# Monocular Depth Estimation with Hierarchical Fusion of Dilated CNNs and Soft-Weighted-Sum Inference

Bo Li, Yuchao Dai, Mingyi He

*Northwestern Polytechnical University, Xi'an, China*

**Abstract**

Monocular depth estimation is a challenging task in complex compositions depicting multiple objects of diverse scales. Albeit the recent great progress thanks to the deep convolutional neural networks (CNNs), the state-of-the-art monocular depth estimation methods still fall short to handle such real-world challenging scenarios.

In this paper, we propose a deep end-to-end learning framework to tackle these challenges, which learns the direct mapping from a color image to the corresponding depth map. First, we represent monocular depth estimation as a multi-category dense labeling task by contrast to the regression based formulation. In this way, we could build upon the recent progress in dense labeling such as semantic segmentation. Second, we fuse different side-outputs from our front-end dilated convolutional neural network in a hierarchical way to exploit the multi-scale depth cues for depth estimation, which is critical to achieve scale-aware depth estimation. Third, we propose to utilize soft-weighted-sum inference instead of the hard-max inference, transforming the discretized depth score to continuous depth value. Thus, we reduce the influence of quantization error and improve the robustness of our method. Extensive experiments on the NYU Depth V2 and KITTI datasets show the superiority of our method compared with current state-of-the-art methods. Furthermore, experiments on the NYU V2 dataset reveal that our model is able to learn the probability distribution of depth.

soft inference, dilated convolution.

## 1. Introduction

Depth estimation aims at predicting pixel-wise depth for a single or multiple images, which is an essential intermediate component toward 3D scene understanding. It has been shown that depth information can benefit tasks such as recognition [1, 2], human computer interaction [3], and 3D model reconstruction [4]. Traditional techniques have predominantly worked with multiple images to make the problem of depth prediction well-posed, which include $N$-view reconstruction, structure from motion (SfM) and simultaneous localization and mapping (SLAM) [].

However depth estimation from a monocular single viewpoint lags far behind its multi-view counterpart. This is mainly due to the fact that the problem is ill-posed and inherently ambiguous: a single image on its own does not provide any depth cue explicitly (*i.e.*, given a color image of a scene, there are infinite number of 3D scene structures explaining the 2D measurements exactly). When specific scene dependent knowledge is available, depth estimation or 3D reconstruction from single images can be achieved by utilizing geometric assumptions such as the "Blocks World" model [5], the "Origami World" model [6], shape from shading [7] and repetition of structures [8]. However, these cues typically work for images with specific structures and may not be applied to general scenarios.

Recently, learning based monocular depth estimation methods that predicting scene geometry directly by learning from data, have gained popularity. Typically, such approaches recast the underlying depth estimation problem in a pixel-level scene labeling pipeline by exploiting relationship between monocular image and depth. Fully-convolutional neural network has been proved to be an effective method to solve these kinds of problems. There have been considerable progress in applying deep convolutional neural network (CNN) to this problem and excellent performances have been achieved [7, 8, 9, 10, 11, 12, 13, 14].

Albeit the above success, state-of-the-art monocular depth estimation meth-

ods still fall short to handle real world challenging complex decompositions depicting multiple objects of diverse scales due to the following difficulties: 1) the serious data imbalance problem due to the perspective effect, where samples with small depths are much more than samples with large depths; 2) there are more rapid changes in depth value compared with other dense predictions tasks such as semantic labeling and 3) long range context information is needed handle the scale ambiguity in depth estimation. Even though there have been various post-processing methods to refine the estimated depth from the deep network map [7, 8, 9, 10, 11, 12, 13, 14], the bottleneck in improving monocular depth estimation is still the specially designed CNN architecture, which is highly desired.

In this paper, we present a deep CNN based framework to tackle the above challenges, which learns the direct mapping from the color image to the corresponding depth map in an end-to-end manner. We recast monocular depth estimation as a multi-category dense labeling as contrast to the widely used regression formulation. Our network is based on the deep residual network [15], where dilated convolution and hierarchical fusion layers are designed to expand the receptive field and to fuse multi-scale depth cues. In order to reduce the influence of quantization error and improve the robustness of our method, we propose to use a soft-weighted-sum inference. Extensive experimental results show that even though we train our network as a standard classification task with the multinomial logistic loss, our network is able to learn the the probability distribution among different categories. The overall flowchart of our framework is illustrated in Fig. 1.

Our main contributions can be summarized as:

- We propose a deep end-to-end learning framework to monocular depth estimation by recasting it as a classification task, where both dilated convolution and hierarchical feature fusion are used to learn the scale-aware depth cues.

- Our network is able to output the probability distribution among differ-

ent depth labels. We propose a soft-weighted-sum inference, which could reduce the influence of quantization error and improve the robustness.

