# DEPT: Depth Estimation by Parameter Transfer with a Lightweight Model for Single Still Images

Hongwei Qin, Student Member, IEEE, Xiu Li, Member, IEEE, Yangang Wang, Yongbing Zhang, Member, IEEE, and Qionghai Dai Senior Member, IEEE

*Abstract*—In this paper, we propose a novel method for automatic depth estimation from color images using parameter transfer. By modeling the correlation between color images and their depth maps with a set of parameters, we get a database of parameter sets. Given an input image, we extract the high-level features to find the best matched image sets from the database. Then the set of parameters corresponding to the best match are used to estimate the depth of the input image. Compared to the past learning-based methods, our trained model only consists of trained features and parameter sets, which occupy little space. We evaluate our depth estimation method on several benchmark RGB-D (RGB + depth) datasets. The experimental results are comparable to the state-of-the-art, while the model size is very small and very suitable for mobile devices, demonstrating the promising performance of our proposed method.

*Index Terms*—depth estimation, parameter transfer, 3D reconstruction.

## I. INTRODUCTION

**MAGES** captured with conventional cameras lose the depth information of the scene. However, scene depth is of great importance for many computer vision tasks. 3D applications like 3D reconstruction for scenes (*e.g.*, Street View on Google Map), robot navigation, 3D videos, and free view video(FVV) [1], [2] all rely on scene depth. Depth information can also be useful for 2D applications like image enhancing [3] and scene recognition [4]. Recent RGB-D imaging devices like Kinect are greatly limited on the perceptive range and depth resolution. Neither can they extract depth for the existing 2D images. Therefore, depth estimation from color images has been a useful research subject.

In this paper, we propose a novel depth estimation method to generate depth maps from single still images. Our method applies to arbitrary color images. We build the connection between image and depth with a set of parameters. A parameter sets database is constructed, and the parameter sets are transferred to input images to get the corresponding depth maps. Some estimation results are shown in Fig. 1.

As a reminder, the paper is organized as follows. In Section II, the related techniques are surveyed. In Section III, we introduce our proposed DEPT (depth estimation by parameter transfer) method in details. We demonstrate our method on the RGB-D benchmark datasets in Section IV. Finally, we conclude our work in Section V.

#### II. RELATED WORKS

In this section, we introduce the techniques related to this paper, which are respectively depth estimation from a single image, and parameter transfer.

#### A. Depth Estimation from Single Images

The reason why depth estimation from a single image is possible lies in that there are some monocular depth cues in a 2D image. Some of these cues are inferred from local properties like color, shading, haze, defocus, texture variations and gradients, occlusions and so on. Global cues are also crucial to inferring depth, as the ability humans have. So, integrating local and global cues of a single image to estimate depth is reasonable.

There are semi-automatic and automatic methods for depth estimation from single images. Horry *et al.* [5] propose *tour into the picture*, where the user interactively adds planes to an image to make animation. The work of Zhang *et al.* [6] requires the user to add constrains manually to images to estimate depth.

Automatic methods for single image depth estimation come up in recent years. Hoiem et al. [7] propose automatic photo pop-up, which reconstructs an outdoor image using assumed planar surfaces of it. Delage et al. [8] develop a Bayesian framework applied to indoor scenes. Saxena et al. [9] propose a supervised learning approach, using a discriminativelytrained Markov Random Field (MRF) that incorporates multiscale local and global image features. Then, they improve this method in [10]. After that, depth estimation from predicted semantic labels is proposed by Liu et al. [11]. A more sophisticated model called Feedback Enabled Cascaded Classification Models (FE-CCM) is proposed by Li et al. [12]. One typical depth estimation method is Depth Transfer, developed by Karsch et al. [13]. This method first builds a large scale RGB-D images and features database, then acquires the depth of the input image by transferring the depth of several similar images after warping and optimizing procedures.

There are several recent works that try to solve the depth estimation problem and semantic segmentation problem unitedly. Ladicky *et al.* [14] propose to predict pixel-wise semantic class labels to improve both depth estimation and semantic segmentation performance. Eigen *et al.* [15], [16] use a multi-scale convolutional architecture to refine local depth prediction with global information. Wang *et al.* [17] propose to decompose the image into local segments for region-level

H. Qin, X. Li, Y. Zhang and Q. Dai are with the Department of Automation, Tsinghua University, Beijing, 100084 China e-mail: li.xiu@sz.tsinghua.edu.cn H. Qin, X. Li and Y. Zhang are with Graduate School at Shenzhen, Tsinghua University, Shenzhen, 518055 China

Y. Wang is with Microsoft Research Asia.



(b) Estimated depth maps by DEPT

Fig. 1. Selected images and corresponding depth maps estimated by DEPT. The darker the red is, the further (from the imaging device) the objects are. The darker the blue is, the closer the objects are.

depth and semantic prediction under the guidance of global layout.

