

2D-to-3D Image Conversion by Learning Depth from Examples

Janusz Konrad, Meng Wang, and Prakash Ishwar

Department of Electrical and Computer Engineering, Boston University
8 Saint Mary's St., Boston, MA USA 02215

{jkonrad, wangmeng, pi}@bu.edu

Abstract

Among 2D-to-3D image conversion methods, those involving human operators have been most successful but also time-consuming and costly. Automatic methods, that typically make use of a deterministic 3D scene model, have not yet achieved the same level of quality as they often rely on assumptions that are easily violated in practice. In this paper, we adopt the radically different approach of “learning” the 3D scene structure. We develop a simplified and computationally-efficient version of our recent 2D-to-3D image conversion algorithm. Given a repository of 3D images, either as stereopairs or image+depth pairs, we find k pairs whose photometric content most closely matches that of a 2D query to be converted. Then, we fuse the k corresponding depth fields and align the fused depth with the 2D query. Unlike in our original work, we validate the simplified algorithm quantitatively on a Kinect-captured image+depth dataset against the Make3D algorithm. While far from perfect, the presented results demonstrate that on-line repositories of 3D content can be used for effective 2D-to-3D image conversion.

1. Introduction

The availability of 3D hardware today (3D TVs, Blu-Ray players, handheld gaming consoles, cell phones, still and video cameras) is not yet matched by 3D content production. Although methods have been proposed to convert 2D images to 3D stereopairs, the most successful approaches are interactive, i.e., involve human operators [4, 1, 7], and, therefore, time-consuming and costly.

The problem of depth estimation from a single 2D image, which is the main step in 2D-to-3D conversion, can be formulated in various ways, for example as a shape-from-shading problem [14]. However, this problem is severely under-constrained; quality depth estimates can be found only for special cases. Other methods, often called multi-view stereo, attempt to recover depth by estimating scene

geometry from multiple images not taken simultaneously. For example, a moving camera permits structure-from-motion estimation [11] while a fixed camera with varying focal length permits depth-from-defocus estimation [10]. Both are examples of the use of multiple images of the same scene captured at different times or under different exposure conditions (e.g., all images of the Statue of Liberty). Although such methods are similar in spirit to our approach, the main difference is that while these methods use images known to depict the same scene as the query image, we use all images available in a large repository and automatically select suitable ones for depth recovery.

Some electronics manufacturers have developed real-time 2D-to-3D converters that rely on stronger assumptions and simpler processing than the methods listed above, e.g., faster-moving or larger objects are closer to the viewer, higher frequency of texture belongs to objects located further away, etc. Although such methods may work well in some cases, in general it is very difficult, if not impossible, to construct a deterministic scene model that covers all possible background and foreground combinations.

Recently, machine learning techniques based on image parsing have been used to estimate the depth map from single monocular images [8, 6]. Such methods have the potential to generate depth maps for 2D visual material, but currently work only on few types of images using carefully-selected training data (precise, laser-scanned depth estimates or manually-annotated semantic depth classes). The algorithm we describe in this paper is somewhat similar to these two methods except that it applies to arbitrary scenes and requires no manual annotation.

The trend to use large image databases for various computer vision tasks, such as object recognition [12] and image saliency detection [13], has recently inspired us to develop a data-driven approach to 2D-to-3D conversion [5]. However, this approach is computationally involved due to the use of SIFT-flow for disparity warping. Furthermore, it has only been tested qualitatively which is subjective.

In this paper, we propose a simplified algorithm that “learns” the scene depth from a large repository of im-

This work was supported by the NSF under grant ECS-0905541.

age+depth pairs and is more efficient computationally than our original algorithm [5]. Furthermore, we validate the simplified algorithm quantitatively on a Kinect-captured image+depth dataset [9] against the Make3D algorithm [8]. While far from perfect, the presented results demonstrate that on-line repositories of 3D content can be used for effective 2D-to-3D image conversion.