- Our method achieves the state-of-the-art performance on both indoor and outdoor benchmarking datasets, NYU V2 and KITTI dataset.
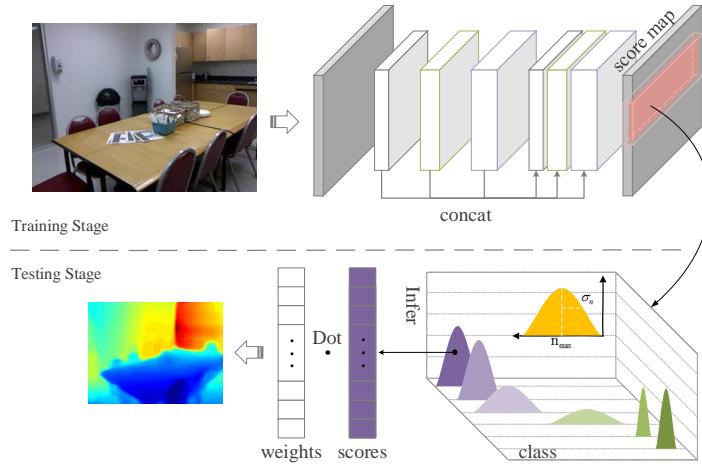


Figure 1: Flowchart of our monocular depth estimation framework, which is built upon deep Residual network [15] and consists of dilated convolution and hierarchical feature fusion. Soft-weighted-sum inference is used to predict continuous depth values from the discrete depth labels. We also illustrate typical probability distribution of labels from the network, which shows that our classification based framework is able to learn the similarity between labels.

## 2. Related work

In this section, we briefly review related works for monocular depth estimation, which can be roughly categories as conventional MRF/CRF based methods and deep learning based methods.

**MRF/CRF Based Methods:** Seminal work by Saxena *et al.* [16, 17] tackles the problem with a multi-scale Markov Random Field (MRF) model, with the parameters of the model learned through supervised learning. Liu *et al.* [18] estimated the depth map from predicted semantic labels, achieving improved

performance with a simpler MRF model. Ladicky *et al.* [19] showed that perspective geometry can be used to improve results and demonstrated how scene labeling and depth estimation can benefit each other under a unified framework, where a pixel-wise classifier was proposed to jointly predict a semantic class and a depth label from a single image. Besides these parametric methods, other works such as [20, 21, 22] recast monocular depth estimation in a nonparametric fashion, where the whole depth map is inferred from candidate depth maps. Liu *et al.* [21] proposed a discrete-continue CRFs, which aims to avoid the over-smoothing and maintain occlusion boundaries. Anirban *et al.* [] proposed a Neural Regression Forest model for this problem. These works provide important insights and cues for single image depth estimation problem, while most of them utilized the hand-crafted features thus limited their performance especially for complex scenarios.

**Deep Learning Based Methods:** Recently, monocular depth estimation has been greatly advanced thanks to deep convolutional neural network (CNN). Eigen *et al.* [23] presented a framework by training a large hierarchical deep CNN. However, partly due to the fully connect layers used in the network model, their network have to be trained with very large scale data. By contrast, Li *et al.* [7] proposed a patch-based CNN framework and a hierarchical-CRF model to post-process the raw estimated depth map, which significantly reduces the number of training image needed. Liu *et al.* [8] proposed a CRF-CNN joint training architecture, which could learn the parameters of the CRF and CNN jointly. Wang *et al.* [9] proposed a CNN architecture for joint semantic labeling and monocular depth prediction. Chen *et al.* [24] proposed an algorithm to estimate metric depth using annotations of relative depth.

Very recently, Laina *et al.* [12] proposed using the Huber loss instead of the $L_2$ loss to deal with the long tail effect of the depth distribution. Cao *et al.* [11] demonstrated that formulating depth estimation as a classification task could achieve better results than regression with $L_2$ loss, while insufficient analysis is given for the success. In addition, different with our method, they used hard-max inference in the testing phase. Xu *et al.* [13] proposed a Multi-Scale

Continuous CRFs to better extract the hierarchical information and improve the smoothness of the final results. Our hierarchical information fusion strategy is much simpler than [13], while we also achieve comparable results.

**Unsupervised monocular depth learning** Besides the above methods using ground truth depth maps to supervise the network learning, there is another group of methods that using novel view synthesis to supervised the network learning by exploiting the availability of stereo images and image sequences [25] [14] [26] [27] citeUnsupervised-Depth-Motion. Garg *et al.*[25] proposed to train a network for monocular depth estimation using an image reconstruction loss, where a Taylor approximation is performed to linearize the loss. Godard *et al.*[14] replaced the use of explicit depth data during training with easier-to-obtain binocular stereo footage, which enforces consistency between the disparities produced relative to both the left and right images, leading to improved performance and robustness compared to existing approaches. Along this pipeline, Zhou *et al.*[26] presented an unsupervised learning framework for the task of monocular depth and camera motion estimation from unstructured video sequences based on image warping to evaluate the image error. Kuznietsov *et al.*[27] learned depth in a semi-supervised way, where sparse ground-truth depth and photoconsistency are jointly used. Ummenhofer *et al.*[28] trained a convolutional network end-to-end to compute depth and camera motion from successive, unconstrained image pairs, where the architecture is composed of multiple stacked encoder-decoder networks.

The key supervision signal for these "unsupervised" methods comes from the task of novel view synthesis: given one input view of a scene, synthesize a new image of the scene seen from a different camera pose. Essentially, pairs of rectified stereo images or consecutive image frames have already encode the depth information implicitly.

Our work is also related to the works on FCN (fully convolutional network) based dense labeling. Long *et al.* [29] proposed the fully convolution neural network for semantic segmentation, which is widely used in other dense labeling problems. Hariharan *et al.* [30] presented that low-level CNN feature is bet-

6

ter to the boundary preserving and object location. Recently, Yu *et al.* [31] demonstrated that dilated convolution could expand the receptive field of the corresponding neuron while keeping the resolution of the feature map. Chen [32] successfully utilized the dilated convolution on the semantic problem and show how to build them on the pre-trained CNN.