Besides, there are several other efforts on depth estimation with unified global and local information. Liu et al. [18] use continuous variables encoding the depth of the superpixels in the input image and discrete variables representing relationships between neighboring superpixels to perform inference through a graphical model. Zhuo et al. [19] propose to use a hierarchical representation of the indoor scene, and refine the depth map guided by global layout. Liu and Shen et al. [20] propose a method to refine depth map predicted by convolutional networks by continuous conditional random field. In the work by Baig et al. [21], [22], they express the global depth map of an image as a linear combination of a depth basis learned from examples. The basis is actually a dictionary of the training dataset, and the images near the cluster centroids are picked as basis elements. Our concurrent and independent work also use cluster controids but with a totally different way, which we will introduce in detail.

Under specific conditions, there are other depth extract methods, such as dark channel prior proposed by He *et al.* [23], proved effective for hazed images.

The method closest to ours is the parametric model developed by Wang *et al.* [24] for describing the correlation between single color images and depth maps. This work treats the color image as a set of patches and derives the correlation with a kernel function in a non-linear mapping space. They get convincing depth map through patch sampling. However, this work only demonstrates the effectiveness of the model, and can't estimate depth with an arbitrary input image. Our improvements are two-fold: we extend this model from one image to many, and we transfer parameter set to an arbitrary input image according to best image set match.

## B. Parameter Transfer

We carry out a survey on transfer methods in the field of depth estimation. The non-parametric scene parsing by Liu *et al.* [25] avoids explicitly defining a parametric model and scales better with respect to the training data size. The Depth Transfer method by Karsch *et al.* [13] leverages this work and assumes that scenes with similar semantics should have similar depth distributions after densely aligned. Their method contains three stages. First, given an input image, they find K best matched images in RGB space. Then, the K images are warped to be densely aligned with the input. Finally, they use an optimization scheme to interpolate and smooth the warped depth values to get the depth of the input.

Our work is different in three aspects. First, instead of depth, we transfer parameter set to the input image, so we don't need post process like warping. Second, our database is composed of parameter sets instead of RGB-D images, so the database occupies little space. Third, the depth values are computed with the transferred parameter set directly, so we don't need an optimization procedure after transfer.

## III. DEPT: DEPTH ESTIMATION BY PARAMETER TRANSFER

In this section, we first introduce the modeling procedure for inferring the correlation between color images and depth maps. Then, we introduce the parameter transfer method in detail.

## A. The Parametric Model

The prior work of Wang *et al.* [24] proposed a model to build the correlation between a single image I and its corresponding depth map D with a set of parameters. We extend this by using a set of similar images IS and their corresponding depth map DS. So the parameters contain information of all the images in the set.

We regard each color image as a set of overlapped fixedsize color patches, of which the size will be discussed later. For each image, we sample the patches  $x_1, x_2, ..., x_p$  and their corresponding depth values from RGB-D image set. To avoid over-fitting, we only sample p patches from each image. In our experiment, we set p as 1000, and the samples account for 0.026% of the total patches in one image. We use a uniform sampling method, *i.e.*, we separate the image into grids and select samples uniformly from all the grids. By denoting Nas the number of images in an image set, totally we sample  $N \times p$  patches. Specially, for single image, N = 1.

1) Modeling the Correlation between Image and Depth: After the sampling procedure, we model the correlation by measuring the sum squared error between the depth  $\hat{d}$  mapped with the sampled color patches and the ground truth depth d. The model is written as

$$E = \sum_{i=1}^{p \times N} |tr(W^T \sum_{j=1}^n \gamma_j \phi(x_i * f_j)) - d_i|^2 , \qquad (1)$$

where E is the sum squared estimation error, p is the number of sample patches per image, N is the number of images in the image set,  $f_j$  is the filters, n is the number of filters and set as 9 in all the experiments. If set larger, the algorithm is expected to get better results, but bring larger cost.  $\phi$  is the kernel function to map the convolved patches and sum them up to *one patch*,  $\gamma_j$  is the weight of each convolved patch, W is the weight matrix, whose size is the same of the *one patch*, aiming at integrating the overall information from each patch.

