware

Laser Scanner

Active Illumination

Time of Flight

Estimating Depth



Train 2 networks

Global coarse-scale network

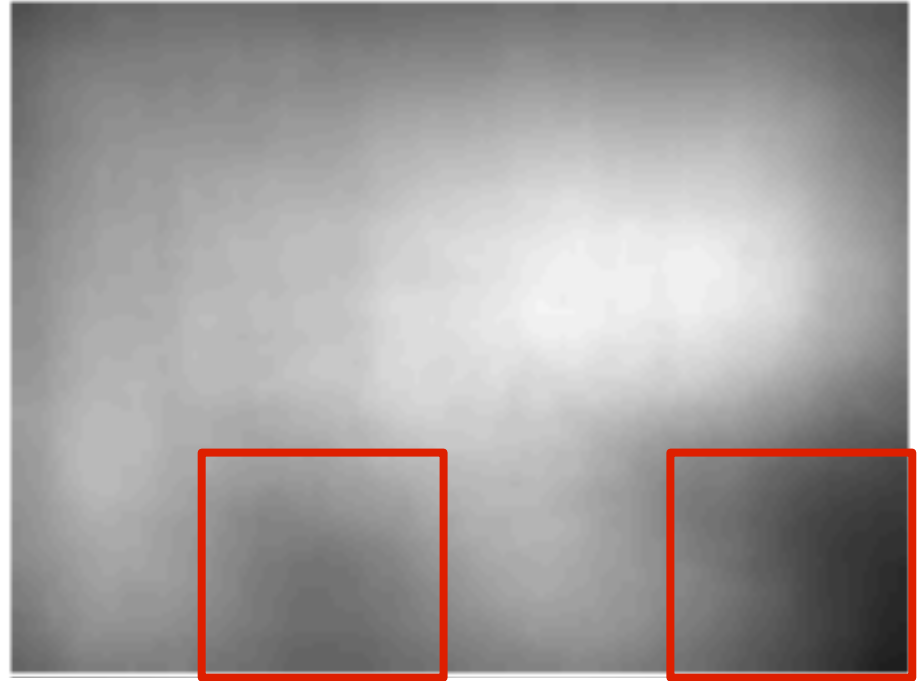
Local fine-scale network

Global coarse-scale network

Learns a coarse depth map







Used as input to local network

Intuition:

Coarse info learnt already

Focus on learning finer info





Scale ambiguity

Scale invariant error function

$$D(y, y^*) = \frac{1}{2n} \sum_{i=1}^n (\log y_i - \log y_i^* + \alpha(y_i, y_i^*))^2$$

$$\alpha(y_i, y_i^*) = \frac{1}{n} \sum_{i=1}^n (\log y_i^* - \log y_i)$$

$$D(ay, ay^*) = \frac{1}{2n} \sum_{i=1}^n (\log ay_i - \log ay_i^* + \alpha(ay_i, ay_i^*))^2$$

$$D(ay, ay^*) = \frac{1}{2n} \sum_{i=1}^n (\log a - \log a + \log y_i - \log y_i^* + \alpha(ay_i, ay_i^*))^2$$

$$D(ay, ay^*) = \frac{1}{2n} \sum_{i=1}^n (\log y_i - \log y_i^* + \log a - \log a + \alpha(y_i, y_i^*))^2$$