## 3. Our Framework

Targeting at handling the real world challenges with the current state-of-the-art methods, we propose a deep end-to-end learning framework to monocular depth estimation, which learns the direct mapping from a color image to the corresponding depth map. Our framework to monocular depth estimation consists of two stages: model training with classification loss and inference with soft-weighted sum. First, by recasting monocular depth estimation as multi-class labeling, we design an hierarchical fusion dilated CNN to learn the mapping from an RGB image to the corresponding depth score map directly. Our network architecture hierarchically fuses multi-scale depth features, which is important to achieve scale-aware monocular depth estimation. Second, we propose a soft-weighted-sum inference as contrast to the hard-max inference, which transfers the discretized depth scores to continuous depth values. In this way, we could reduce the influence of quantization error and improve the robustness.

### 3.1. Network Architecture

Our CNN architecture is illustrated in Fig. 2, in which the weights are initialized from a pre-trained 152 layers deep residual CNN (ResNet) [15]. Different from existing deep network [33], ResNet [15] explicitly learns residual functions with reference to the layer inputs, which makes it easier to optimize with higher accuracy from considerably increased network depth. ResNet [15] was originally designed for image classification. In this work, we re-purpose it to make it suitable to our depth estimation task by

- Removing all the fully connect layers. In this way, we greatly reduce the number of model parameters as most of the parameters are in the

7

fully connect layers [10]. Although preserving the fully connect layers is beneficial to extract long range context information, our experiments show that it is unnecessary in our network thanks to dilated convolution.

- Taking advantage of the dilated convolution [31]. Dilated convolution could expand the receptive field of the neuron without increasing the number of model parameters. Furthermore, with the dilated convolution, we could remove some pooling layers without decreasing the size of receptive field of correspondent neurons. In addition, we could keep the resolution of the feature map and final results, *i.e.*, the output resolution has been increased.

- Hierarchal fusion. We concatenate intermediate feature maps with the final feature map directly. This skip connection design could benefit the multi-scale feature fusion and boundary preserving.
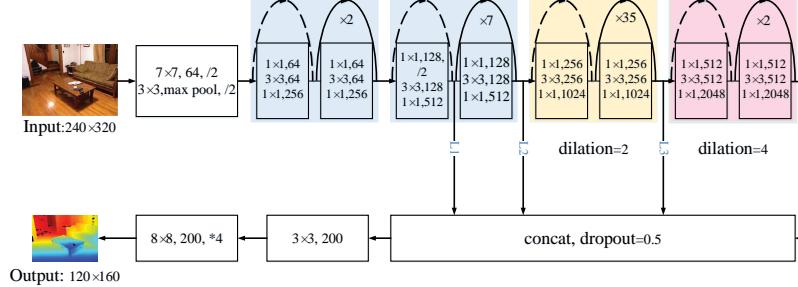


Figure 2: Illustration of our network architecture. The detail of the basic residual block could be referred to [15]. $\times n$ means the block repeats $n$ times. We present all the hyperparameters of convolution and pooling layers. All the convolution layers are followed by batch normalization layer except for the last one. /2 means the layer's stride is 2. *4 means the deconv layer's stride is 4. Dilation shows the dilated ratio of the corresponding parts. $L1, \cdots L6$ are our skip connection layers.

**Dilated Convolution**: Recently, dilated convolution [31] has been successfully utilized in deep convolutional neural network, which could expands the field of perception without increasing the number of model parameters. Spe-

cially, let $F : \mathcal{Z}^2 \rightarrow \mathcal{R}$ be a discrete function. Let $\Omega_r = [r, r]^2 \cap \mathcal{Z}^2$ and let $k : \Omega_r \rightarrow \mathcal{R}$ be a discrete filter of size $(2r+1)^2$. The discrete convolution filter $*$ can be expressed as

$$(F * k)(\mathbf{p}) = \sum_{\mathbf{s}+\mathbf{t}=\mathbf{p}} F(\mathbf{s})k(\mathbf{t}). \tag{1}$$

We now generalize this operator. Let $l$ be a dilation factor and let $*_l$ be defined as

$$(F *_l k)(\mathbf{p}) = \sum_{\mathbf{s}+l\mathbf{t}=\mathbf{p}} F(\mathbf{s})k(\mathbf{t}). \tag{2}$$

We refer to $*_l$ as a dilated convolution or an $l$-dilated convolution. The conventional discrete convolution $*$ is simply the 1-dilated convolution. An illustration of dilated convolution could be found in Fig. 3.
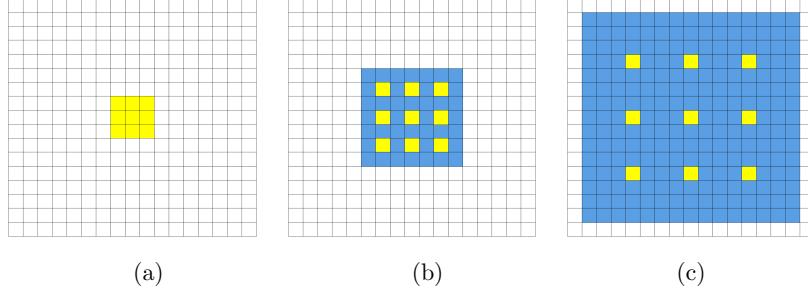


(a)  (b)  (c)

Figure 3: Systematic dilation supports exponential expansion of the receptive field without loss of resolution or coverage. (a), (b), (c) are 1-dilated, 2-dilated, 4-dilated convolution respectively. And the corresponding receptive fields are $3 \times 3$, $7 \times 7$, and $15 \times 15$. The receptive field grows exponentially while the number of parameters is fixed.