Eq. 1 can be rewritten as

$$E = \sum_{i=1}^{p \times N} |\mathbf{w}^T \phi(X_i F) \gamma - d_i|^2 , \qquad (2)$$

where  $X_i$  is a matrix reshaped from patch  $x_i$ . The row size of  $X_i$  is the same as  $f_i$ , while  $F = [f_1, f_2, ..., f_n]$ ,  $\gamma = [\gamma_1, \gamma_2, ..., \gamma_n]^T$ . w is the result of concatenating all the entries of W.

At the image level, F describes the texture gradient cues of the RGB image by extracting the frequency information.  $\gamma$ describes the variance of filters. We use Principle Component Analysis (PCA) to initialize F, and optimize it afterwards. As for the size of filter, we need to balance between efficiency 2) Estimating Model Parameters: First, we rewrite Eq. 2 as

$$E = \|M\phi(XF)\gamma - \mathbf{d}\|_2^2 , \qquad (3)$$

and

$$E = \|\mathbf{\Gamma}\phi(F^T\hat{X})\mathbf{w} - \mathbf{d}\|_2^2 , \qquad (4)$$

where X is got by concatenating all the  $X_i$  in Eq. 2.  $\hat{X}$  is got by concatenating all the  $X_i^T$ . Each row of M is  $\mathbf{w}^T$ , and each row of  $\Gamma$  is  $\gamma^T$ . So Eq. 3 is a least square problem of  $\gamma$ , and Eq. 4 is a least square problem of  $\mathbf{w}$ . Then we minimize E by optimizing the filters F. Finally we get a set of parameters, consisting of F,  $\gamma$ , and  $\mathbf{w}$ . The detailed method for solving this can be found in our previous work [24].

#### B. Parameter Transfer

Our parameter transfer procedure, outlined in Fig. 2, has three stages. First, we build a parameter set database using training RGB-D images. Second, given an input image, we find the most similar image sets using high-level image features, and transfer the parameter set to the input image. Third, we compute the depth of the input image.

1) Parameter Set Database Building: Given a RGB-D training dataset, we compute high-level image features for each image. Here, we use GIST [26] features, which can be used to measure similarities of images. Then, we categorize the training images to N sets, using k-means cluster method. Next, we get the central GIST feature for each image set. For each image set, the corresponding parameter set is obtained using our parameter estimate model. The central GIST features and corresponding parameter sets compose our parameter set database. Actually, this database is so small as to occupy much less space compared to the RGB-D datasets.

2) Image Set Matching: Given an input image, we compute its GIST feature and find the best matched central GIST feature from our trained database. Then the parameter set corresponding to the best matched central GIST feature (*i.e.* the central GIST feature of the most similar image set) is transferred to the input image. We define the best match as

$$G_{best} = \min_{i=1,2,...,N} \|G_{input} - G_i\| , \qquad (5)$$

where  $G_{input}$  denotes the GIST feature of the input image, and  $G_i$  denotes the central GIST feature of each image set.

As the most similar image set matches the input closely in feature space, the overall semantics of the scenes are similar. At the low level, the cues such as the texture gradient, texture variation, and color are expected to be roughly similar to some extent. With the model above, the parameters connecting the images and depth maps should be similar. So, it is reasonable to transfer the parameter set to the input image.

#### JOURNAL OF ${\rm LATE} X$ CLASS FILES, VOL. XX, NO. XX, XX XXXX



Fig. 2. Our pipeline for estimating depth. First we build a parameter set database, then the parameter set is transferred to the input image according to the best matched GIST feature. Finally, the parameter set is used to estimate the depth.

*3) Depth Estimation:* We use the color patches of the input image and the transferred parameter set to map the estimation depth. The computational formula is:

$$\hat{\mathbf{d}} = M\phi(XF)\gamma , \qquad (6)$$

where X is the patches, F is the filters.  $\gamma$  is the weight to balance the filters. M is the weight matrix. These parameters are all from the parameter set.

#### IV. EXPERIMENT

In this section, we evaluate the effectiveness of our DEPT method on single image RGB-D datasets.

## A. RGB-D Datasets

We use the Make3D Range Image Dataset [27]. The dataset is collected using 3D scanner and the corresponding depth maps using lasers. There are totally 534 images separated into two parts, which are the training part containing 400 images and the testing part containing 134 images, respectively. The color image resolution is  $2272 \times 1704$ , and the ground truth depth map resolution is  $55 \times 305$ . Befor training, we resize the depth map resolution to the same size of the color image, so RGB and D (Depth) have pixel-wise correspondence.