$$D(ay, ay^*) = D(y, y^*)$$

Loss Function

Scale invariant

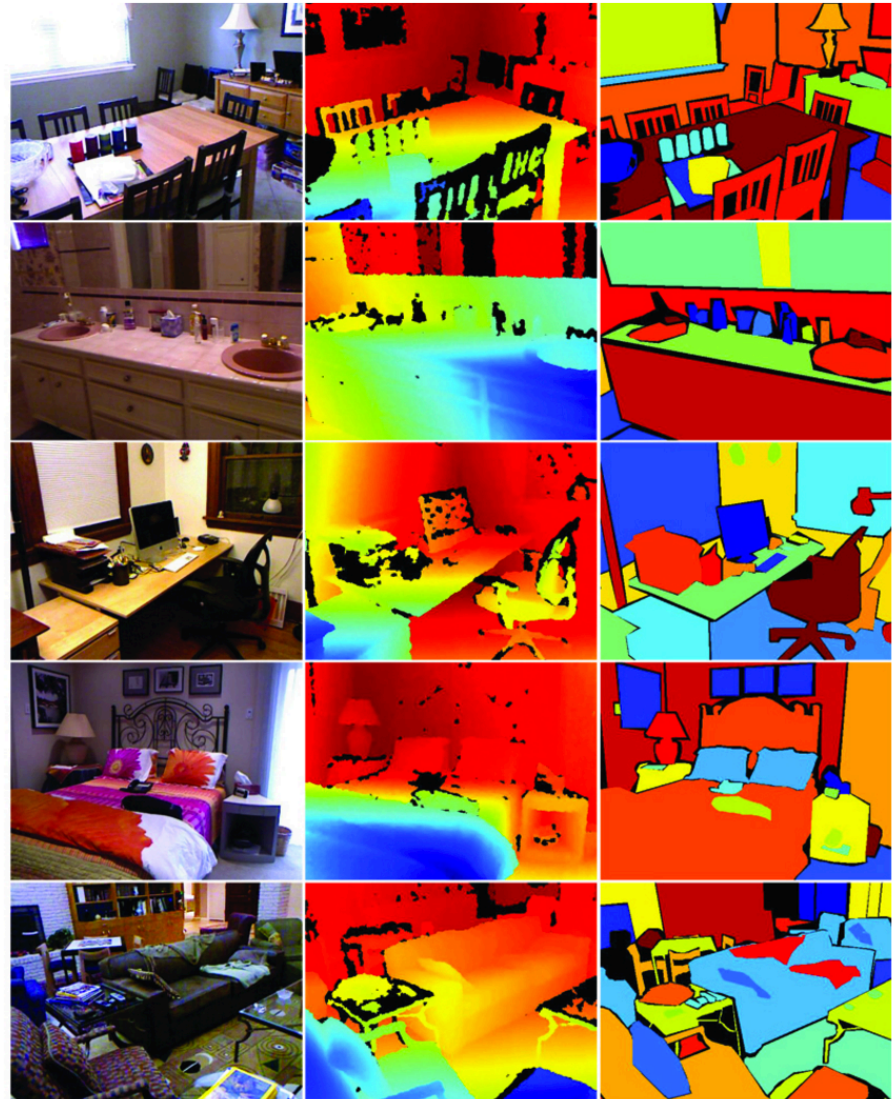
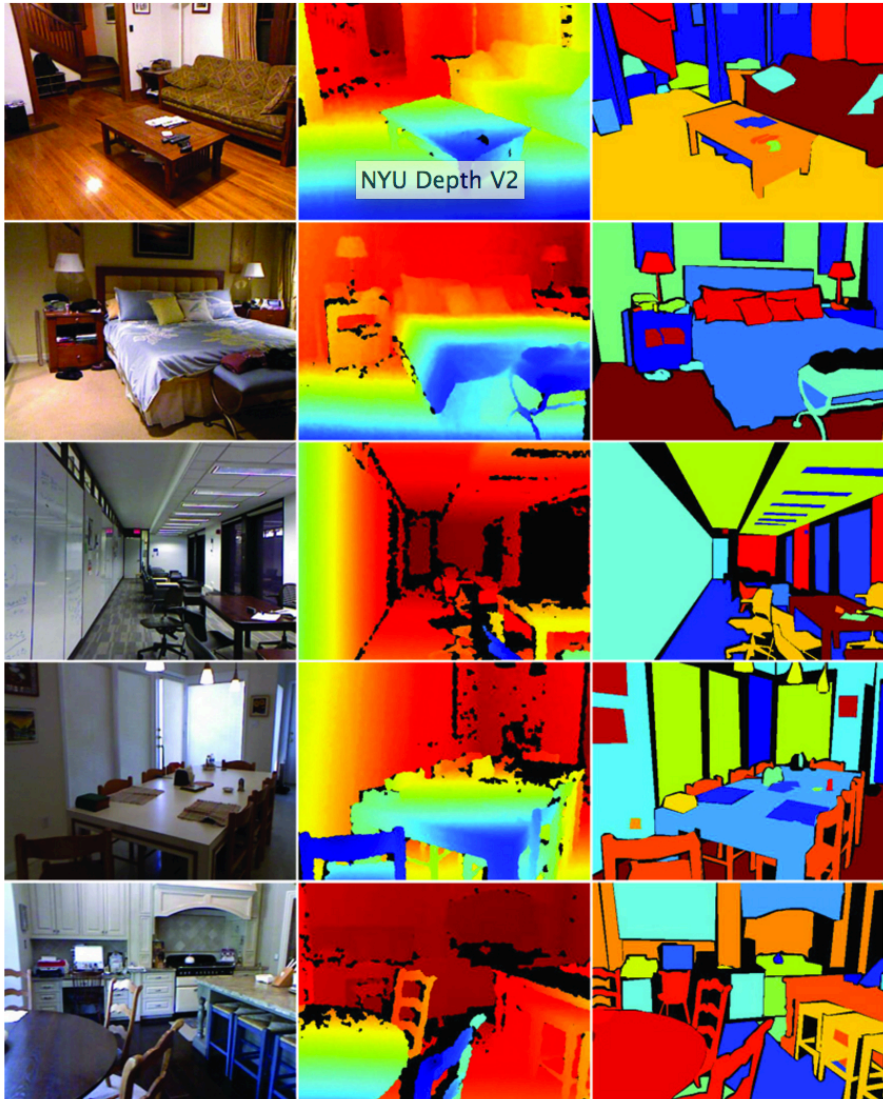
$$L(y, y^*) = \frac{1}{n} \sum_{i=1}^n d_i^2 - \frac{\lambda}{n^2} \left(\sum_{i=1}^n d_i \right)^2$$

$$d_i = \log y_i - \log y_i^*$$

2 Datasets

NYUDepthV2

Indoor Rooms



KITTI

Outdoor images taken on a car



How do you get ground truth?