2. Proposed approach

Our approach is built upon a key observation and an assumption [5]. The key observation is that among millions of image+depth pairs available on-line, there likely exist many pairs whose 3D content matches that of a 2D input (query). The assumption is that two images that are photometrically similar are likely to have similar 3D structure (depth). This is not unreasonable since photometric properties are often correlated with 3D content (depth, disparity). For example, edges in a depth map almost always coincide with photometric edges. We rely on the above observation and assumption to “learn” depth from a dictionary of image+depth pairs and render a stereopair in the following steps:

1. **k nearest-neighbor (k NN) search:** finding k image+depth pairs that are photometrically most similar to the 2D query,
2. **depth fusion:** median filtering of the k depth fields,
3. **cross-bilateral depth filtering:** smoothing of the median-fused depth field to remove spurious variations, while preserving depth discontinuities,
4. **stereo rendering:** generation of the right image of the stereopair using the 2D query (left) image and smoothed median depth field followed by suitable processing of occlusions and newly-exposed areas.

Although the original algorithm [5] includes an additional step of depth (disparity) warping *via* SIFT-flow to better align k NN depth (disparity) fields with the 2D query, this step is computationally demanding while bringing only small quality improvement to the depth estimates. We forgo this step for the sake of computational efficiency without sacrificing much performance.

Fig. 1 shows the block diagram of our approach. The sections below provide a description of each step and some high-level mathematical detail. In these sections, Q is the query image for which a right image Q_R is being sought. We assume that a database $\mathcal{I} = \{(I^1, d^1), (I^2, d^2), \dots\}$ of image+depth pairs (I^k, d^k) is available. Note that a database of stereoscopic videos, such as YouTube 3D, could be processed to extract image+depth pairs. The goal is to find a depth estimate \hat{d} and then a right-image estimate \hat{Q}_R given the 3D database \mathcal{I} .

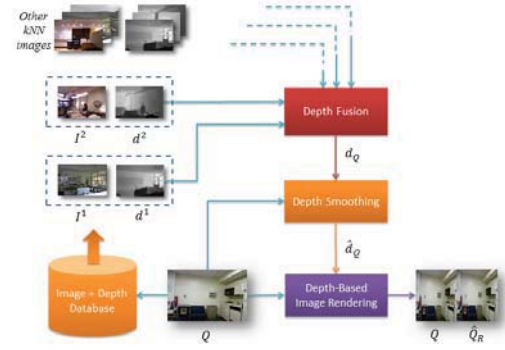


Figure 1. Block diagram of the overall algorithm; algorithmic details for each block are provided in the sections below.

2.1. k NN search

There exist two types of images in a large 3D image repository: those that are relevant for determining depth in a 2D query image, and those that are irrelevant. Images that are not photometrically similar to the 2D query need to be rejected because they are not useful for estimating depth (as per our assumption). Note that although we might miss some depth-relevant images, we are effectively limiting the number of irrelevant images that could potentially be more harmful to the 2D-to-3D conversion process. The selection of a smaller subset of images provides the added practical benefit of computational tractability when the size of the dictionary is very large.

Our 2D query image Q is the left image from a stereopair whose right image Q_R is unknown. We assume that a database of 3D images or videos \mathcal{I} , such as the NYU depth database [9] or YouTube 3D, is available, and that for each RGB image I^i in the database the corresponding depth field d^i is either known or can be computed from a stereopair.

One method for selecting a useful subset of depth-relevant images from a large dictionary is to select only the k images that are closest to the input where closeness is measured by some distance function capturing global image properties such as color, texture, edges, etc. As this distance function, we use the Euclidean norm of the difference between histograms of oriented gradients (HOGs) [2] computed from two images. Each HOG consists of 144 real values (4×4 blocks with 9 gradient direction bins) that can be efficiently computed. This image closeness measure is significantly less computationally complex than the weighted Hamming distance between binary hashes of features that we used originally [5].

We perform a search for top matches to our 2D query among all 3D images in the database \mathcal{I} . The search returns an ordered list of image+depth pairs, from the most to the least photometrically similar *vis-à-vis* the 2D query. We discard all but the top k matches (k NNs) from this list.

Fig. 2 shows search results for four 2D query images (office, bedroom, dining room and kitchen). An examination