## B. Image Cluster

We compute the GIST features for each image in the training dataset. Then we use k-means algorithm to cluster the images into N sets, here we set N as 30. The images are well separated according to the scene semantics. The silhouette plot in Fig. 3 measures how well-separated the resulting image sets are. Lines on the right side of 0 measure how distant that image is from neighboring image sets. Lines on the left of



Fig. 3. Silhouette plot of the k-means cluster result. Each line represents an image. Lines on the right side of 0 measure how distant that image is from neighboring image sets. Lines on the left of 0 indicate that image is probably assigned to the wrong set. The vertical axis indicates different clusters (image sets).

0 indicate that image is probably assigned to the wrong set. The vertical axis indicates different clusters (image sets). As we can see, most of the images are well clustered. As for the choosing of N, initially we choose it by observing the silhouette plot, then we try a series of values with a step of 10. The results around 30 are close, and 30 is the best. The cluster number can also be set according to existing pattern classification methods (*e.g.* methods to find best k in k-means algorithm [28], [29]). We believe N should not be too large or too small. Too large N may set similar scenes apart while too small N may result in large scene variety in one set.

An example image set is shown in Fig. 5. It can be seen



(a) One clustered image set



(b) The corresponding depth maps

Fig. 5. One example image set after image cluster procedure. (a) is a clustered image set, containing 18 semantic similar images, (b) are their corresponding depth maps. The depth distributions in the images are roughly similar.



Fig. 4. Energy decline curves of the 30 image sets. E is on a ln scale.

that the clustered images have roughly similar semantic scene. The depth distributions also seem similar, as are shown in the color images as well as the depth maps.

## C. Parameter Sets Estimation

For each image set, we estimate the corresponding model parameters. The overlapped patch size is set  $15 \times 15$ . The filter size is set as  $3 \times 3$ . We separate each image into grids and uniformly sample 1000 patches per image. So for an N sized image set, totally  $1000 \times N$  patches are sampled, which occupy 0.026% of the whole image set. We initialize the filters with PCA method, and optimize all the parameters using warmstart gradient descent method. The iteration stop condition is  $E < 10^{-6}$ . In our experiment, the energy (i.e., the sum squared errors E) declines as Fig. 4 shows. As can be seen, most of the curves come to a steady state after about 1000 iterations. The smaller the steady energy is, the more similar the images in that set are.

For each image set, we obtain one optimized parameter set. The 30 parameter sets and the corresponding cluster centroids (the center of the GIST features in each image set) make up the parameter sets database.

#### D. Depth Estimation by Parameter Transfer

For each of the testing 134 images, we find the best matched image set from the parameter sets database and compute the depth maps using the computational formula of Eq. 6.

1) Quantitative Comparison with Previous Methods: We calculate three common error metrics for the estimated depth. Denoting  $\hat{\mathbf{D}}$  as the estimated depth and  $\mathbf{D}$  as the ground truth depth, we calculate **RE** (*relative error*):

$$\mathbf{RE} = \frac{|\hat{\mathbf{D}} - \mathbf{D}|}{\mathbf{D}} , \qquad (7)$$

**LE** ( $\log_{10} error$ ):

$$\mathbf{LE} = |\log_{10}(\mathbf{D}) - \log_{10}(\mathbf{D})| , \qquad (8)$$

and **RMSE** (root mean squared error):

$$\mathbf{RMSE} = \sqrt{\sum_{i=1}^{P} (\hat{\mathbf{D}}_i - \mathbf{D}_i)^2 / P} , \qquad (9)$$

TABLE I Average error and database size comparison of various estimate methods.

Method	RE	LE	RMSE	Trained Database
Depth MRF [9]	0.530	0.198	16.7	-
Make3D [27]	0.370	0.187	-	-
Feedback Cascades [12]	-	-	15.2	-
Deep CNN Fields [20]	0.314	0.119	8.60	140 MB
Depth Transfer [13]	0.361	0.148	15.1	2.44 GB
DEPT with GIST(ours)	0.489	0.182	16.9	1.47 MB
DEPT with CNN(ours)	0.421	0.172	16.7	1.25 MB

where P is the pixel number of a depth map.

Error measure for each image is the average value of all the pixels on the ground truth resolution scale ( $55 \times 305$ ). Then the measures are averaged over all the 134 images to get final error metrics, which are listed in Table I.