**Hierarchical Fusion:** As the CNN is of hierarchical structure, which means high-level neurons have larger receptive field and more abstract features, while the low-level neurons have smaller receptive field and more detail information. Thus, combining multi-scale informations for pixel-level prediction tasks have received considerable interests.

We propose to concatenate the high-level feature map and the intermediate feature map. The skip connection structure benefits both the multi-scale fusion and boundary preserving. In our network, the $L1, L2, L3, L4$ layers are of the same size, we concatenate them directly.

In conclusion, we briefly summarize our final network design. Typically, the pre-trained residual network is consisted of 4 parts. We remove the max-pooling layer in the last 2 parts and expand the corresponding convolution kernel with dilation 2 and 4 respectively. Then, a concatenation layer is added to fuse the hierarchical multi-scale informations from layers $L1 - L4$. The last two layers of our network are convolution layer and deconvolution layer. The parameters setting is presented in Fig. 2.

### 3.2. Soft-Weighted-Sum Inference

We reformulate depth estimation as classification task by equally discretizing the depth value in log space. Specifically,

$$l = round((\log(d) - \log(d_{min}))/q), \tag{3}$$

where $l$ is the quantized label, $d$ is the continuous depth value, $d_{min}$ is the minimum depth value in the dataset or set to be a small value like 0.1. $q$ is the width of the quantization bin.

With the quantization label, we train our network with the multinomial logistic loss.

$$L(\theta) = -\left[\sum_{i=1}^{N}\sum_{k=1}^{K} 1\{y^{(i)} = k\} \log \frac{\exp(\theta^{(k)T}x^{(i)})}{\sum_{i=1}^{K}\exp(\theta^{(j)T}x^{(i)})}\right], \tag{4}$$

where $N$ is the number of training samples, $\exp(\theta^{(k)T}x^{(i)})$ is the probability of label $k$ of sample $i$, and $k$ is the ground truth label.

In the testing stage, we propose to use the soft-weighted-sum inference. It is worth noting that, this method transforms the predicted score to the continuous depth value in a natural way. Specifically:

$$\hat{d} = \exp\{\mathbf{w}^T\mathbf{p}\}, w_i = \log(d_{min}) + q * i, \tag{5}$$

where $\mathbf{w}$ is the weight vector of depth bins. $\mathbf{p}$ is the output score. In our experiments, we set the number of bins to 200.

10

*3.3. Data Augmentation*

Although the training dataset is of ten thousands images, we still find the data augmentation is important to improve the final performance. In this work, we augment our dataset by 4 times for both the NYU v2 and the KITTI dataset. The augmentation methods we utilized include:

- Color: Color channels are multiplied by a factor $c \in [0.9, 1.1]$ randomly.

- Scale: We scale the input image by a factor of $s \in [1.3, 1.5]$ randomly and crop the center patch of images to match the network input size.

- Left-Right flips: We flip left and right images horizontally.

- Rotation: We rotate the input image randomly by a factor of $r \in [-5, 5]$.

*3.4. Implementation details*

Before proceeding to the experimental results, we give implementation details of our method. Our implementation is based on the efficient CNN toolbox: caffe [34] with an NVIDIA Tesla Titian X GPU.

The proposed network is trained by using stochastic gradient decent with batch size of 1 (This size is too small, thus we average the gradient of 8 iterations for one back-propagation), momentum of 0.9, and weight decay of 0.0004. Weights are initialized by the pre-trained model from ResNet [15]. The network is trained with iterations of 50k by a fixed learning rate 0.001 in the first 30k iterations, then divided by 10 every 10k iterations.

## 4. Experimental Results

In this section, we report our experimental results on monocular depth estimation for both outdoor and indoor scenes. We used the NYU V2 dataset, and the KITTI dataset, as they are the the largest open dataset we can access at present. We compared our method with the state-of-the-art methods published recently.

For quantitative evaluation, we report errors obtained with the following metrics, which have been extensively used in [16, 18, 23, 19, 21].

- Threshold: % of $d_i$ s.t. $\max\left(\frac{\hat{d}_i}{d_i}, \frac{d_i}{\hat{d}_i}\right) = \delta < thr$

- Mean relative error (Rel): $\frac{1}{|T|} \sum_{d \in T} |\hat{d} - d|/d$

- Mean $\log_{10}$ error ($\log_{10}$): $\frac{1}{|T|} \sum_{d \in T} |\log_{10} \hat{d} - \log_{10} d|$

- Root mean squared error (Rms): $\sqrt{\frac{1}{|T|} \sum_{d \in T} \|\hat{d} - d\|^2}$

where $d$ is the ground truth depth, $\hat{d}$ is the estimated depth, and $T$ denotes the set of all points in the images.

### 4.1. NYU V2 dataset

The NYU V2 dataset [4] contains around 240k RGB-depth image pairs, of which comes from 464 scenes, captured with a Microsoft Kinect. The official split consists of 249 training and 215 testing scenes. We equally sampled frames out of each training sequence, resulting in approximately 24k unique images. After off-line augmentations, our dataset comprises of approximately 96k RGB-D image pairs. We fill in the invalid pixels of the raw depth map with the "colorization" method, which is provided in the toolbox of NYU V2 dataset [4].

