# A fully end-to-end deep learning approach for real-time simultaneous 3D reconstruction and material recognition

Cheng Zhao, Li Sun and Rustam Stolkin
Extreme Robotics Lab, University of Birmingham, UK.
IRobotCheng@gmail.com

*Abstract*—This paper addresses the problem of *simultaneous* 3D reconstruction and material recognition and segmentation. Enabling robots to recognise different materials (concrete, metal etc.) in a scene is important for many tasks, e.g. robotic interventions in nuclear decommissioning. Previous work on 3D semantic reconstruction has predominantly focused on recognition of everyday domestic objects (tables, chairs etc.), whereas previous work on material recognition has largely been confined to single 2D images without any 3D reconstruction. Meanwhile, most 3D semantic reconstruction methods rely on computationally expensive post-processing, using Fully-Connected Conditional Random Fields (CRFs), to achieve consistent segmentations. In contrast, we propose a deep learning method which performs 3D reconstruction while simultaneously recognising different types of materials and labeling them at the pixel level. Unlike previous methods, we propose a fully end-to-end approach, which does not require hand-crafted features or CRF post-processing. Instead, we use only learned features, and the CRF segmentation constraints are incorporated inside the fully end-to-end learned system. We present the results of experiments, in which we trained our system to perform real-time 3D semantic reconstruction for 23 different materials in a real-world application. The run-time performance of the system can be boosted to around 10Hz, using a conventional GPU, which is enough to achieve real-time semantic reconstruction using a 30fps RGB-D camera. To the best of our knowledge, this work is the first real-time end-to-end system for simultaneous 3D reconstruction and material recognition.

*Index Terms*—3D semantic reconstruction, material recognition, real-time, fully end-to-end, deep neural network

## I. INTRODUCTION

Real-time 3D semantic reconstruction is required in many robotics applications, such as autonomous navigation or grasping and manipulation. While a variety of well-known methods [1][2][3] can reconstruct accurate 3D maps at real-time frame rates, the resulting point-clouds contain no semantic-level understanding of the observed scenes. Hence, the problem of 3D semantic reconstruction is attracting increasing attention in the robotics research community. Recent methods [4][5][6][7][8] not only generate a 3D point cloud map, but also simultaneously assign a semantic label to each point in the cloud. However, these methods are typically designed to search for everyday domestic 3D objects (e.g. "table", "chair" etc.) in domestic (non-industrial) scenes.

In contrast, an ability to recognize different kinds of *materials* could play a very important role in numerous robotics applications. Understanding material properties (e.g. friction or deformability) of objects can be used to inform grasp planning and manipulation. Rescue robots should understand
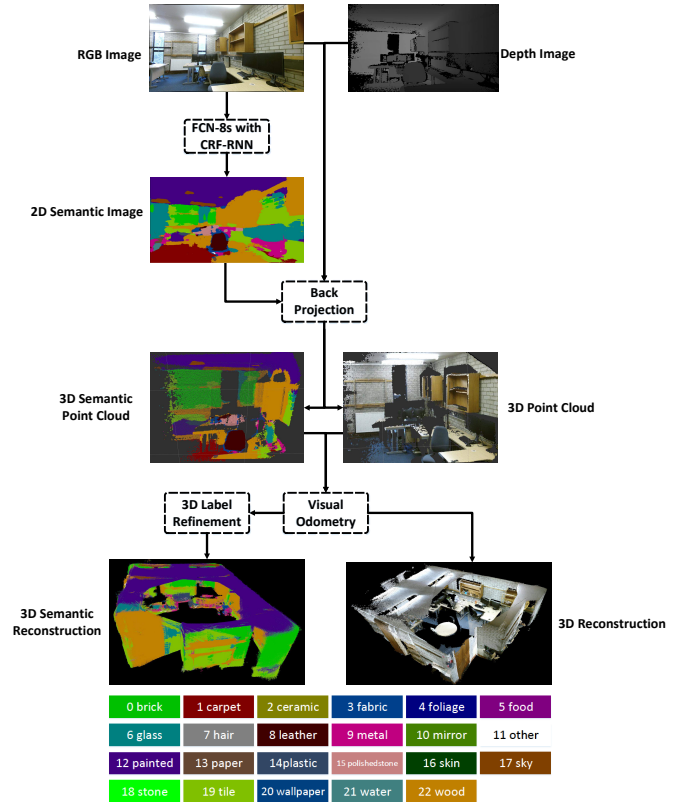