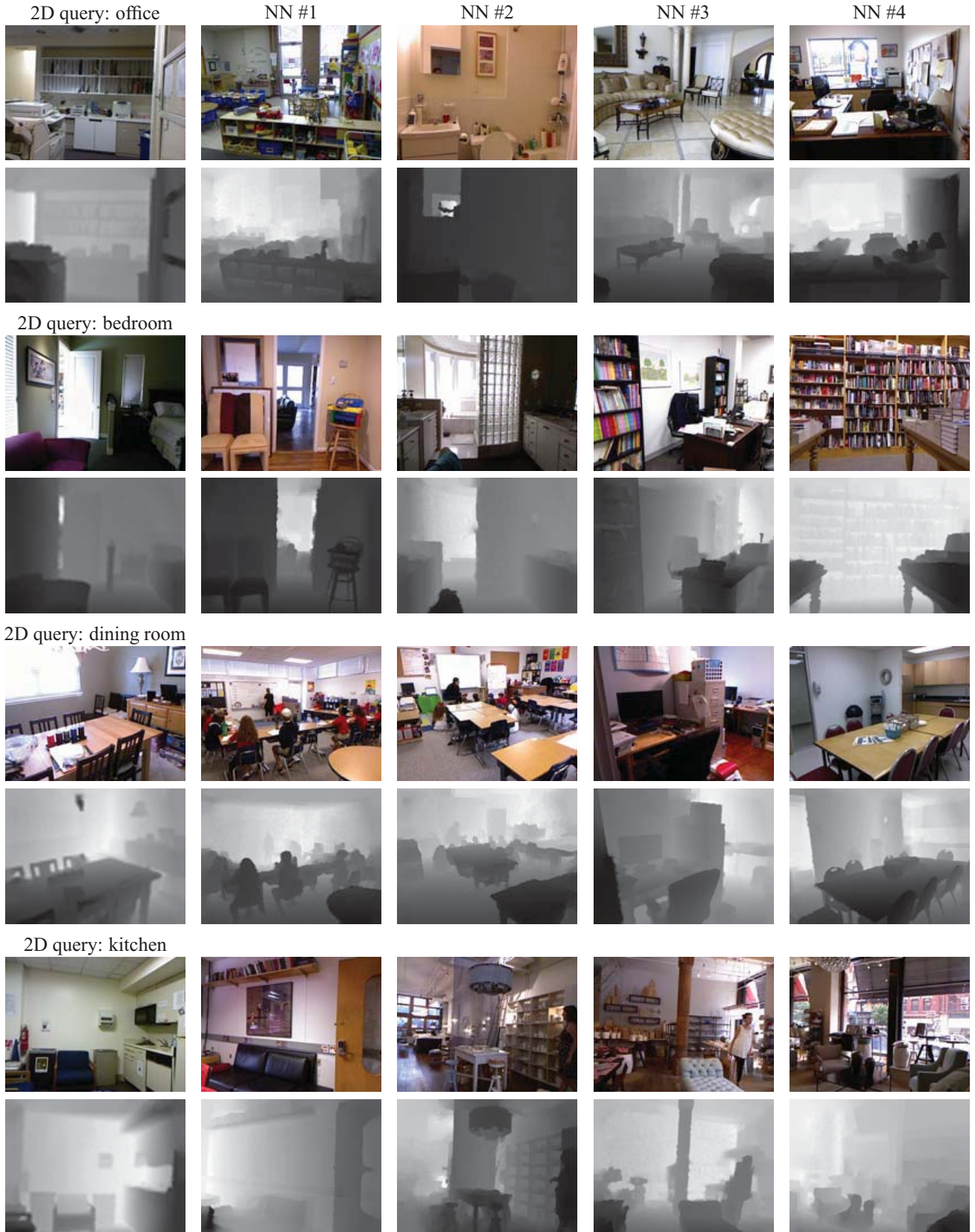


Figure 2. RGB image and depth field of four 2D queries (left column), and their four nearest neighbors (columns 2-5) retrieved using the Euclidean norm on the difference between histograms of gradients [2]. All image+depth pairs are from NYU depth dataset (see Section 3).

of the search results shows that the majority of retained images share a very similar global 3D structure with the query image - a large table in the dining room or a slanted wall on right in the kitchen. Although none of the k NNs perfectly matches the corresponding 2D query, the general underlying depth is closely related to that expected in the query.

The average similarity between a query and its k -th nearest neighbor usually decays with increasing k . While for large databases, larger values of k may be appropriate, since there are many good matches, for smaller databases this may not be true. Therefore, a judicious selection of k is important. For now, we denote by \mathcal{K} the set of indices i of image+depth pairs that are the top k photometrically-nearest neighbors of the query Q .