The original image resolution is $480 \times 640$. We downsampled the images to $240 \times 320$ as our network input. The resolution of the our network output is $120 \times 160$, which is half of the input size. In this dataset, we quantize the depth value into 200 bins.

In Table 1, we compared our method with the recent published state-of-the-art methods [10, 11, 12, 13].

In Fig. 4, we provide a qualitative comparison of our method with [12] and [23]. (We compare with these methods due to they published their results and they are the state-of-the-art methods at present). From Fig.4, it is clear to observe that our results are of high visual quality, although we have not applied any post-processing methods.

12

Table 1: Depth estimation results on the NYU v2 dataset, ∗ represent the results are only calculated on the valuable pixels

| Method | Train Num | Accuracy (higher is better) | | | Error (lower is better) | | |
|---|---|---|---|---|---|---|---|
| | | $\delta < 1.25$ | $\delta < 1.25^2$ | $\delta < 1.25^3$ | Rel | log10 | Rms |
| Eigen *et al.* [10] | 120K | 76.9% | 95.0% | 98.8% | 0.158 | - | 0.641 |
| Cao *et al.* [11] | 120k | 80.0% | 95.6% | 98.8% | 0.148 | 0.063 | 0.615 |
| Laina *et al.* [12] | 12k | 81.1% | 95.3% | 98.8% | 0.127 | 0.055 | 0.573 |
| Xu *et al.* [13] | 95k | 81.1% | 95.4% | 98.7% | **0.121** | **0.052** | 0.586 |
| Ours | 12k | **82.0%** | **96.0%** | **98.9%** | 0.139 | 0.058 | **0.505** |

### *4.2. KITTI dataset*

The KITTI dataset [35] consists of a large number of outdoor street scene images of the resolution $376 \times 1242$. We utilized the "Eigen" training/testing split, which consists of 22600 training images and 697 testing images. We fill in the invalid pixel of the raw depth map with the "colorization" method, which is provided in the toolbox of NYU V2 dataset [4]. For the error calculation, We only consider the lower crop of the image of size $256 \times 1242$. While in the training phase, we input the entire image to the network for more context information. We compare with the state-of-the-art methods Eigen *et al.* [23], Garg *et al.* [25] and Godard *et al.* [14].

The original image resolution is $376 \times 1240$. We downsampled the images to $188 \times 620$ as our network input. The resolution of the our network output is $94 \times 310$, which is half of the input size. For this dataset, we quantize the depth value into 50 bins.

In Table 2, we compared our method with the recent published state-of-the-art methods [23, 25, 14].

In Fig. 5, we provide a qualitative comparison of our method with [10] and [14]. (We compare with these methods due to they published their results and they are the state-of-the-art methods at present). From Fig.5, it is clear to observe that our results are of high visual quality, although we have not applied any post-processing methods.
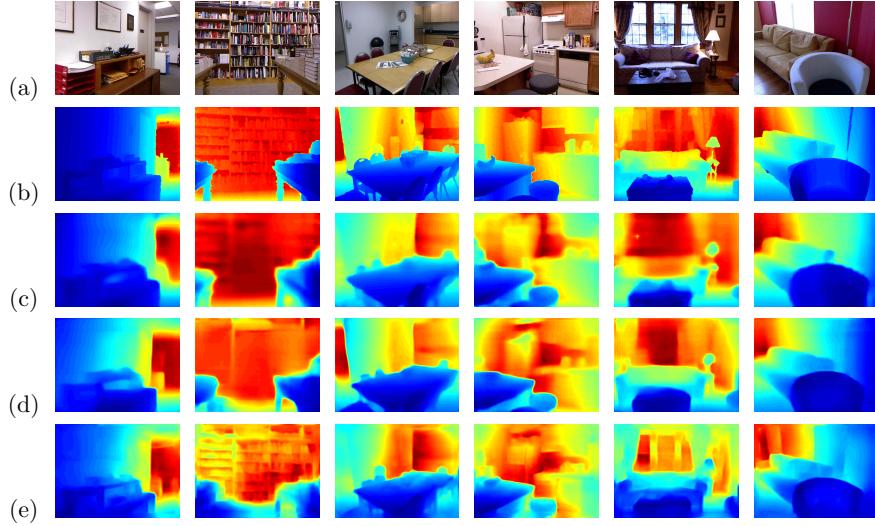
Figure 4: Qualitative comparison of the estimated depth map on the NYU V2 dataset with our method and some state-of-the-art methods. Color indicates depth (red is far, blue is close). (a) RGB. (b) The Ground Truth. (c) Results of our proposed method. (d) Results of [12]. (e) Results of [10]

## 5. Performance Analysis

In this section, we present more analysis of our model, where the experiments are conducted on the NYU V2 dataset. First, we present a component analysis of our network architecture design, *i.e.*, the contribution of each component. Second, we analyze the distribution of our network output, which demonstrates the necessary of our soft-weighted-sum inference strategy.

Table 2: Depth estimation results on the KITTI dataset.