Fig. 1: **Pipeline of proposed simultaneous 3D reconstruction and material recognition system**. Firstly, FCN-8s with CRF-RNN is employed for 2D material recognition using the RGB image from RGB-D camera. Then the semantically labeled RGB image, and the corresponding depth image, are combined together through back-projection to generate a semantic point cloud for each key frame. Finally, all semantic point clouds are combined incrementally using visual odometry, and Bayesian update is employed for label probability refinement.

their surrounding materials when planning movements through precarious rubble under a collapsed building. In nuclear decommissioning, robots must enter hazardous zones (in very old legacy buildings with significant uncertainty) to perform "characterisation". Understanding which materials make up the scene is a critical aspect of characterisation, and will inform subsequent interventions, e.g. cutting, dismantling, cleaning,

manipulating. Unfortunately, previous work [9][10][11][12] only addresses the problems of detecting and segmenting materials in a single RGB image, and does not perform 3D material reconstruction.

In this paper, we present a fully end-to-end system, which performs real-time 3D reconstruction while simultaneously recognizing and labeling each pixel according to its material, Fig.1. The main contributions of this paper can be summarized as follows:

- To the best of our knowledge, this is the first system to perform simultaneous 3D reconstruction and material recognition.
- Hand-crafted features or post-processing CRF optimization are not required. In contrast, the system is fully end-to-end learned, and this helps to deliver real-time performance, as well as generality for different applications.
- The run-time performance of the whole system can be boosted, using a conventional GPU, to around 10Hz, which is enough to achieve real-time semantic reconstruction using a 30fps RGB-D camera.
- We demonstrate our method in a real-world application, reconstructing a room while simultaneously recognizing and labeling 23 different materials.

## II. RELATED WORK

In this section, we firstly review real-time 3D semantic reconstruction in Section II-A and material recognition in Section II-B. Then we will give a discussion in Section II-C.

### A. Real-time 3D semantic reconstruction

Recent real-time 3D reconstruction and SLAM approaches, e.g. [1] (Visual-features with RANSAC), [2] (Direct image alignment based on optimization), or [3](Point cloud alignment based on ICP) can effectively generate dense or semi-dense 3D maps, but they have no understanding of the observed scenes and objects.

The more complex problem of 3D semantic reconstruction remains an open research problem. Recent approaches can be grouped into two main categories: 3D semantic reconstruction base on 3D template matching [4][13][5], and 3D semantic reconstruction base on 2D semantic segmentation [6][7].

The former methods rely on 3D template matching, so can only be used in situations with many repeated and identical objects. Only known 3D objects can be recognised, and semantic labeling for the remainder of the scene is not possible. The latter methods typically employ visual features, combined with a classifier such as random forest, for 2D semantic segmentation. The visual features are hand-crafted and also requires a non-linear transformation between local and global descriptors. Unlike CNN features, which can be learned from training data in an end-to-end fashion, such methods require application-specific, human-designed components.

The work most closely related to ours is [8], which performs dense, 3D semantic mapping of indoor scenes, using deconvolutional neural networks[14]. Real-time frame-rates of about 25 Hz are achieved *without* CRF optimization post-processing.

But most such methods [6][7][8] rely on using fully connected CRF[15] optimization as an offline post-processing step, following online 3D reconstruction, i.e. these methods do not actually achieve semantic mapping in real-time. Additionally, these methods are not fully end-to-end trainable, and there is no interaction between classifier learning and CRF learning. The parameters of classifier and CRF cannot be jointly learned in a united framework. Furthermore, all of the above methods are focused on semantic *object* recognition. In contrast, our work is the first method that achieves simultaneous 3D reconstruction with semantic *material* labeling, and we achieve both 3D reconstruction *and* semantic labeling simultaneously in real-time.

### B. Material recognition

Materials recognition is a challenging research topic due to wide variation in appearance within categories. Previous material recognition research predominantly focused on material classification, and did not achieve pixel-wise material segmentation. Most previous work employed hand-crafted visual features, e.g. reflectance-based edge features [16], variances of oriented gradients [17], and pairwise local binary patterns [18]. Recently CNN features [19][20][21] have been employed to achieve the state-of-the-art results of material classification in many public material datasets. In addition to the 2D features, [9] combined 3D geometry (surface normals, camera intrinsic and extrinsic parameters) with 2D features (texture and color) to improve material classification.

For pixel-wise material segmentation, [10] convert patch-based trained CNN classifiers into an efficient fully convolutional framework combined with a fully connected CRF to perform pixel-wise material recognition. [11] combined local appearance with separately recognized global contextual cues including objects and places, which can lead to a superior result. They employed fully convolutional network(FCN) [22] followed by recurrent neural network(RNN) for dense pixel-wise material segmentation. [12] proposed a novel CNN architecture trained on 4D light-field images and employ FCN for per-pixel material recognition. However, in contrast to our work, none of these methods perform material recognition simultaneously with 3D reconstruction.