2.2. Depth fusion

In general, none of the NN image+depth pairs $(I^i, d^i), i \in \mathcal{K}$ match a query Q accurately (Fig. 2). However, the location of some objects (e.g., furniture) and parts of the background (e.g., walls) is quite consistent with those in the query. If a similar object (e.g., table) appears at a similar location in several k NN images, it is likely that such an object also appears in the query and the depth field being sought should reflect this. We compute this depth field by applying the median operator across the k NN depths at each spatial location \mathbf{x} as follows:

$$d[\mathbf{x}] = \text{median}\{d^i[\mathbf{x}], \forall i \in \mathcal{K}\}. \quad (1)$$

Examples of the fused depth fields d are shown in the central column of Fig. 3. Although these depths are overly smooth, they provide a globally-correct, although coarse, assignment of distances to various areas of the scene.

2.3. Cross-bilateral depth filtering

While the median-based fusion helps make depth more consistent globally, the fused depth is overly smooth and locally inconsistent with the query image due to:

1. misalignment of edges between the fused depth field and query image,
2. lack of fused depth edges where sharp object boundaries occur,
3. lack of fused depth smoothness where smooth depth changes are expected.

In order to correct this, similarly to Agnot *et al.* [1], we apply a variant of bilateral filtering to the fused depth d with the RGB query image as a reference. Bilateral filtering is an edge-preserving image smoothing method that applies anisotropic diffusion controlled by the local image content [3]. We apply bilateral filtering to the fused depth with two goals: alignment of the depth edges with those of

the query image Q and local noise/granularity suppression in the fused depth d . This is implemented as follows:

$$\begin{aligned} \hat{d}[\mathbf{x}] &= \frac{1}{\gamma[\mathbf{x}]} \sum_{\mathbf{y}} d[\mathbf{y}] h_{\sigma_s}(\mathbf{x} - \mathbf{y}) h_{\sigma_e}(Q[\mathbf{x}] - Q[\mathbf{y}]), \\ \gamma[\mathbf{x}] &= \sum_{\mathbf{y}} h_{\sigma_s}(\mathbf{x} - \mathbf{y}) h_{\sigma_e}(Q[\mathbf{x}] - Q[\mathbf{y}]), \end{aligned} \quad (2)$$

where \hat{d} is the filtered depth field and $h_{\sigma}(\mathbf{x}) = \exp(-\|\mathbf{x}\|^2/2\sigma^2)/2\pi\sigma^2$ is a Gaussian weighting function. Note that the directional smoothing of d is controlled by the query image via the weight $h_{\sigma_e}(Q[\mathbf{x}] - Q[\mathbf{y}])$. For large discontinuities in Q , the weight $h_{\sigma_e}(Q[\mathbf{x}] - Q[\mathbf{y}])$ is small and thus the contribution of $d[\mathbf{y}]$ to the output is small. However, when $Q[\mathbf{y}]$ is similar to $Q[\mathbf{x}]$ then $h_{\sigma_e}(Q[\mathbf{x}] - Q[\mathbf{y}])$ is relatively large and the contribution of $d[\mathbf{y}]$ to the output is larger. In essence, depth filtering (smoothing) is happening along (and not across) query edges.

Fig. 3 compares the fused depth before cross-bilateral filtering (d) and after (\hat{d}). The filtered depth preserves the global properties captured by the unfiltered depth field d , and is smooth within objects and in the background. At the same time it keeps edges sharp and aligned with the query image structure.

2.4. Stereo rendering

In order to generate an estimate of the right image \hat{Q}_R from the 2D query Q , we need to compute the disparity δ from the estimated depth \hat{d} . Assuming that the fictitious image pair (Q, \hat{Q}_R) was captured by parallel cameras with baseline B and focal length f , the disparity is simply $\delta[x, y] = Bf/\hat{d}[\mathbf{x}]$, where $\mathbf{x} = [x, y]^T$. We forward-project the 2D query Q to produce the right image:

$$\hat{Q}_R[x + \delta[x, y], y] = Q[x, y] \quad (3)$$

while rounding the location coordinates $(x + \delta[x, y], y)$ to the nearest sampling grid point. We handle occlusions by depth ordering: if $(x_i + \delta[x_i, y_i], y_i) = (x_j + \delta[x_j, y_j], y_j)$ for some i, j , we assign to the location $(x_i + \delta[x_i, y_i], y_i)$ in \hat{Q}_R an RGB value from that location (x_i, y_i) in Q whose disparity $\delta[x_i, y_i]$ is the largest. In newly-exposed areas, i.e., for x_j such that no x_i satisfies $(x_j, y_i) = (x_i + \delta[x_i, y_i], y_i)$, we apply simple inpainting using `inpaint_nans` from *MatlabCentral*.