As can be seen, our results are better than Depth MRF [9] in view of RE and LE, better than Make3D [27] in view of LE. Totally speaking, the results of DEPT are comparable with the state-of-the-art learning based automatic methods. Especially, DEPT only requires a very small sized database, and once the database is built, we can compute the depth directly. Built from the 400 training RGB-D images that occupy 628MB space, our database size is only 188KB (0.03%). As a contrast, the trained database of Depth Transfer [13] occupies 2.44GB<sup>1</sup> (about 4 times of the original dataset size). Though our method has disadvantage in average errors over the Depth Transfer [13], we have large advantages in database space consuming and computer performance requirement(in [13], the authors claim Depth Transfer requires a great deal of data (GB scale) to be stored concurrently in memory in the optimization procedure), which are especially crucial when the database grows in real applications. Recent deep CNN based depth estimation methods get lower errors. Essentially, our convolutional operation and optimization method is similar to convolutional neural network with only one layer. From this point of view, our method uses much less parameters with good results, if implemented on high-end GPU, as deep CNNs are, our method would gain much higher efficiency.

Further more, our method also has advantages in some of the estimation effects, as is detailed in the following qualitative evaluation.

2) Qualitative Evaluation: A qualitative comparison of our estimated depth maps, depth maps estimated by Depth Transfer [13] and the ground truth depth maps are demonstrated in Fig. 6 and Fig. 7. As can be seen, our estimated depth maps are visually reasonable and convincing, especially in the details like texture variations (e.g., the tree in the second column of Fig. 6) and relative depth (e.g., the pillars' depth in the last column of Fig. 6 is well estimated by our DEPT method, while Depth Transfer [13] estimates wrong). Actually, some of our results are even more accurate than the ground truth (e.g., in the third column in Fig. 7, there is a large part of wrong

<sup>&</sup>lt;sup>1</sup>Implemented with the authors' public codes at http://research.microsoft.com/en-us/downloads/29d28301-1079-4435-9810-74709376bce1/



(a) Test images



(b) Ground truth depth maps



(c) Estimated depth maps by DEPT (our method)



(d) Estimated depth maps by Depth Transfer [13]

Fig. 6. Performance comparison: scenes of streets, squares and trees. (a) show some test images containing streets, squares or trees, (b) are corresponding ground truth depth maps, (c) are estimated depth maps by DEPT (our method), (d) are estimated depth maps by Depth Transfer [13]



(a) Test images



(b) Ground truth depth maps



(c) Estimated depth maps by DEPT (our method)



(d) Estimated depth maps by Depth Transfer [13]

Fig. 7. Performance comparison: scenes of buildings. (a) show some test images containing buildings, (b) are corresponding ground truth depth maps, (c) are estimated depth maps by DEPT (our method), (d) are estimated depth maps by Depth Transfer [13]



(d) Estimated depth maps by Depth Transfer [13]

Fig. 8. Performance comparison: indoor scenes. (a) show some test images containing buildings, (b) are corresponding ground truth depth maps, (c) are estimated depth maps by DEPT (our method), (d) are estimated depth maps by Depth Transfer [13]

depth in the building area of the ground truth depth map). The ground truth maps have some scattered noises, which may result from the capturing device. While the noises in our depth maps are less because of the using of overall information in the image set. But we must point out that the sky areas in our depth maps are not as pleasing, which may result from the variation of sky color and texture among various images in a set, especially when the cluster result is biased. This may result in the increase of average error in the previous metrics. However, as the increasing of RGB-D images acquired by depth imaging devices, our database can expand easily due to the extremely small space consuming, which means we may get more and more accurate matched parameter sets for existing RGB images and video frames.

#### E. Evaluation on Indoor Datasets

We also implement an experiment on the NYU Depth V2 Dataset [30], which consists of 1449 indoor RGB-D images captured with Kinect. We use the *Labeled Dataset*<sup>2</sup>, i.e., 1449 densely labeled pairs of aligned RGB and depth images.

The dataset is partitioned into 795 training images and 654 testing images. When training DEPT, we cluster the training dataset to 80 sets, guided by k-means silhouette plot and linear search. One example of the cluster set is shown in Figure 9. Quantitative results are shown in Table. II. In addition to the three standard metrics, we also report the metrics used in [14], defined as

$$\frac{1}{N} \sum_{p=1}^{N} \left[ \left[ max(\frac{d_p}{g_p}, \frac{g_p}{d_p}) = \delta < t \right] \right] \times 100\%,$$
(10)

where  $g_p$  is the ground-truth of pixel p,  $d_p$  is the corresponding estimated depth, N is the number of pixels,  $t = 1.25, 1.25^2, 1.25^3$  is the threshold, and [[·]] denotes the indicator function. We can observe that DEPT achieves comparable quantitative results with much less space and time consumption.

