

Single Image Depth Estimation Using a Multi-Scale Convolutional Neural Network

Xiaoyu Zhang, Mingxiang Fan, Mengzheng Pan, Lili Zhu

Abstract—Compared with stereo image depth estimation, finding depth relations from a single image is less straightforward. Moreover, the mapping between a single image and the depth map is inherently ambiguous, and requires both global and local information. In this project, we present a Multi-Scale Deep Convolutional Neural Network for single image depth estimation. The method we used in this project employed two deep network stacks: a coarse global prediction based on the entire image, and another to refine this prediction locally. This method is evaluated on the NYU-depth v2 dataset and compares with several previous works including network structures of AlexNet and ResNet.

Index Terms—Convolutional Neural Network, Depth Estimation

I. MOTIVATION

Depth prediction from RGB images is a crucial topic in robotics, virtual reality, and 3D modeling because it is beneficial for understanding geometric relations within a scene. In turn, such relations help extract more information from objects and their environment, usually leading to enhancements in current recognition projects, as well as speed up the development of further applications, such as 3D modeling, physics and support models, robotics, and potentially reasoning about occlusions.

Nowadays, although many researchers have done much research on estimating depth based on stereo images or motion, there has been relatively little on predicting depth from a single 2D image. However, in real-world practice, the monocular case arises more often and it is reasonable to expect a wide and convenient usage of single image depth estimation. For example, images distributed on the web and social media outlets, real estate listings, and shopping sites, are all monocular cases. Therefore, we decided to focus on depth prediction for monocular cases.

Depth prediction for the monocular image is more difficult than stereo ones. Provided accurate image correspondences, depth can be recovered deterministically in the stereo case. With stereoscopic images, depth can be computed from local correspondences and triangularization, while estimating the geometry relationship of the camera positions will help with the accuracy. By contrast, methods for inferring depth from a single image have involved image segmentation, texture variations, texture gradients, interposition and shading. It is barely possible to get good result without machine learning algorithms to automatically learn these complex tasks. Thus

in our work, we performed supervised learning using convolutional neural network to achieve acceptable performance.

In this project, we present an approach for estimating depth from a single image. We directly regress on the depth using a neural network with two components: one that first estimates the global structure of the scene, then a second that refines it using local information. The network is trained using a loss that combines l_2 norm and scale-invariant error. We used raw data of NYU-Depth V2 [1], which is shown in Fig.1, to train our model, and validated our model on another dataset which was provided by Professor Matthew Johnson-Roberson.

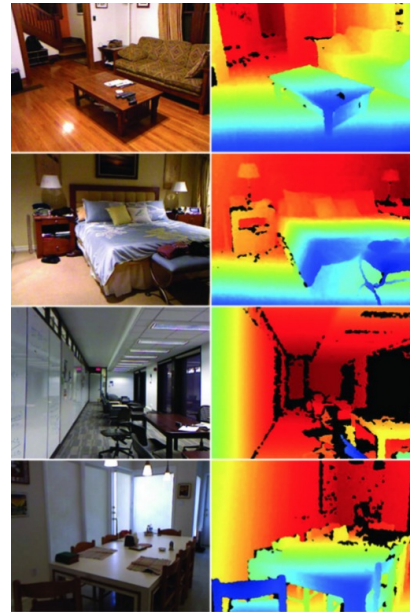


Fig. 1: Example of NYU dataset

II. PREVIOUS WORK

Stereo depth estimation has been a broadly used and extensively investigated approach of recovering depth information, using image pairs of the same scene to reconstruct 3D shapes. Scharstein et al. [2] have provided a survey and evaluation of many methods for 2-frame stereo correspondence. They use techniques like matching, aggregation and optimization to organize the system. Snavely et al. [3] apply a multiview stereo method to a creative application, matching across views of many uncalibrated consumer photographs of the same scene to create accurate 3D reconstructions of common landmarks.

In order to obtain better results and relax the need for careful camera alignment [4] [5] [6] [7], some machine learning

X. Zhang, M. Fan, M. Pan and L. Zhu were with the Department of Electrical and Computer Engineering, University of Michigan, Ann Arbor, MI, 48109, USA e-mail: zhxiaoyu@umich.edu

where \hat{y}_i is predicted depth map, y is ground truth, i is the index of each n pixels. We should notice that since there exists unavailable depth data points in raw data from our dataset, only those pixels with available depth information are considered as computing the error. Setting $\lambda = 0$ reduces the loss to element-wise l_2 norm, while $\lambda = 1$ is the scale-invariant error exactly. In practice, we choose $\lambda = 0.5$ to achieve a balance between both errors.