NYU Depth V2

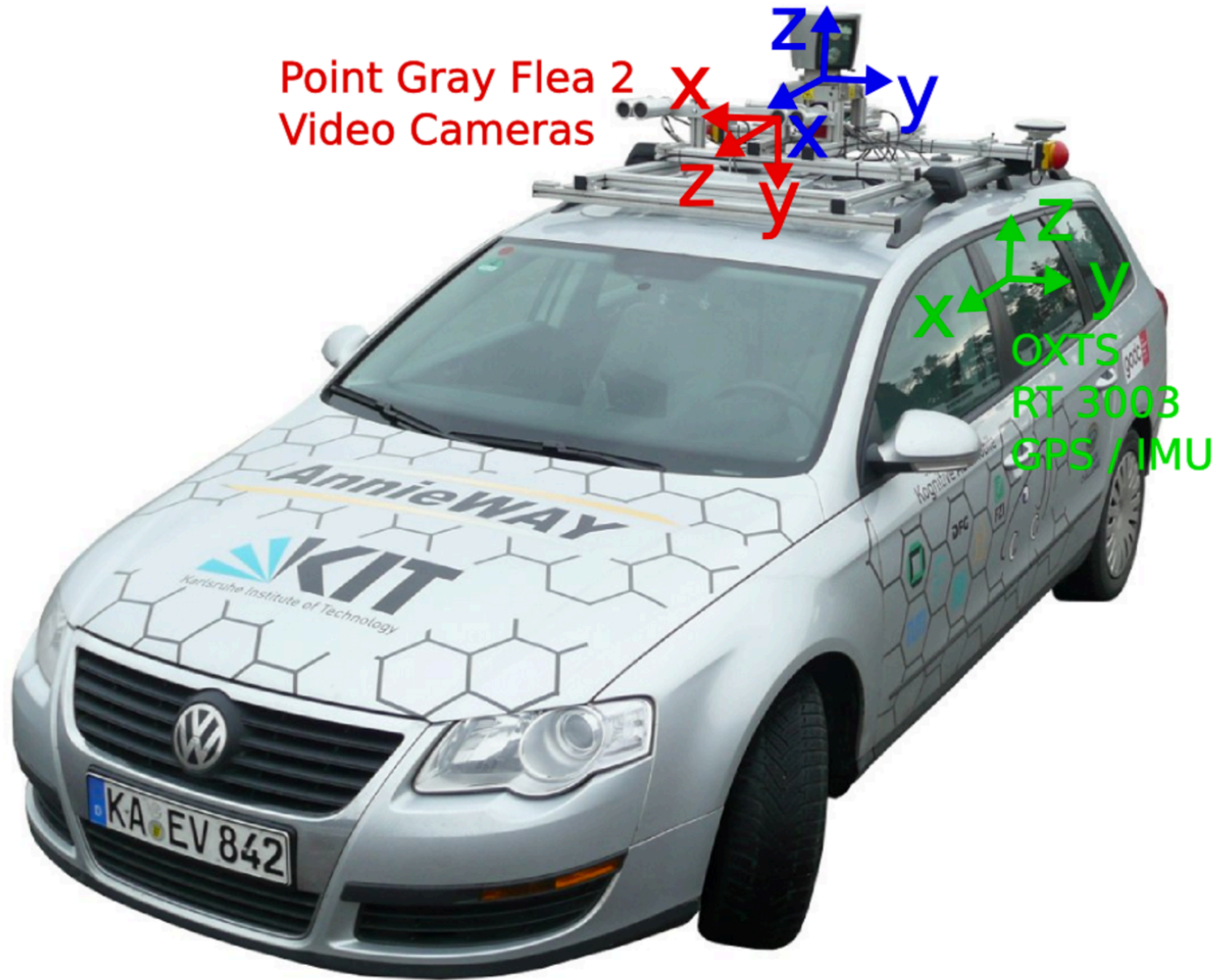
Kinect



KITTI

Velodyne HDL-64E Laserscanner

Point Gray Flea 2
Video Cameras



Time of Flight

Times how long light travels

From light source to camera

Results



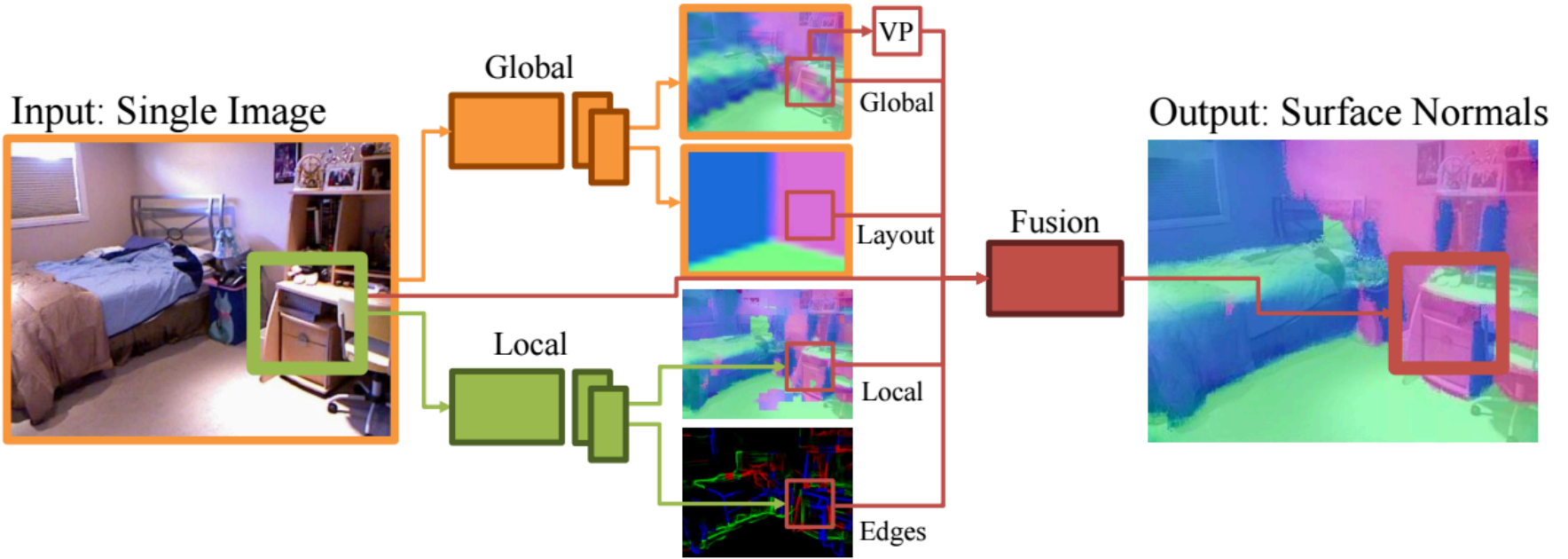


	Mean	Make3D	Ladicky&al	Karsch&al	Coarse	Coarse + Fine	
threshold $\delta < 1.25$	0.418	0.447	0.542	–	0.618	0.611	higher is better
threshold $\delta < 1.25^2$	0.711	0.745	0.829	–	0.891	0.887	
threshold $\delta < 1.25^3$	0.874	0.897	0.940	–	0.969	0.971	
abs relative difference	0.408	0.349	–	0.350	0.228	0.215	lower is better
sqr relative difference	0.581	0.492	–	–	0.223	0.212	
RMSE (linear)	1.244	1.214	–	1.2	0.871	0.907	
RMSE (log)	0.430	0.409	–	–	0.283	0.285	
RMSE (log, scale inv.)	0.304	0.325	–	–	0.221	0.219	