Qualitative results are shown in Figure. 8. We can see DEPT gets not so smooth results (we did not use smoothing operation as Depth Transfer did), but infers more details on the edges. This may be useful when an application cares more about edges of depth map.

<sup>&</sup>lt;sup>2</sup>http://cs.nyu.edu/~silberman/datasets/nyu\_depth\_v2.html



Fig. 9. One example image set after image cluster procedure on NYU. The clustered image set contains 9 semantic similar images.

In addition, the testing procedure (654 images) consumes about 4 hours with DEPT on our computer (Intel Xeon E3-1330 V2 CPU, 16GB RAM, 64bit Windows 7, without any algorithm optimization), while it takes about 45 hours with Depth Transfer [13].

## F. Replace GIST With Deep CNN

Following prior work [31], [32], we observe that convolutional neural networks have good scene descriptions for images. Thus, we follow the method of [32] to compute the representations of the RGB images. The CNN feature extraction process is illustrated in Fig. 10. For each of the training images, the representation is computed as follows:

$$v = CNN_{\theta_c}(I), \tag{11}$$

where  $CNN_{\theta_c}(I)$  transforms the pixels of image I into a 4096-dimensional activation of the fully connected layer immediately before the classifier, i.e., the 1000-way softmax layer. The CNN parameters  $\theta_c$  contain approximately 60 million parameters and the architecture closely follows the network of Krizhevsky *et al.* [33], but we chop off the final 1000-way softmax layer. In this way, after network forwarding, each image is represented as a 4096-dimensional vector.

This vector can be treated as features of the image, which will be referred to below as CNN features. We replace GIST features with CNN features in the previous framework. We carry out experiments with the new framework. The result is listed in Table. I and Table. II. We observe decrease of all the error indicators. This performance is better than the originally proposed method in the conference version of this work [34]. In the mean time, though CNN features can improve the performance, it increases time consuming and model size. Because for now, most of the CNN based methods rely on high performance GPUs. It is too slow on personal computers, not to mention mobile devices. So, we need to balance the performance, speed and model size in real applications. However, we can expect more improvement of DEPT when better algorithms for semantic scene matching are proposed.

#### V. CONCLUSION AND FUTURE WORKS

In this paper, we propose a lightweight, effective and fully automatic technique to restore depth information from single still images. Our depth estimation by parameter transfer (DEP-T) method is novel in that we use clustered scene semantics similar image sets to model the correlation between RGB information and D (depth) information, obtaining a database of parameter sets and cluster centers. DEPT only requires the trained parameter sets database which occupies much less space compared with previous learning based methods. Experiments on RGB-D benchmark datasets show quantitatively comparable to the state-of-the-art and qualitatively good

Method	RE	LE	RMSE	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$	Database	Time
Depth Transfer [13]	0.374	0.134	1.12	49.81%	79.46%	93.75%	1.14 GB	45 hours
DC-Depth [18]	0.335	0.127	1.06	51.55%	82.32%	95.00%	-	-
Zhuo et. al [19]	0.305	0.122	1.04	52.50%	83.77%	96.16%	-	-
DEPT with GIST (ours)	0.392	0.151	1.19	48.50%	79.02%	93.12%	465 KB	4 hours
DEPT with CNN (ours)	0 353	0.130	1 1 1	51 24%	80.62%	94 35%	460 KB	4 hours

 TABLE II

 EXPERIMENTAL RESULTS ON NYU INDOOR DATASET.



Fig. 10. An illustration of the CNN feature extraction architecture. A  $224 \times 224$  crop of an image (RGB) is presented as the input. It is convolved with 96 different filters, each of size  $7 \times 7$ , using a stride of 2 in both x and y. The resulting feature maps are then: (i) passed through a rectified linear function (not shown), (ii) pooled (max within  $3 \times 3$  regions, using stride 2) and (iii)contrast normalized across feature maps to give 96 different  $55 \times 55$  element feature maps. Similar operations are repeated in layers 2,3,4,5. The last two layers are fully connected, taking features from the top convolutional layer as input in vector form ( $6 \times 6 \times 256 = 9216$  dimensions). The output of layer 7 are our CNN features in vector form(4096 dimensions). The final layer is a 1000-way softmax function, whose output is 1 predicted class out of 1000.

results. The estimated depth maps are visually reasonable and convincing, especially in the details like texture variations and relative depth. Further more, as the increasing of RGB-D images acquired by depth imaging devices, our database can expand easily due to the extremely small space consuming. As our model is only about one MB, it is very suitable to use on mobile devices (The code will be released upon publication). In the future work, we would like to improve the cluster accuracy by exploring more accurate similarity metrics that are applicable to our image and depth correlation model. We plan to build a larger RGB-D image dataset as more data brings better performance with our method. Finally, we suppose it is also meaningful to improve the depth estimation performance for video frames by using optical flow features or other features related to time coherence.