The following equivalent forms provides an additional ways to view metric.

$$\begin{aligned} D(y, \hat{y}_i) &= \frac{1}{n} \sum_i d_i^2 - \frac{\lambda}{n^2} \left(\sum_i d_i \right)^2 \\ &= \frac{1}{n} \sum_i d_i^2 - \frac{\lambda}{n^2} \sum_{i,j} d_i d_j \end{aligned}$$

The equation above expresses the error by comparing relationships between pairs of pixels i, j in the output: the difference between each pair of pixels should be a bit of different in prediction and ground truth for lower error. Our error function added a penalty $-\frac{\lambda}{n^2} \sum_{i,j} d_i d_j$ to the original l_2 error, which credits when two pixels with same direction are estimated to be opposite. In conclusion, if a prediction has mistakes which is consistent with another, it is an imperfect prediction.

In addition to the combined scale-invariant error, we also measure the performance according to several error metrics as comparison.

C. Training Loss

In addition to performance evaluation, our combined error function could also be used as training loss. We define training loss of each sample as the following equation:

$$L(y, \hat{y}_i) = \frac{1}{n} \sum_i d_i^2 - \frac{\lambda}{n^2} \left(\sum_i d_i \right)^2$$

where $d_i = \hat{y}_i - y_i$ and $\lambda \in [0, 1]$, $\log y$ is log of prediction. We can change λ to modify the weight of elementwise l_2 , while $\lambda = 1$ is the scale-invariant error exactly. We take the average, $\lambda = 0.5$, finding that this produces good absolute-scale predictions while slightly improving qualitative output.

V. EXPERIMENTS

In this section, we give qualitative results of our models and quantitative metric evaluations. All experiments are implemented on a desktop computer with GTX 1080 graphic card, core i7 6700k CPU and 32GB memory.

Since our model was originally inspired by Eigen's Network, we compared our result with the network proposed in their paper. Besides, we also implemented a network based on ResNet-18 in Fig.4 and a network based on Alexnet in Fig.3, and compared our result with these networks. Although ResNet can be as deep as 152 layers, we only compared our method with ResNet-18 because they have similar model complexity.

TABLE I: metric evaluation results

	$t < 1.25$	Abs rel diff	RMSE
AlexNet	0.9070	0.0899	0.1051%
ResNet18	0.8954	0.1101	0.1076%
Eigen	0.8951	0.1033	0.1009%
Eigen-modified	0.9010	0.0877	0.0920%

A. Dataset

We use NYU-Depth v2 Dataset [1] for this task. This data set is composed of video sequences from a variety of indoor scenes as recorded by both the RGB and Depth cameras from the Microsoft Kinect. We train our model on a subset of this NYU-Depth v2 dataset. Then, we also validated our model on another dataset which was provided by Professor Matthew Johnson-Roberson.

B. Evaluation Metrics

For ground truth depth images y and predicted images \hat{y} , we evaluate our method on three different metrics:

Percentage of pixels with relative error (larger means better performance)

$$t = \max\left(\frac{y}{\hat{y}}, \frac{\hat{y}}{y}\right) < 1.25$$

Absolute Relative Difference:

$$e = \left\| \frac{y - \hat{y}}{y} \right\|$$

Root Mean squared Error(RMSE):

$$e = \sqrt{\frac{1}{n} \|y - \hat{y}\|^2}$$

As shown in the Metric evaluation results Table.I, our modified-Eigen network is slightly better than the original Eigen network according to the evaluation, while surpassing the performance of Alexnet and resnet.

C. Strengths

We can also see from Fig.7 that our modified Eigen network gives best prediction result, taking smooth and accuracy into account.

All the predicted results could sometimes cover the invalid pixels of the raw Kinect data, since the network can learn to predict the depth information from other parts of the dataset, which enables the neural networks to fix the raw depth image from kinect.

Our model uses less parameter than original Eigen network, and gets better performance, which indicates that our network structure is more suitable for predicting depths.

Besides, we also trained our model and original Eigen model on another dataset from Professor Matthew Johnson-Roberson from the University of Michigan. This dataset contains complex outdoor scenes from computer simulation. The predicted depth images and the original depth images are shown in Fig.8. As we calculate the error of our prediction versus Eigen's network, our work have better performance on this more complex dataset, which indicates that our network structure is more reasonable.