Table 1: Comparison on the NYUDepth dataset

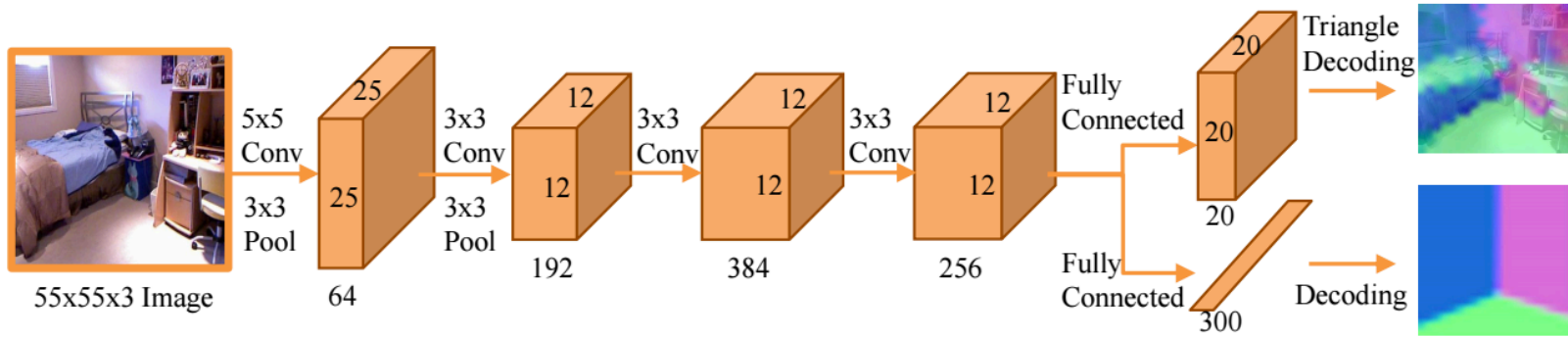
Estimating Surface Normals

Similar to Eigen

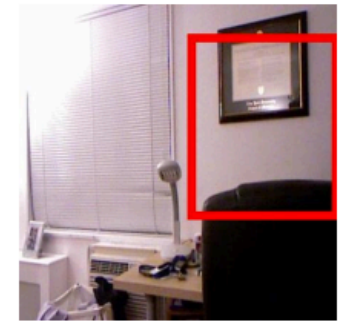
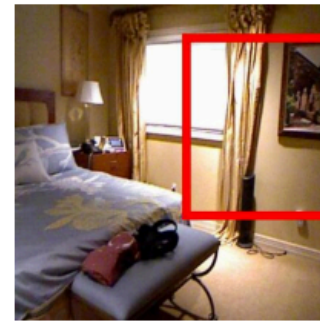
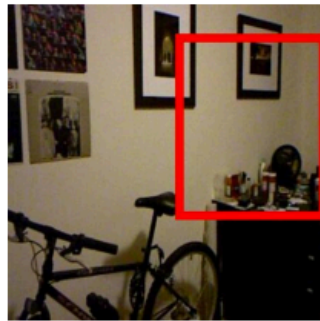
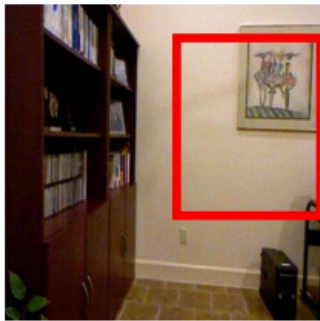
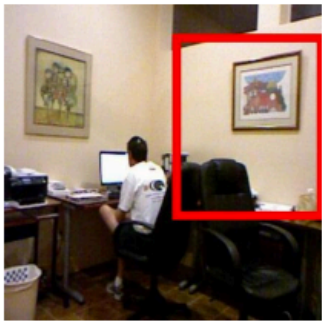
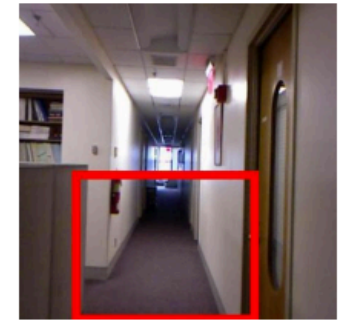
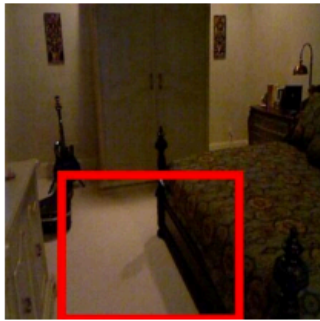
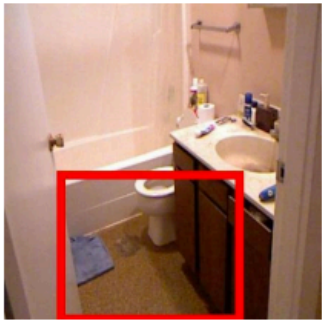
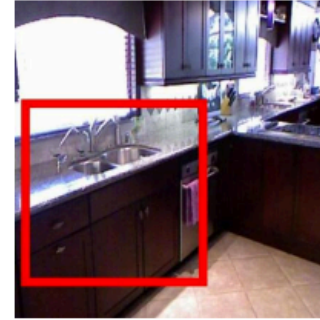
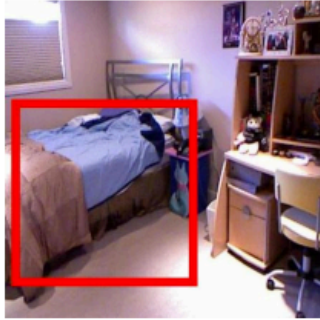
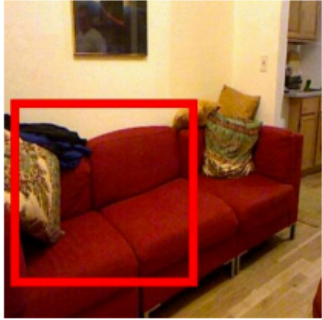