3. Experimental Results

We have tested our approach on a database of indoor scenes captured by the Microsoft Kinect camera [9] that contains 1449 pairs of RGB images and corresponding depth fields. Kinect cameras use structured (infrared) light to provide an accurate depth map of the captured scene, but do not work well in daylight and at large distances; the above database is limited to indoor scenes.

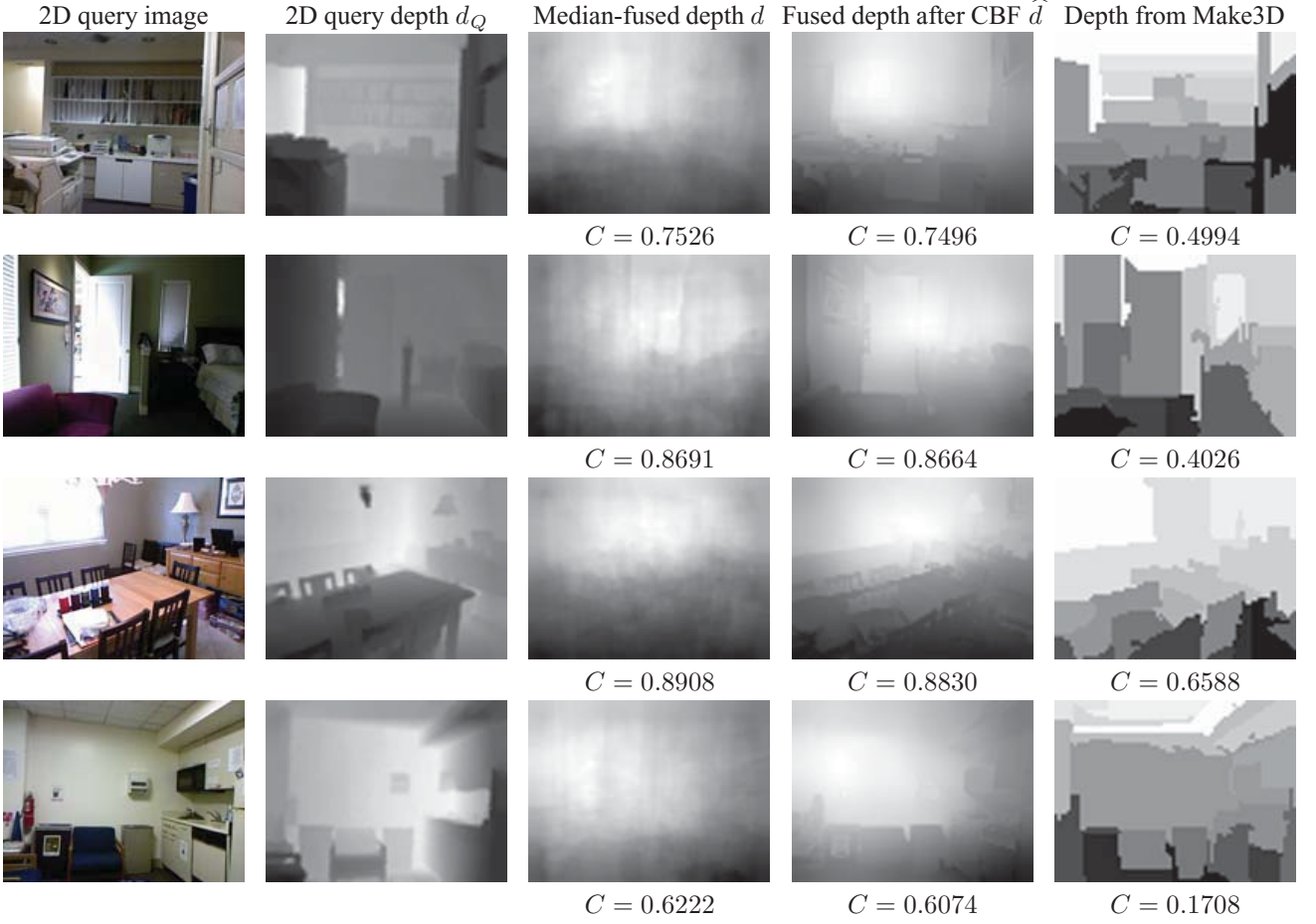


Figure 3. RGB image and depth field of queries from Fig. 2 as well as the estimated depth after fusion and cross-bilateral filtering (CBF) and depth computed using the Make3D algorithm. Normalized depth cross-covariances are included under each estimated depth field.