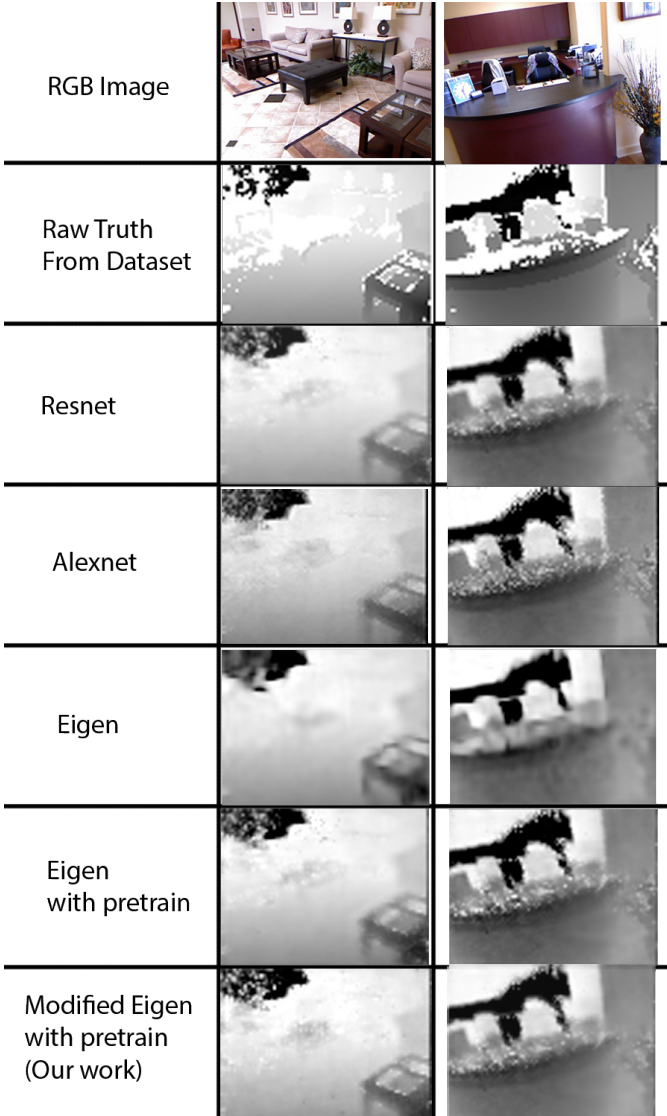


Fig. 7: Prediction results

D. Weakness

Estimating depth from a single image is a challenge task for the ambiguity along with different complicated lighting and shading conditions of a single image. Also, different datasets may have different kinds of assumptions which may increase the difficulty of directly applying one model trained on one dataset to another dataset and give as good results. In our experiment, we train our model on a subset of the NYU-Depth V2 dataset, which breaks many assumptions from general cases. Besides, NYU-Depth V2 dataset consists mainly indoor images while there are datasets and real-life scenes that contains mainly outdoor images with complicated lighting conditions. Therefore, our model may have relatively poor performance on some other datasets.

Apart from applying to different datasets, our result still has spaces of improvement compared to ground truth depth image. Although we achieved better performance than some traditional convolutional neuron networks such as VGG, AlexNet, our result can still overlook some small details showed in

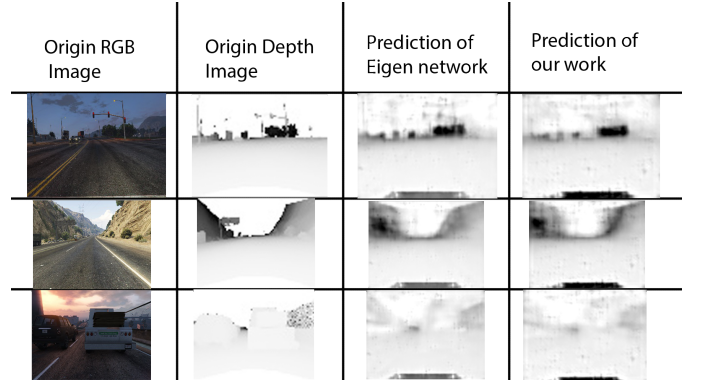


Fig. 8: prediction result on another dataset

ground truth map. On the other hand, because of using a fine-scale network to extract local details, sometimes the output also include some texture edges.

Besides, because we have limited calculation power, only small-scale networks could be examined. To achieve state-of-the-art performance, we may need more complex structure for depth estimation, and our network structure may not be suitable to expand.