Trains 3 networks

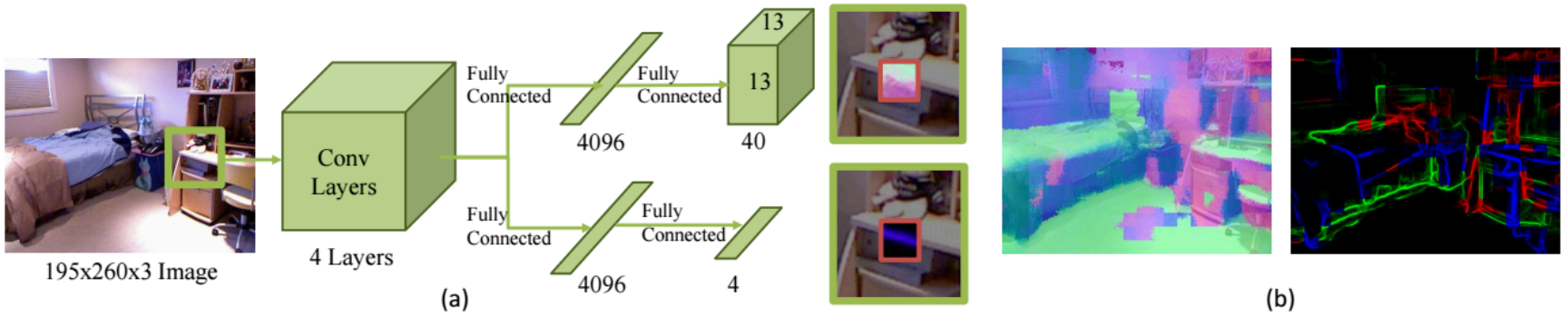
Global coarse-scale network



Trains for room layout as well

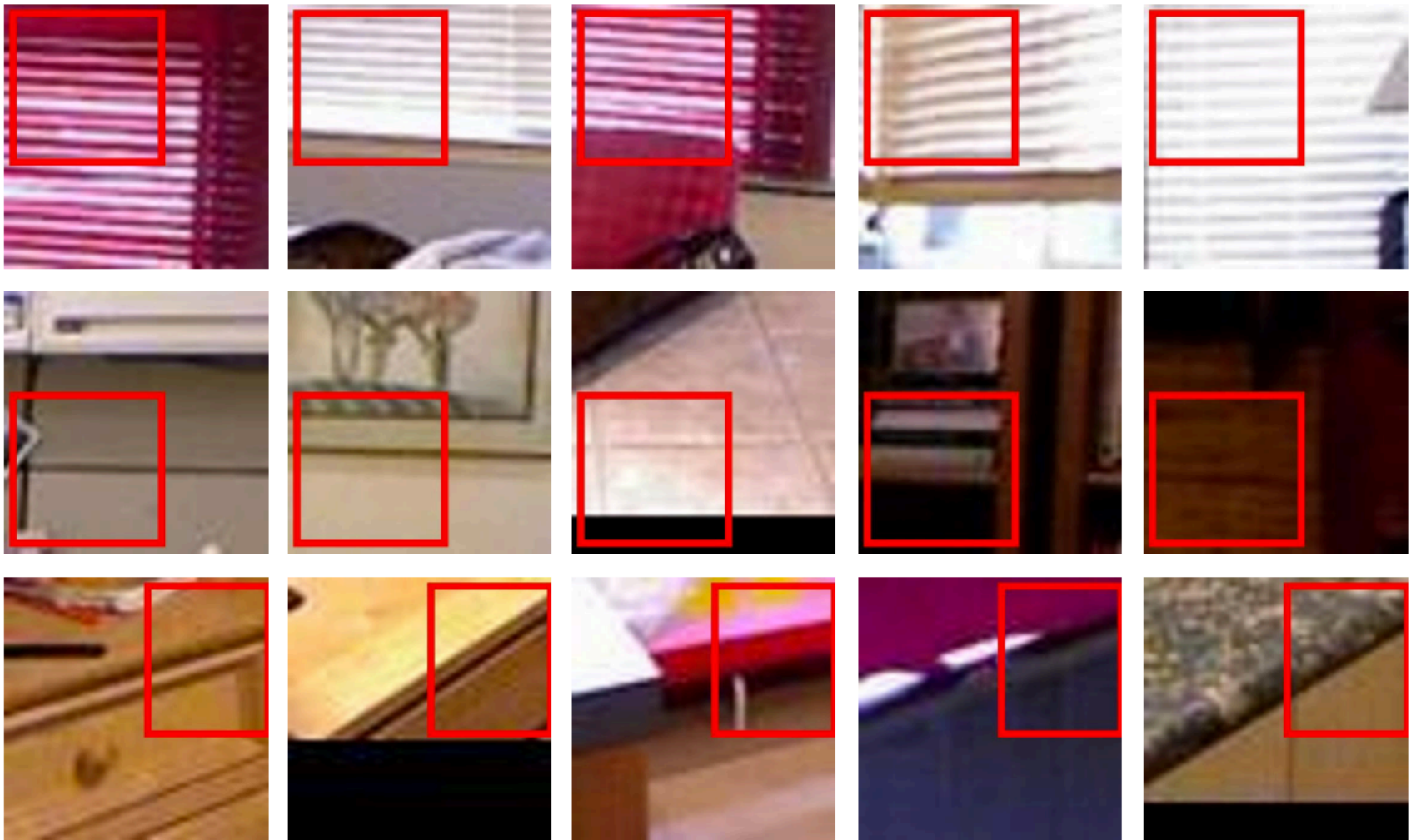


Local fine-scale network



Trains for edge labels as well

Convex, concave, occlusion, N/A



Difference: Global and Local