In order to evaluate the performance of the proposed algorithm quantitatively, we applied leave-one-out cross-validation (LOOCV) as follows. We selected one image+depth pair from the database as the 2D query (Q, d_Q) treating the remaining pairs as the dictionary based on which a depth estimate \hat{d} and a right-image estimate \hat{Q}_R are computed. As the quality metric, we used normalized cross-covariance between the estimated depth \hat{d} and the ground-truth depth d_Q defined as follows:

$$C = \frac{1}{N\sigma_{\hat{d}}\sigma_{d_Q}} \sum_{\mathbf{x}} (\hat{d}[\mathbf{x}] - \mu_{\hat{d}})(d_Q[\mathbf{x}] - \mu_{d_Q}) \quad (4)$$

where N is the number of pixels in \hat{d} and d_Q , $\mu_{\hat{d}}$ and μ_{d_Q} are the empirical means of \hat{d} and d_Q , respectively, while $\sigma_{\hat{d}}$ and σ_{d_Q} are the corresponding empirical standard deviations. The normalized cross-covariance C takes values between -1 and +1 (for values close to +1 the depths are very similar and for values close to -1 they are complementary).

In order to select a suitable value of k , we ran the LOOCV test for each image in the Kinect database for all k

	Proposed algorithm		Make3D
	with warping	no warping	
Average C	0.71	0.71	0.45
Median C	0.76	0.75	0.48
Processing time	16h	5s	12h

Table 1. Average and median normalized cross-covariance C and average processing time (3.4GHz CPU) obtained in LOOCV tests on the Kinect depth dataset using the proposed algorithm with warping [5] and without, and also using Make3D [8].

from 1 to 120 and averaged the resulting cross-covariance C across all tests. The average C rapidly rose for small k , achieved maximum at $k = 45$ and then gently rolled-off. Therefore, in all experiments below we used $k = 45$.

Table 1 shows the average and median of cross-covariance C obtained from 1449 LOOCV tests using the proposed algorithm with and without warping. The warping of each k NN depth d^i via SIFT-flow to better align foreground objects was originally proposed in the disparity-learning algorithm [5] at significant computational cost. We are advocating here conversion without warping to reduce complexity. In both cases, an image similarity metric based

on HOGs rather than hashes of features [5] was used. As can be seen in Table 1, the algorithm without warping is several orders of magnitude faster with, basically, no loss of average depth estimation fidelity.

Since most of the fully-automatic 2D-to-3D conversion methods have been developed by 3D equipment manufacturers, the employed algorithms are proprietary. The only automatic 2D-to-3D conversion method for which we were able to find a run-time code was Make3D developed by Saxena *et al.* [8]. Make3D estimates 3D scene structure from a single still image of an unstructured environment by supervised learning of 3D position and orientation of small homogeneous patches in the image. The original Make3D algorithm was trained on images and associated laser-scanned depth maps of mostly architectural structures. Admittedly, it was not optimized for indoor scenes that the Kinect depth dataset is composed of, however we were unable to re-train Make3D for indoor data and it was the only algorithm available for comparison. As can be seen in Table 1, Make3D achieves normalized cross-covariance C of about 0.45-0.48, significantly less than C for the proposed algorithm.

In terms of the computational complexity, our algorithm has a significant edge as well. Applied to all 1449 images of the Kinect database with $k = 45$, it required only 5 seconds as opposed to 12 hours for Make3D. We must note at this point that, due to Make3D's complexity, the depth learning step was performed on reduced-resolution images and depth fields (80×60) as opposed to full-resolution (640×480). Had we used full-resolution data, we would have to wait over 4 weeks for Make3D output. We believe that depth learning at low resolution is acceptable if depth edges are aligned with photometric boundaries, because depth varies smoothly within objects and background. The estimated depth fields \hat{d} were interpolated to full resolution prior to the right-image rendering.

We would like to point out that although C values shown in Fig. 3 are slightly lower for depth fields after cross-bilateral filtering, the depth edge alignment with the query and the high piece-wise depth smoothness are both perceptually beneficial in 3D viewing. In Fig. 4, we show anaglyph images constructed from (Q, \hat{Q}_R) image pairs for the ground truth depth d_Q , and the estimated depths \hat{d} using the proposed approach and Make3D. Although neither conversion is flawless, errors on the bulletin board in the office image and under the chest of drawers in the dining room image produced by Make3D cause significant visual discomfort. For comparison, Fig. 5 shows the kitchen image converted using YouTube 3D. While no quantitative comparison is possible since YouTube does not provide any depth field, visually one is left with the sensation of a cardboard effect; there is no gradual increase of disparity towards the viewer, unlike in the image converted by the proposed method (bottom of the middle column in Fig. 4).