### C. Discussion

In summary, previous real-time semantic 3D reconstruction methods have focused on object recognition, and not on material recognition. Predominantly, such methods require post-processing with a fully-connected CRF, and are not fully end-to-end. While there is literature on per-pixel material recognition in 2D images, no previously reported methods perform 3D material reconstruction. To our best knowledge, this work is the first real-time end-to-end system for simultaneous 3D reconstruction and material recognition.

## III. METHODS

### A. Overview

The pipeline of simultaneous 3D reconstruction and material recognition comprises three units as illustrated in Figure 1: a

real-time 3D reconstruction unit based on RGB-D SLAM[1], a 2D material recognition unit based on FCN-8s[22] with CRF-RNN[23], and a 3D semantic reconstruction unit based on Bayesian update. Firstly, the FCN-8s with CRF-RNN is employed for 2D material recognition using the RGB image from RGB-D camera. Then the semantically labeled RGB image, and the corresponding depth image, are combined together through back-projection to generate a semantic point cloud for each key frame. Finally, all semantic point clouds are combined incrementally using visual odometry, and Bayesian update is employed for label probability refinement.

### B. RGB-D SLAM Mapping

We use the RGB-D SLAM method of [1] for real-time 3D reconstruction. It is a graph-based SLAM system which includes a front-end system to processes the RGB-D sensor data to calculate geometric relationships through visual features based on RANSAC and ICP. Subsequently, the back-end system registers pairs of image frames to construct a pose graph. G2O[24] is employed for graph optimization to obtain a maximum likelihood solution for the camera trajectory. Finally, RGB-D sensor data is combined together to generate a 3D point cloud.

In our system, RGB-D SLAM plays two important roles: 1) it can provide the transformation information between two adjacent semantic point clouds, enabling incremental semantic label fusion; 2) the visual odometry is used to combine all semantic point clouds to generate a global semantic map.

### C. 2D material recognition

Our neural network is implemented in caffe [25] framework and employs FCN-8s [22] followed by CRF-RNN [23] architecture.

*1) FCN:* FCN is the first end-to-end and pixel-to-pixel semantic segmentation architecture which can take an input of arbitrary size and generate correspondingly-sized output images. This architecture is based on the VGG 16-layer[26] network. The learned representations in the VGG-16 network can be transferred through fine-tuning our network, using the extracted patches in public material dataset MINC[10]. Next, we transplant the fully connected VGG network into a fully convolutional VGG network and inherit the weights of the fine-tuning network. Finally, the FCN-32s, FCN-16s and FCN-8s networks are trained using MINC dataset sequentially. FCN defines a skip architecture which can combine semantic information from a deep, coarse layer with shape information from a shallow, fine layer. The output of the deep layer has rich semantic information but loses most of the shape information, while the shallow layer has rich shape information but lacks semantic information. Therefore FCN improve the accuracy of semantic segmentation by fusing the outputs from both deep and shallow layers.

FCN has convolutional filters with large receptive fields and 5 pooling layers. It does not incorporate smoothness constraints between neighbouring pixels. Hence, it can only generate coarse pixel-wise predictions with blob-like shapes.

In our system, FCN-8s is employed as the first part of the network to provide unary potentials to the CRF-RNN.

*2) CRF-RNN:* CRF-RNN, following FCN, combines the strengths of both FCN and fully-connected CRF into a single end-to-end unified framework. Fully-connected CRF accounts for contextual information by minimising the energy $E(x)$ function in the Gibbs distribution, to generate the most likely label assignment $x$.

Energy function $E(x)$ consists of a unary data term and pairwise smoothness term, Equation 1. Unary term $\psi_u(x_i)$ is obtained from the FCN-8s, which predicts pixel labels without considering smoothness or consistency of label assignments. The pairwise term $\psi_p(x_i, x_j)$ encourages similar labeling of pixels with similar properties, while penalizing similar pixels which have different labels.

$$E(x) = \sum_i \psi_u(x_i) + \sum_{i<j} \psi_p(x_i, x_j) \qquad (1)$$

Pairwise potentials are modeled as a linear combination of $M$ Gaussian edge potential kernels as shown in equation 2. $f_i$ is the feature vector of pixel $i$ (e.g. spatial or colour information). $k_G^{(m)}$ is a Gaussian kernel applied to feature vectors. $\omega^{(m)}$ is the linear combination weight. The Potts model $\mu(x_i, x_j) = [x_i \neq x_j]$ is the label compatibility function.

$$\psi_p(x_i, x_j) = \mu(x_i, x_j) \sum_{m=1}^{M} \omega^{(m)} k_G^{(m)}(f_i, f_j) \qquad (2)$$