#### ACKNOWLEDGMENT

This work is supported by National Natural Science Foundation of China (Grant No. 71171121/61033005) and National 863 High Technology Research and Development Program of China (Grant No. 2012AA09A408).

#### REFERENCES

- Q. Liu, Y. Yang, R. Ji, Y. Gao, and L. Yu, "Cross-view down/upsampling method for multiview depth video coding," *Signal Processing Letters, IEEE*, vol. 19, no. 5, pp. 295–298, 2012.
- [2] Y. Liu, Q. Dai, and W. Xu, "A point-cloud-based multiview stereo algorithm for free-viewpoint video." *IEEE Transactions on Visualization* and Computer Graphics, vol. 16, no. 3, pp. 407–418, 2010.

- [3] F. Li, J. Yu, and J. Chai, "A hybrid camera for motion deblurring and depth map super-resolution," in *Computer Vision and Pattern Recognition*, 2008. CVPR 2008. IEEE Conference on. IEEE, 2008, pp. 1–8.
- [4] A. Torralba and A. Oliva, "Depth estimation from image structure," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 24, no. 9, pp. 1226–1238, 2002.
- [5] Y. Horry, K.-I. Anjyo, and K. Arai, "Tour into the picture: using a spidery mesh interface to make animation from a single image," in *Proceedings* of the 24th annual conference on Computer graphics and interactive techniques. ACM Press/Addison-Wesley Publishing Co., 1997, pp. 225–232.
- [6] L. Zhang, G. Dugas-Phocion, J.-S. Samson, and S. M. Seitz, "Singleview modelling of free-form scenes," *The Journal of Visualization and Computer Animation*, vol. 13, no. 4, pp. 225–235, 2002.
- [7] D. Hoiem, A. A. Efros, and M. Hebert, "Automatic photo pop-up," in ACM Transactions on Graphics (TOG), vol. 24, no. 3. ACM, 2005, pp. 577–584.
- [8] E. Delage, H. Lee, and A. Y. Ng, "A dynamic bayesian network model for autonomous 3d reconstruction from a single indoor image," in *Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on*, vol. 2. IEEE, 2006, pp. 2418–2428.
- [9] A. Saxena, S. H. Chung, and A. Y. Ng, "Learning depth from single monocular images," in Advances in Neural Information Processing Systems, 2005, pp. 1161–1168.
- [10] —, "3d depth reconstruction from a single still image," *International Journal of Computer Vision*, vol. 76, no. 1, pp. 53–69, 2008.
- [11] 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. IEEE, 2010, pp. 1253–1260.
- [12] C. Li, A. Kowdle, A. Saxena, and T. Chen, "Towards holistic scene understanding: Feedback enabled cascaded classification models," in *Advances in Neural Information Processing Systems*, 2010, pp. 1351– 1359.
- [13] K. Karsch, C. Liu, and S. B. Kang, "Depthtransfer: Depth extraction from video using non-parametric sampling," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 2014.
- [14] L. Ladicky, J. Shi, and M. Pollefeys, "Pulling things out of perspective," in Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on. IEEE, 2014, pp. 89–96.