4. Conclusions

We have proposed a simplified data-driven 2D-to-3D conversion method and have objectively validated its performance against state-of-the-art Make3D algorithm. The proposed algorithm compares favorably in terms of both estimated depth quality and computational complexity. Admittedly, the validation was limited to a database of indoor scenes on which Make3D was not trained. The generated anaglyph images produce a comfortable 3D perception but are not completely void of distortions. With the continuously increasing amount of 3D data on-line and with the rapidly growing computing power in the cloud, the proposed algorithm seems a promising alternative to operator-assisted 2D-to-3D conversion.

References

- [1] L. Agnot, W.-J. Huang, and K.-C. Liu. A 2D to 3D video and image conversion technique based on a bilateral filter. In *Proc. SPIE Three-Dimensional Image Processing and Applications*, volume 7526, Feb. 2010. 1, 4
- [2] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *Proc. IEEE Conf. Computer Vision Pattern Recognition*, pages 886–893, 2005. 2, 3
- [3] F. Durand and J. Dorsey. Fast bilateral filtering for the display of high-dynamic-range images. *ACM Trans. Graph.*, 21:257–266, July 2002. 4
- [4] M. Guttman, L. Wolf, and D. Cohen-Or. Semi-automatic stereo extraction from video footage. In *Proc. IEEE Int. Conf. Computer Vision*, pages 136–142, Oct. 2009. 1
- [5] J. Konrad, G. Brown, M. Wang, P. Ishwar, C. Wu, and D. Mukherjee. Automatic 2D-to-3D image conversion using 3D examples from the Internet. In *Proc. SPIE Stereoscopic Displays and Applications*, volume 8288, Jan. 2012. 1, 2, 5, 6
- [6] B. Liu, S. Gould, and D. Koller. Single image depth estimation from predicted semantic labels. In *Proc. IEEE Conf. Computer Vision Pattern Recognition*, pages 1253–1260, June 2010. 1
- [7] R. Phan, R. Rzeszutek, and D. Androutsos. Semi-automatic 2D to 3D image conversion using scale-space random walks and a graph cuts based depth prior. In *Proc. IEEE Int. Conf. Image Processing*, Sept. 2011. 1
- [8] A. Saxena, M. Sun, and A. Ng. Make3D: Learning 3D scene structure from a single still image. *IEEE Trans. Pattern Anal. Machine Intell.*, 31(5):824–840, May 2009. 1, 2, 5, 6
- [9] N. Silberman and R. Fergus. Indoor scene segmentation using a structured light sensor. In *Proc. Int. Conf. on Computer Vision - Workshop on 3D Represent. and Recogn.*, 2011. 2, 4
- [10] M. Subbarao and G. Surya. Depth from defocus: A spatial domain approach. *Intern. J. Comput. Vis.*, 13:271–294, 1994. 1
- [11] R. Szeliski and P. H. S. Torr. Geometrically constrained structure from motion: Points on planes. In *European Workshop on 3D Structure from Multiple Images of Large-Scale Environments*, pages 171–186, 1998. 1



Figure 4. Anaglyph images generated using the ground-truth depth and depths estimated by the proposed and Make3D algorithm.

- [12] A. Torralba, R. Fergus, and W. T. Freeman. 80 million tiny images: A large data set for nonparametric object and scene recognition. *IEEE Trans. Pattern Anal. Machine Intell.*, 30(11):1958–1970, Nov. 2008. 1
- [13] M. Wang, J. Konrad, P. Ishwar, K. Jing, and H. Rowley. Image saliency: From intrinsic to extrinsic context. In *Proc. IEEE Conf. Computer Vision Pattern Recognition*, pages 417–424, June 2011. 1
- [14] R. Zhang, P. S. Tsai, J. Cryer, and M. Shah. Shape-from-shading: A survey. *IEEE Trans. Pattern Anal. Machine Intell.*, 21(8):690–706, Aug. 1999. 1



Figure 5. Anaglyph image from YouTube 2D-to-3D conversion.