trained separately

Fusion Network

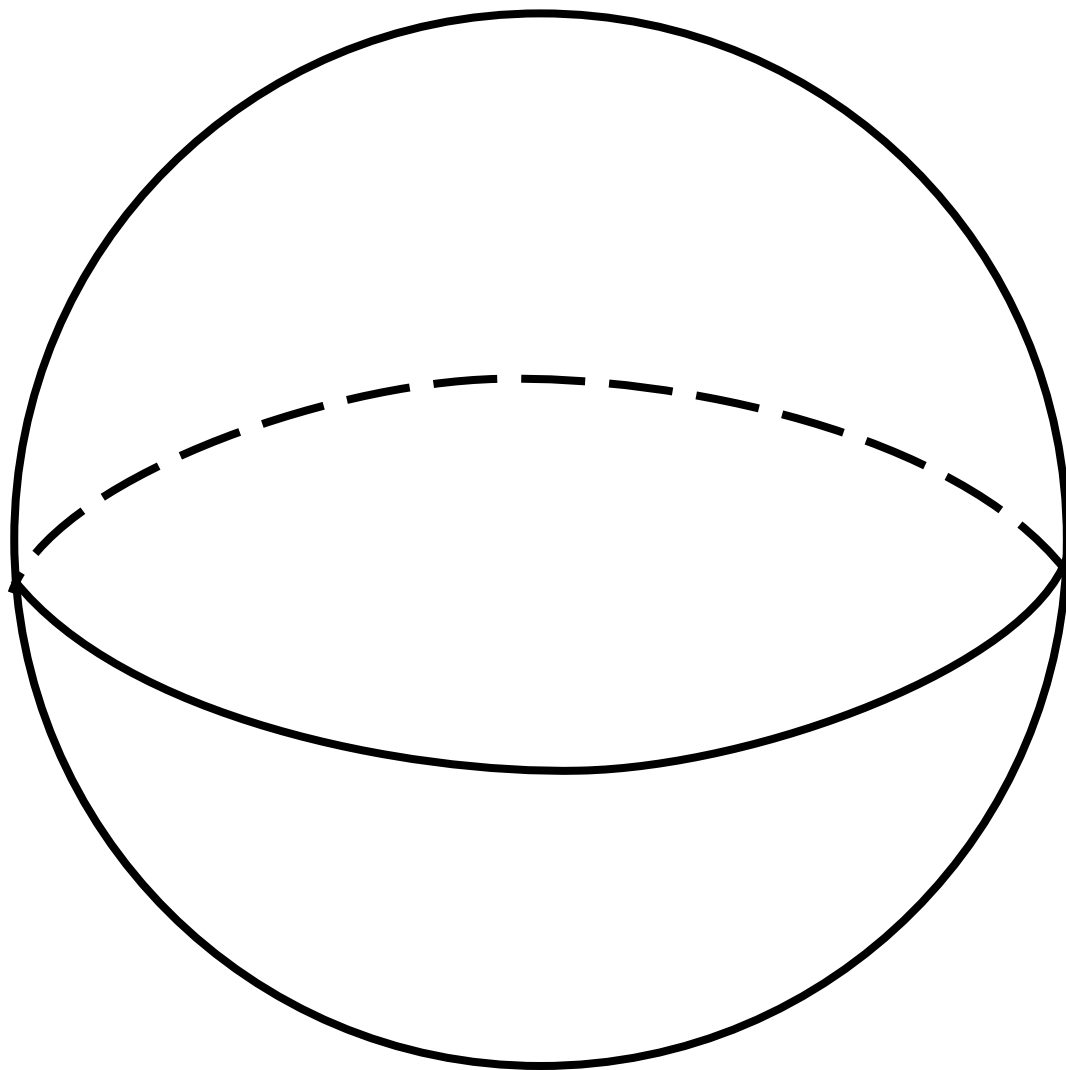
Combines both networks

How to represent normals

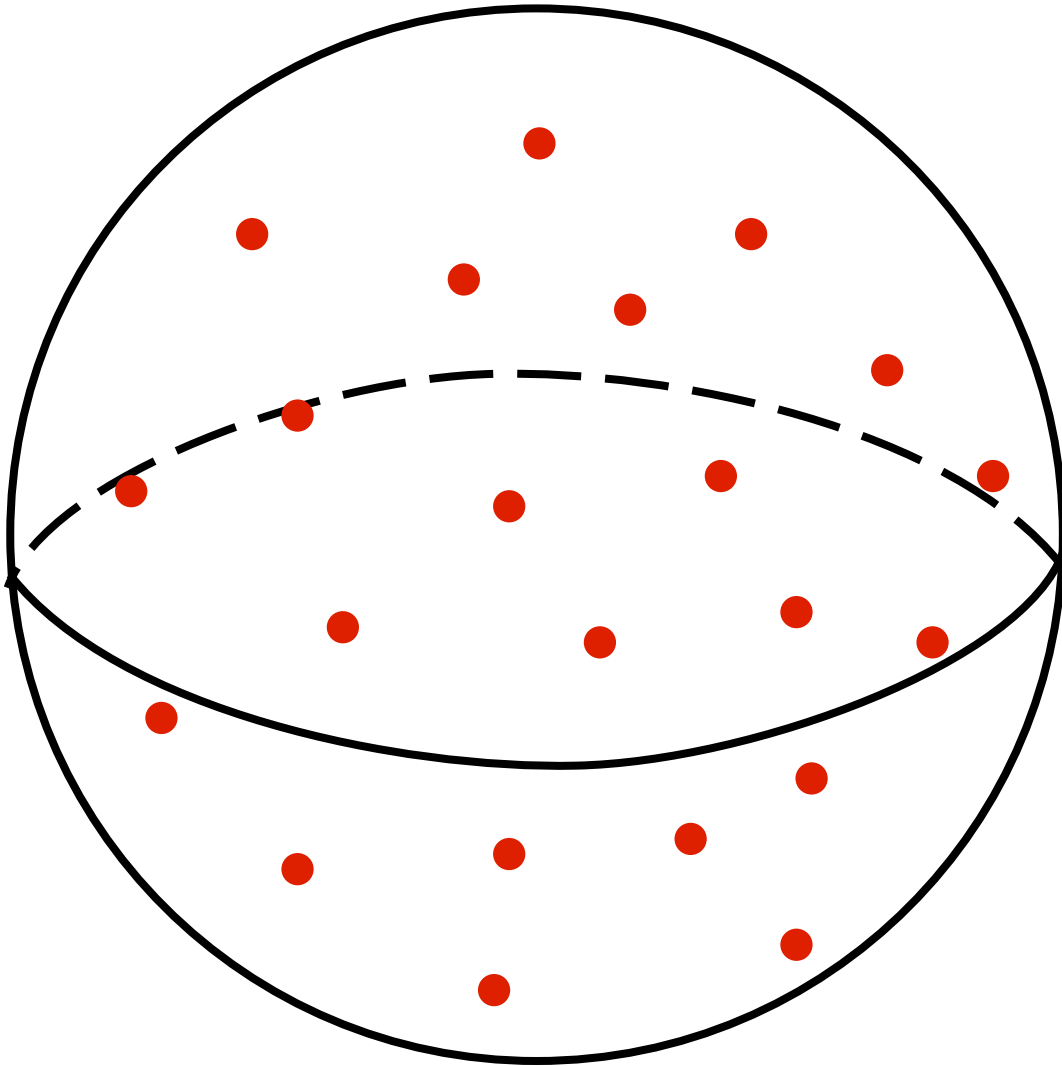
Normals lie in continuous space

Regression as Classification

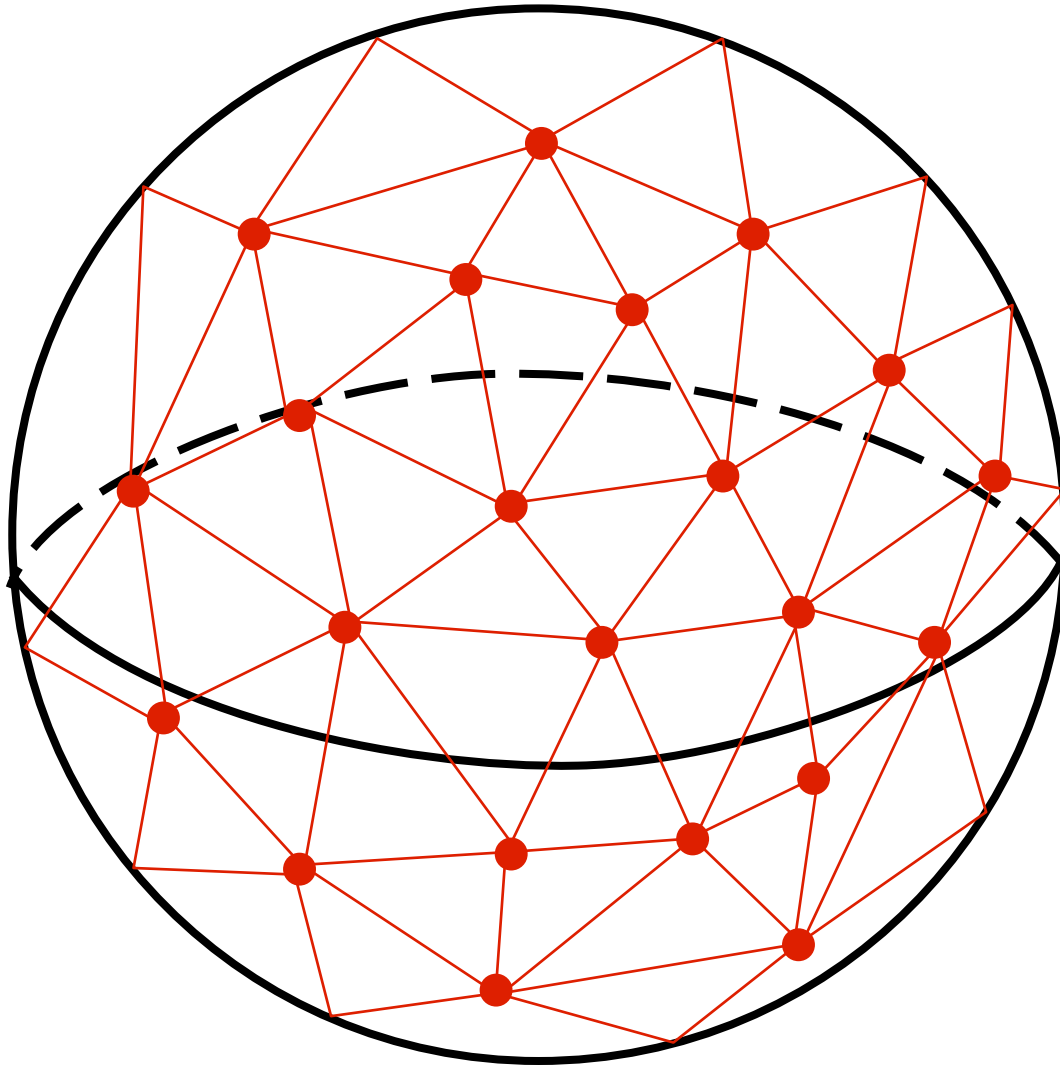
Surface normal triangular coding



Codebook with k-means



Delaunay Triangulation cover



Triangles as classes

Represent Surface Normals

Weighted sum of triangle corners

Loss Function

$$L(I, Y) = - \sum_{i=1}^{M \times M} \sum_{k=1}^K (1(y_i = k) \log F_{i,k}(I))$$