- [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*, 2014, pp. 2366–2374.
  [16] D. Eigen and R. Fergus, "Predicting depth, surface normals and semantic
- [16] D. Eigen and R. Fergus, "Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture," *arXiv* preprint arXiv:1411.4734, 2014.
- [17] P. Wang, X. Shen, Z. Lin, S. Cohen, B. Price, and A. L. Yuille, "Towards unified depth and semantic prediction from a single image," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 2800–2809.
- [18] M. Liu, M. Salzmann, and X. He, "Discrete-continuous depth estimation from a single image," in *Computer Vision and Pattern Recognition* (CVPR), 2014 IEEE Conference on. IEEE, 2014, pp. 716–723.
- [19] W. Zhuo, M. Salzmann, X. He, and M. Liu, "Indoor scene structure analysis for single image depth estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 614– 622.
- [20] F. Liu, C. Shen, and G. Lin, "Deep convolutional neural fields for depth estimation from a single image," in *Proceedings of the IEEE Conference* on Computer Vision and Pattern Recognition, 2015.
- [21] M. H. Baig, V. Jagadeesh, R. Piramuthu, A. Bhardwaj, W. Di, and N. Sundaresan, "Im2depth: Scalable exemplar based depth transfer," in *Applications of Computer Vision (WACV), 2014 IEEE Winter Conference* on. IEEE, 2014, pp. 145–152.
- [22] M. H. Baig and L. Torresani, "Coarse-to-fine depth estimation from a single image via coupled regression and dictionary learning," *arXiv* preprint arXiv:1501.04537, 2015.
- [23] K. He, J. Sun, and X. Tang, "Single image haze removal using dark channel prior," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 33, no. 12, pp. 2341–2353, 2011.
- [24] Y. Wang, R. Wang, and Q. Dai, "A parametric model for describing the correlation between single color images and depth maps," *Signal Processing Letters*, 2014.
- [25] C. Liu, J. Yuen, and A. Torralba, "Nonparametric scene parsing via label transfer," *Pattern Analysis and Machine Intelligence, IEEE Transactions* on, vol. 33, no. 12, pp. 2368–2382, 2011.
- [26] A. Oliva and A. Torralba, "Modeling the shape of the scene: A holistic representation of the spatial envelope," *International journal of computer vision*, vol. 42, no. 3, pp. 145–175, 2001.
- [27] A. Saxena, M. Sun, and A. Y. Ng, "Make3d: Learning 3d scene structure from a single still image," *Pattern Analysis and Machine Intelligence*, *IEEE Transactions on*, vol. 31, no. 5, pp. 824–840, 2009.
- [28] A. Rodriguez and A. Laio, "Clustering by fast search and find of density peaks," *Science*, vol. 344, no. 6191, pp. 1492–1496, 2014.
- [29] G. Hamerly and C. Elkan, "Learning the k in k-means," Advances in neural information processing systems, vol. 16, p. 281, 2004.
- [30] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus, "Indoor segmentation and support inference from rgbd images," in *Computer Vision–ECCV* 2012. Springer, 2012, pp. 746–760.
- [31] K. Kang and X. Wang, "Fully convolutional neural networks for crowd segmentation," *Eprint Arxiv*, 2014.
- [32] A. Karpathy and L. Fei-Fei, "Deep visual-semantic alignments for generating image descriptions," *arXiv preprint arXiv:1412.2306*, 2014.
- [33] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [34] X. Li, H. Qin, Y. Wang, Y. Zhang, and Q. Dai, "Dept: Depth estimation by parameter transfer for single still images," *Asian Conference on Computer Vision (ACCV)*, 2014.



Xiu Li received her Ph.D. degree in computer integrated manufacturing in 2000. Since then, she has been working in Tsinghua University. Her research interests include data mining, business intelligence systems, knowledge management systems and decision support systems, etc.



Yangang Wang is currently an associate researcher at Microsoft Research Asia (MSRA). In 2014, he received my Ph.D. at Tsinghua University (THU) in Beijing, under the supervision of Prof. Qionghai Dai. In 2009, he received his B.E. from the school of Instrument Science and Engineering at Southeast University (SEU). His research interests involve in the area of computer vision, computer graphics and computational photography.



Yongbing Zhang received the B.A. degree in English and the M.S. and Ph.D degrees in computer science from the Harbin Institute of Technology, Harbin, China, in 2004, 2006, and 2010, respectively. He is currently with the Graduate School at Shenzhen, Tsinghua University, Shenzhen, China. He has authored or coauthored over 30 conference and journal papers. His current research interests include video processing, image and video coding, video streaming, and transmission.



**Hongwei Qin** received the B.S. degree from Tsinghua University, Beijing, China, in 2012, where he is currently working toward the Ph.D. degree in the Department of Automation.



Qionghai Dai (SM05) received the B.S. degree from Shanxi Normal University, Shanxi, China, in 1987, and the M.E. and Ph.D. degrees from Northeastern University, Shenyang, China, in 1994 and 1996, respectively. Since 1997, he has been with the faculty of Tsinghua University, Beijing, China, and is currently a Professor and the Director of the Broadband Networks and Digital Media Laboratory, Beijing. His current research interests include video communication, computer vision, and computational photography.