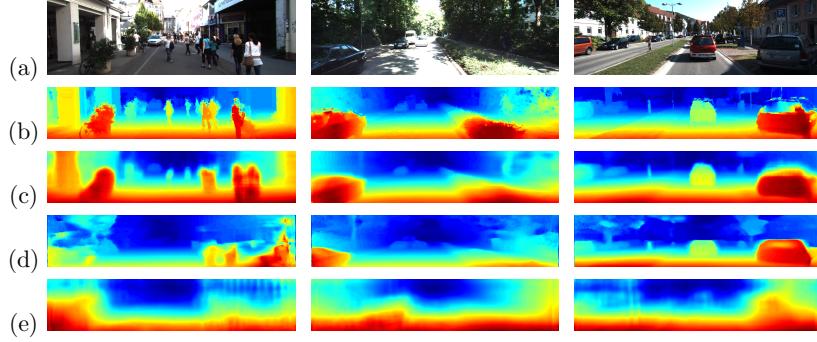| Method | Train Num | Accuracy (higher is better) | | | Error (lower is better) | | |
|---|---|---|---|---|---|---|---|
| | | $\delta < 1.25$ | $\delta < 1.25^2$ | $\delta < 1.25^3$ | Rel | log10 | Rms |
| Eigen *et al.* [23] | 22600 | 70.2% | 89.0% | 95.8% | 0.203 | - | 6.307 |
| Godard *et al.* [14] | 22600 | 80.3% | 92.2% | 96.4% | 0.148 | - | 5.927 |
| Godard *et al.* [14] cap 50m | 22600 | 81.8% | 93.1% | 96.9% | 0.140 | - | 4.471 |
| Ours | 22600 | **85.6%** | **96.2%** | **98.8%** | **0.113** | **0.049** | **4.687** |
| Ours cap 50m | 22600 | **86.4%** | **96.4%** | **98.9%** | **0.109** | **0.047** | **3.906** |

Figure 5: Qualitative comparison of the depth map estimated on KITTI dataset. Color indicates depth (red is close, blue is far). (a) RGB. (b) The Ground Truth. (c) Results of proposed method. (d) Results of [14]. (e) Results of [10]

### 5.1. Effect of Architecture Design

In order to explore the effectiveness of our hierarchical fusion dilated CNN, we conduct the following component analyze experiments. First, we utilize the normal convolution kernel instead of the dilated convolution kernel in the last 2 parts of our network. Second, we remove the skip connection structure. At last, we use the network without dilated convolution and skip connection. The corresponding experimental results are presented in Tab. 3. As we can see, both dilated convolution and hierarchical fusion play important roles in achieving improved performance.

Table 3: Component evaluation for our CNN architecture design and soft-weighted sum inference.

| method | Accuracy (%) | | | Error | | |
|---|---|---|---|---|---|---|
| | $\delta < 1.25$ | $\delta < 1.25^2$ | $\delta < 1.25^3$ | Rel | log10 | Rms |
| no dilation | 78.02% | 94.61% | 98.52% | 0.157 | 0.066 | 0.559 |
| no concat layer | 81.64% | 95.9% | 98.8% | 0.141 | 0.059 | 0.509 |
| ours hard-max | 81.82% | 95.53% | 98.53% | 0.142 | 0.060 | 0.531 |
| ours soft-weighted sum | **82.0%** | **96.0%** | **98.9%** | **0.139** | **0.058** | **0.505** |

*5.2. Effect of Soft-Weighted-Sum Inference*

One important contribution of this work is the proposed soft-weighted-sum inference. Here, we would like to elaborate the necessity and effectiveness of it.

Firstly, we give the probability distribution variation of randomly selected positions along the training in Fig. 6. The most interesting thing is that: In the training phase, we utilize the multinomial logistic loss, which means we don't specially discriminate the distance between the "nearby" and "further" classes. While, the probability distribution is rather clustered. More interestingly, the probability distribution roughly follow the Gaussian distribution, which means it is symmetric. At last, as the training goes on, the distribution of probability is becoming more concentrated, but always maintains symmetry similar to that of Gaussian distribution.

Secondly, we use the hard-max inference and give the confusion matrix in Fig. 7. The confusion matrix presents a kind of diagonal dominant and symmetric structure, which means most of the error prediction occurs in nearby classes.

Thirdly, we increase the number of depth bins. Under the same training setting, we present the variation of "pixel accuracy" and the relative errors in Tab. 4. With the increase of number of bins, the "pixel accuracy" drop dramatically, while the relative error keeps stable. This trend presents that: 1) At present, the network cannot distinguish the very detailed distance variation even we train it with the very detailed "label". In other words, "depth perception" ability of the network is limited.

Table 4: Pixel-wise accuracy and Rel w.r.t. number of bins.