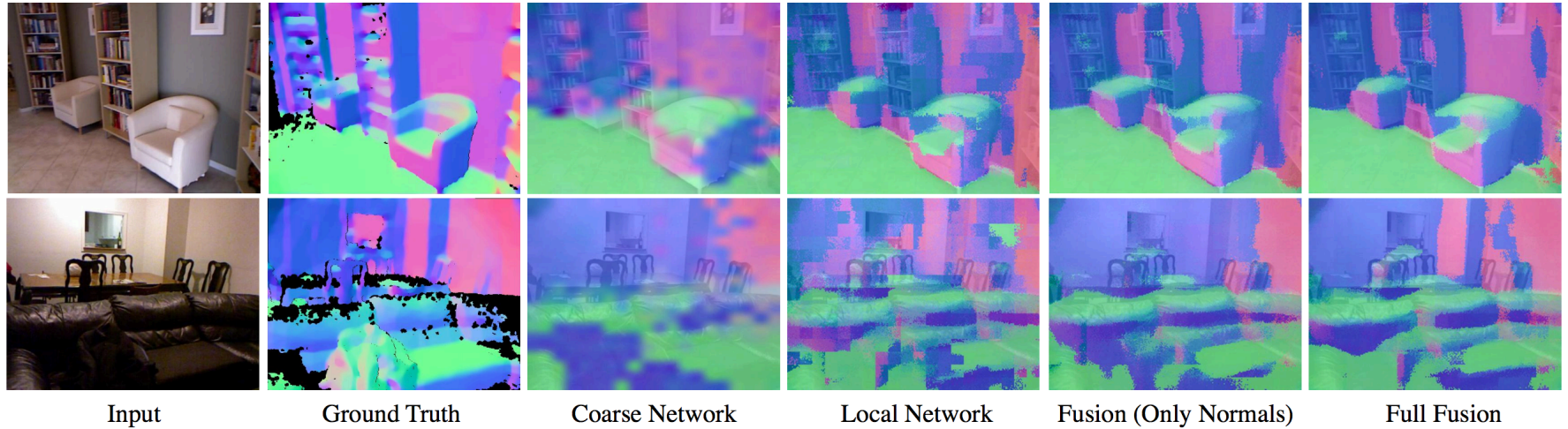


Table 2: Ablative Analysis

	Mean	Median	RMSE	11.25°	22.5°	30°
Full	25.0	13.8	35.9	44.2	63.2	70.3
Full (Soft)	24.2	17.3	32.2	36.8	58.5	68.7
Fusion (+VP)	25.3	14.4	35.9	42.7	62.5	69.9
Fusion (+Edge)	25.8	15.3	36.0	40.0	61.6	69.7
Fusion (+Layout)	25.8	14.9	36.3	41.1	61.9	69.5
Fusion	26.0	15.5	36.2	39.5	61.3	69.3
Bottom-up	32.2	23.5	42.0	27.2	48.5	58.5
Top-down	29.0	19.8	38.3	32.7	53.8	62.4
Eigen et al.(Fusion)	26.8	19.3	35.2	32.6	55.3	65.5
Eigen et al.(Coarse)	27.9	23.4	34.5	25.5	48.4	60.6

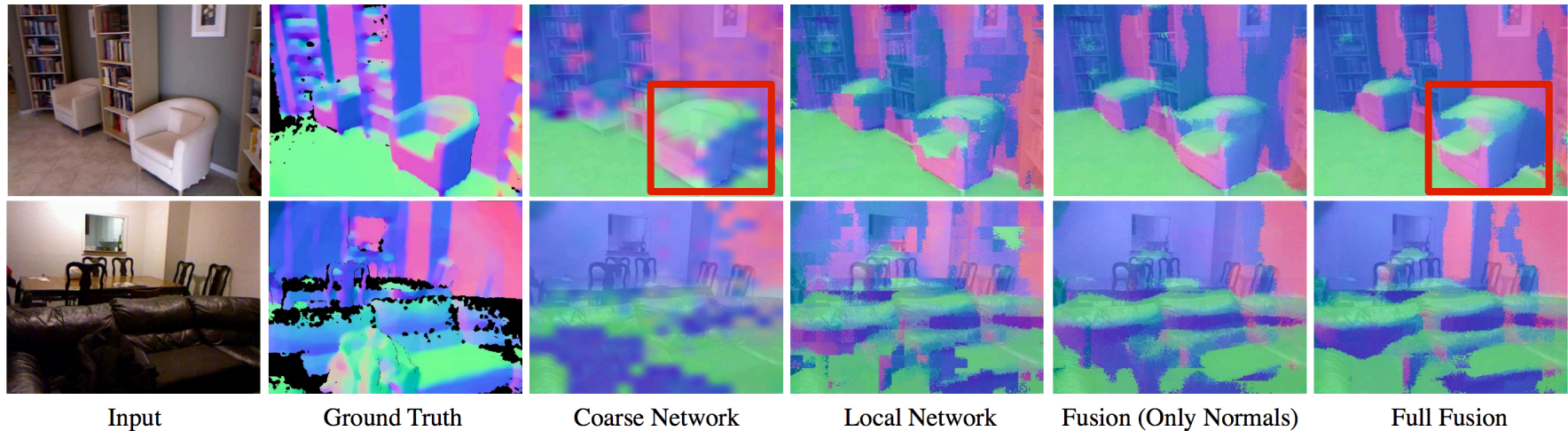


Table 2: Ablative Analysis

	Mean	Median	RMSE	11.25°	22.5°	30°
Full	25.0	13.8	35.9	44.2	63.2	70.3
Full (Soft)	24.2	17.3	32.2	36.8	58.5	68.7
Fusion (+VP)	25.3	14.4	35.9	42.7	62.5	69.9
Fusion (+Edge)	25.8	15.3	36.0	40.0	61.6	69.7
Fusion (+Layout)	25.8	14.9	36.3	41.1	61.9	69.5
Fusion	26.0	15.5	36.2	39.5	61.3	69.3
Bottom-up	32.2	23.5	42.0	27.2	48.5	58.5
Top-down	29.0	19.8	38.3	32.7	53.8	62.4
Eigen et al.(Fusion)	26.8	19.3	35.2	32.6	55.3	65.5
Eigen et al.(Coarse)	27.9	23.4	34.5	25.5	48.4	60.6

Thoughts

Do not address bas-relief

Incorporate Computer Graphics

Inverse problem

Given surface normals

How should the scene look?

What is the correct image?

Incorporate image
formation model

Why depth from single image