E. Spaces of Improvement

To achieve a more state-of-the-art result, we can extend our method to add more scales to the system, such as successive finer-scaled local networks and extract the depth information which we cannot get from this model. We can also incorporate other 3D geometry information such as semantic segmentation and surface normal to help improve overall performance. Also, we can feed in our training system with images from different datasets which include complicated outdoor lighting conditions after taking proper preprocessing step to get better performance on different datasets.

VI. CONCLUSION

Our system accomplishes depth prediction from single images through the use of two deep networks, one that estimates the global depth structure, and another that refines it locally at finer resolution. We achieve better results than some neural network structures from previous work with similar amount of parameters. In future work, we plan to extend our method to incorporate further 3D geometry information, such as surface normals. We'll also try to explore the performance of multi-scale neural networks on different types of tasks, and investigate some more complex network structures with stronger hardwares, if possible.

ACKNOWLEDGMENT

This project is funded through the generous donation from parents of the team members in the form of tuition fee and computers(contribution shows in Fig.9). This project is advised by Prof.Pilanci from University of Michigan, and we got precious advises from the GSIs of EECS545. We would like to thank them for their love, advice and support.

	Mengzhen Pan	Lili Zhu	Mingxiang Fan	Xiaoyu Zhang
Pre-processing of datasets			√	
Network structure design		√		√
Tuning of hyper-parameters	√		√	
Training and testing	√	√	√	√
Evaluate result	√			
Paperworks		√	√	√

Fig. 9: contribution table

REFERENCES

- [1] P. K. Nathan Silberman, Derek Hoiem and R. Fergus, “Indoor segmentation and support inference from rgbd images,” in *ECCV*, 2012.
- [2] D. Scharstein and R. Szeliski, “A taxonomy and evaluation of dense two-frame stereo correspondence algorithms,” *International journal of computer vision*, vol. 47, no. 1-3, pp. 7–42, 2002.
- [3] N. Snavely, S. M. Seitz, and R. Szeliski, “Photo tourism: exploring photo collections in 3d,” in *ACM transactions on graphics (TOG)*, vol. 25, pp. 835–846, ACM, 2006.
- [4] K. Konda and R. Memisevic, “Unsupervised learning of depth and motion,” *arXiv preprint arXiv:1312.3429*, 2013.
- [5] R. Memisevic and C. Conrad, “Stereopsis via deep learning,” in *NIPS Workshop on Deep Learning*, vol. 1, p. 2, 2011.
- [6] K. Yamaguchi, T. Hazan, D. McAllester, and R. Urtasun, “Continuous markov random fields for robust stereo estimation,” *Computer Vision–ECCV 2012*, pp. 45–58, 2012.
- [7] F. H. Sinz, J. Q. Candela, G. H. Bakir, C. E. Rasmussen, and M. O. Franz, “Learning depth from stereo,” in *Joint Pattern Recognition Symposium*, pp. 245–252, Springer, 2004.
- [8] A. Saxena, S. H. Chung, and A. Y. Ng, “Learning depth from single monocular images,” in *Advances in neural information processing systems*, pp. 1161–1168, 2006.
- [9] A. Saxena, M. Sun, and A. Y. Ng, “Learning 3-d scene structure from a single still image,” in *Computer Vision, 2007. ICCV 2007. IEEE 11th International Conference on*, pp. 1–8, IEEE, 2007.
- [10] B. Liu, S. Gould, and D. Koller, “Single image depth estimation from predicted semantic labels,” in *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pp. 1253–1260, IEEE, 2010.
- [11] D. Hoiem, A. A. Efros, and M. Hebert, “Automatic photo pop-up,” *ACM transactions on graphics (TOG)*, vol. 24, no. 3, pp. 577–584, 2005.
- [12] L. Ladicky, J. Shi, and M. Pollefeys, “Pulling things out of perspective,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 89–96, 2014.
- [13] K. Karsch, C. Liu, and S. B. Kang, “Depth extraction from video using non-parametric sampling,” in *European Conference on Computer Vision*, pp. 775–788, Springer, 2012.
- [14] C. Liu, J. Yuen, A. Torralba, J. Sivic, and W. T. Freeman, “Sift flow: Dense correspondence across different scenes,” in *European conference on computer vision*, pp. 28–42, Springer, 2008.
- [15] D. Eigen, C. Puhrsch, and R. Fergus, “Depth map prediction from a single image using a multi-scale deep network,” in *Advances in neural information processing systems*, pp. 2366–2374, 2014.
- [16] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in neural information processing systems*, pp. 1097–1105, 2012.
- [17] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.