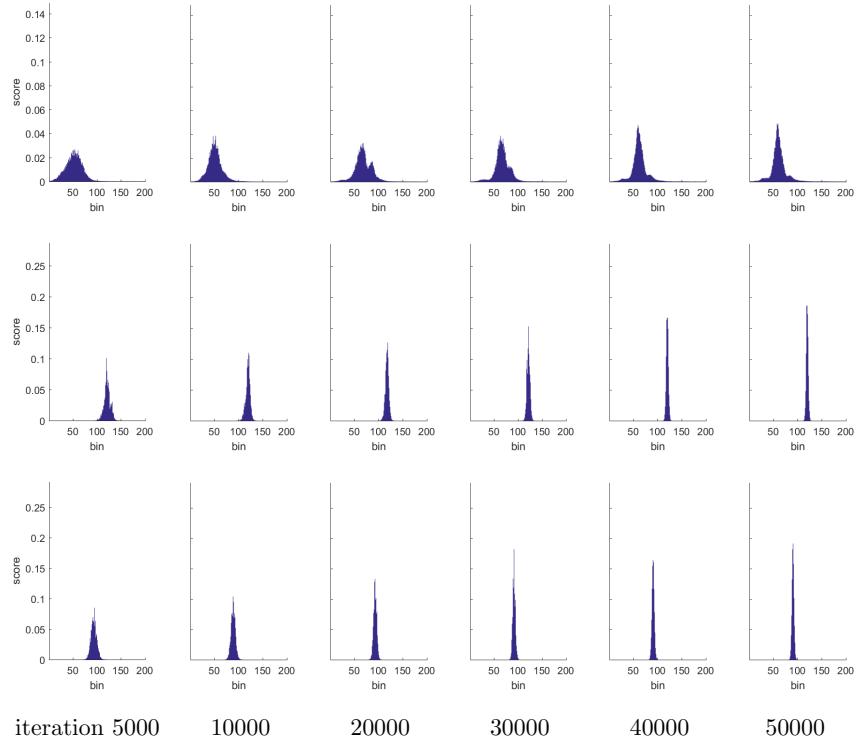| num of bins | 50 | 100 | 200 | 500 | 1000 |
|---|---|---|---|---|---|
| pixel accuracy (%) | 67 | 41 | 25 | 12 | 7 |
| Rel | 0.182 | 0.142 | 0.139 | 0.138 | 0.140 |

Figure 6: Typical score distribution variation of our network output. The points are randomly selected from NYU2 dataset.

## 6. Conclusions

In this paper, we have proposed a deep end-to-end classification based framework to monocular depth estimation. By using both dilated convolution and hierarchical fusion of multi-scale features, our framework is able to deal with the real world difficulties in multi-scale depth estimation. Extensive experiments on both indoor and outdoor benchmarking datasets show the superiority of our method compared with the current state-of-the-art methods. More importantly, experiments also demonstrate that our model is able to learn a probability distribution among different depth labels, which inspires the proposed soft-weighted-sum inference.
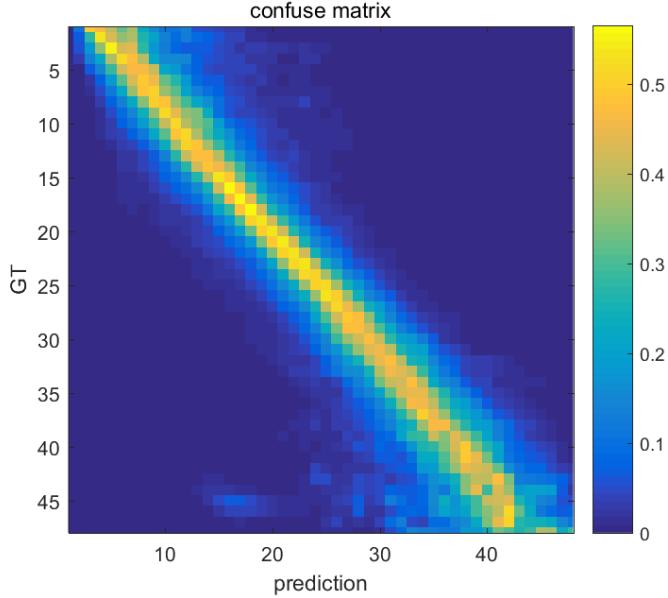
Figure 7: Confusion matrix on the NYU2 dataset. Here, we merge the 200 bins to 50 for better illustration.

## References

[1] X. Ren, L. Bo, D. Fox, RDB-D scene labeling: Features and algorithms, in: Proc. IEEE Conf. Comp. Vis. Patt. Recogn., 2012, pp. 2759–2766.

[2] J. Shotton, T. Sharp, A. Kipman, A. Fitzgibbon, M. Finocchio, A. Blake, M. Cook, R. Moore, Real-time human pose recognition in parts from single depth images, Communications of the ACM 56 (1) (2013) 116–124.

[3] S. R. Fanello, C. Keskin, S. Izadi, P. Kohli, D. Kim, D. Sweeney, A. Criminisi, J. Shotton, S. B. Kang, T. Paek, Learning to be a depth camera for close-range human capture and interaction, ACM T. Graphics 33 (4) (2014) 86.

[4] N. Silberman, D. Hoiem, P. Kohli, R. Fergus, Indoor segmentation and support inference from rgbd images, in: European Conference on Computer Vision, 2012, pp. 746–760.

[5] A. Gupta, A. Efros, M. Hebert, Blocks world revisited: Image understanding using qualitative geometry and mechanics, in: Proc. Eur. Conf. Comp. Vis., 2010, pp. 482–496.

[6] D. Fouhey, A. Gupta, M. Hebert, Unfolding an indoor origami world, in: Proc. Eur. Conf. Comp. Vis., 2014, pp. 687–702.

[7] B. Li, C. Shen, Y. Dai, A. van den Hengel, M. He, Depth and surface normal estimation from monocular images using regression on deep features and hierarchical crfs, in: Proc. IEEE Conf. Comp. Vis. Patt. Recogn., 2015, pp. 1119–1127.

[8] F. Liu, C. Shen, G. Lin, I. Reid, Learning depth from single monocular images using deep convolutional neural fields, TPAMI 38 (10) (2016) 2024–2039.

[9] P. Wang, X. Shen, Z. Lin, S. Cohen, B. Price, A. L. Yuille, Towards unified depth and semantic prediction from a single image, in: Proc. IEEE Conf. Comp. Vis. Patt. Recogn., 2015, pp. 2800–2809.

[10] D. Eigen, R. Fergus, Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture, in: ICCV, 2015, pp. 2650–2658.

[11] Y. Cao, Z. Wu, C. Shen, Estimating depth from monocular images as classification using deep fully convolutional residual networks, [Online]. Avaliable: https://arxiv.org/abs/1605.02305.

[12] I. Laina, C. Rupprecht, V. Belagiannis, F. Tombari, N. Navab, Deeper depth prediction with fully convolutional residual networks, in: 3D Vision (3DV), 2016 Fourth International Conference on, IEEE, 2016, pp. 239–248.

[13] D. Xu, E. Ricci, W. Ouyang, X. Wang, N. Sebe, Multi-scale continuous crfs as sequential deep networks for monocular depth estimation.

[14] C. Godard, O. M. Aodha, G. J. Brostow, Unsupervised monocular depth estimation with left-right consistency.

[15] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: CVPR, 2016, pp. 770–778.

[16] A. Saxena, M. Sun, A. Y. Ng, Make3d: Learning 3d scene structure from a single still image, IEEE Trans. Pattern Anal. Mach. Intell. 31 (5) (2009) 824–840.

[17] A. Saxena, J. Schulte, A. Y. Ng, Depth estimation using monocular and stereo cues, in: Proc. IEEE Int. Joint Conf. Artificial Intell., Vol. 7, 2007.

[18] B. Liu, S. Gould, D. Koller, Single image depth estimation from predicted semantic labels, in: Proc. IEEE Conf. Comp. Vis. Patt. Recogn., 2010, pp. 1253–1260.

[19] L. Ladicky, J. Shi, M. Pollefeys, Pulling things out of perspective, in: Proc. IEEE Conf. Comp. Vis. Patt. Recogn., IEEE, 2014, pp. 89–96.

[20] K. Karsch, C. Liu, S. B. Kang, Depth extraction from video using non-parametric sampling, in: Proc. Eur. Conf. Comp. Vis., Springer, 2012, pp. 775–788.

[21] M. Liu, M. Salzmann, X. He, Discrete-continuous depth estimation from a single image, in: Proc. IEEE Conf. Comp. Vis. Patt. Recogn., 2014, pp. 716–723.

[22] J. Konrad, M. Wang, P. Ishwar, 2d-to-3d image conversion by learning depth from examples, in: Proc. IEEE Conf. Computer Vis. & Pattern Recogn. Workshops, IEEE, 2012, pp. 16–22.

[23] D. Eigen, C. Puhrsch, R. Fergus, Depth map prediction from a single image using a multi-scale deep network, in: Proc. Adv. Neural Inf. Process. Syst., 2014.

[24] W. Chen, Z. Fu, D. Yang, J. Deng, Single-image depth perception in the wild.

[25] R. Garg, K. B. G. Vijay, G. Carneiro, I. Reid, Unsupervised cnn for single view depth estimation: Geometry to the rescue.

[26] T. Zhou, M. Brown, N. Snavely, D. G. Lowe, Unsupervised learning of depth and ego-motion from video, in: Proc. IEEE Conf. Comp. Vis. Patt. Recogn., 2017.

[27] Y. Kuznietsov, J. Stückler, B. Leibe, Semi-supervised deep learning for monocular depth map prediction, in: Proc. IEEE Conf. Comp. Vis. Patt. Recogn., 2017.
URL http://arxiv.org/abs/1702.02706

[28] B. Ummenhofer, H. Zhou, J. Uhrig, N. Mayer, E. Ilg, A. Dosovitskiy, T. Brox, Demon: Depth and motion network for learning monocular stereo, in: Proc. IEEE Conf. Comp. Vis. Patt. Recogn., 2017.
URL http://lmb.informatik.uni-freiburg.de//Publications/2017/UZUMIDB17

[29] E. Shelhamer, J. Long, T. Darrell, Fully convolutional networks for semantic segmentation, IEEE Trans. Pattern Anal. Mach. Intell. 39 (4) (2017) 640–651. doi:10.1109/TPAMI.2016.2572683.

[30] B. Hariharan, P. Arbelaez, R. Girshick, J. Malik, Hypercolumns for object segmentation and fine-grained localization, in: Proc. IEEE Conf. Comp. Vis. Patt. Recogn., 2015, pp. 447–456.

[31] F. Yu, V. Koltun, Multi-scale context aggregation by dilated convolutions, in: ICLR, 2016, pp. 1–10.

[32] L. C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, A. L. Yuille, Semantic image segmentation with deep convolutional nets and fully connected crfs, Computer Science (4) (2014) 357–361.

[33] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, in: Proc. Int. Conf. Learning Representations, 2015.

[34] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, T. Darrell, Caffe: Convolutional architecture for fast feature embedding, in: Proc. ACM Int. Conf. Multimedia, 2014, pp. 675–678.

[35] A. Geiger, P. Lenz, R. Urtasun, Are we ready for autonomous driving? the kitti vision benchmark suite, in: Proc. IEEE Conf. Comp. Vis. Patt. Recogn., 2012, pp. 3354–3361.