A bilateral appearance potential and a spatial smoothing potential($M = 2$) are employed in the pairwise potentials as shown in equation 3. $p_i$ and $p_j$ are the $x, y, z$ spatial information and $I_i$ and $I_j$ are the $R, G, B$ colour information. $\theta_\alpha, \theta_\beta$ and $\theta_\gamma$ are the parameters of Gaussian kernels.

$$k(f_i, f_j) = \omega^{(1)} exp(-\frac{|p_i - p_j|^2}{2\theta_\alpha^2} - \frac{|I_i - I_j|^2}{2\theta_\beta^2}) + \omega^{(2)} exp(-\frac{|p_i - p_j|^2}{2\theta_\gamma^2})$$
$$(3)$$

Because fully-connected CRF considers the pairwise potentials over all pairs of pixels in the image, minimising the energy function exactly is intractable. Therefore, a mean-field approximation is employed for approximating maximum posterior marginal inference. The CRF distribution $P(X)$ is approximated by a simpler distribution $Q(X)$ by minimizing the KL-divergence $D(Q||P)$. This can be written as the product of independent marginal distributions, i.e. $Q(X) = \prod_i Q_i(X_i)$.

In CRF-RNN, one iteration of the mean-field algorithm can be formulated as a stack of common CNN layers, as shown in algorithm 1. Multiple mean-field iterations can be implemented by repeating the above stack of layers. In other words, the repeated mean-field inference can be formulated as a Recurrent Neural Network(RNN). Then this CRF-RNN layer can be inserted as a part of the deep neural network after FCN-8s. During the training process, the error differentials of

CRF-RNN can be passed to FCN-8s during backward propagation, so that the FCN-8s can generate better unary values for CRF-RNN optimization during the forward propagation. Meanwhile, the CRF parameters, such as the weights of the label compatibility function and Gaussian kernels, can be learned.

---

**Algorithm 1:** Formulate one mean-field iteration as a stack of common CNN layers in [23]. The red annotations are the CNN layers related to corresponding steps in the mean-file iteration.

1   $Q_i \leftarrow \dfrac{1}{Z_i} exp(U_i(l))$ for all $i$

2              $\triangleright$ Initialization, <span style="color:red">$U$ from FCN-8s</span>

3 **while** *not converged* **do**

4    $\tilde{Q}_i^{(m)}(l) \leftarrow \sum_{j\neq i} k^{(m)}(f_i, f_j)Q_j(l)$ for all $m$

5          $\triangleright$ Message Passing, <span style="color:red">Bilateral layer</span>

6    $\check{Q}_i(l) \leftarrow \sum_m w^{(m)} \tilde{Q}_i^{(m)}(l)$

7          $\triangleright$ Weighting Filter Outputs, <span style="color:red">Convolutional layer</span>

8    $\hat{Q}_i(l) \leftarrow \sum_{l' \in L} \mu(l, l') \check{Q}_i(l)$

9          $\triangleright$ Compatibility Transform, <span style="color:red">Convolutional layer</span>

10   $\check{Q}_i(l) \leftarrow U_i(l) - \hat{Q}_i(l)$

11         $\triangleright$ Adding Unary Potentials, <span style="color:red">Concatenated layer</span>

12   $Q_i \leftarrow \dfrac{1}{Z_i} exp(\check{Q}_i(l))$

13         $\triangleright$ Normalizing <span style="color:red">Softmax layer</span>

14 **end**

---

### D. 3D label refinement

Following [6][8], Bayesian update is employed to fuse label hypotheses from the semantic point clouds in different views. Each voxel in a semantic point clouds stores the label information and the corresponding discrete probability. Using the camera projection and visual odometry, voxels from different viewpoints can be transformed to a common coordinate frame. This enables us to update the voxel's label probability distribution by means of a recursive Bayesian update, as shown in equation 4.

$$P(x = l_i | I_{1,...,k}) = \frac{1}{Z} P(x = l_i | I_{1,...,k-1}) P(x = l_i | I_k) \quad (4)$$

where $l_i$ is the label prediction, $I_k$ is the $k^{th}$ image and $Z$ is a constant for distribution normalization.

## IV. EXPERIMENTS

In this section, we firstly introduce data preprocessing on a public material dataset MINC[10]. Next the pipeline of network training is described. Finally, we present qualitative and quantitative evaluations of two different experiments: 2D material recognition on the MINC dataset, and 3D semantic reconstruction in a real-world application, respectively.

### A. Data preprocessing

The large-scale public material dataset Materials in Context (MINC)[10] is employed for training our neural network. MINC is diverse and well-sampled across 23 categories, including wood, glass, metal, brick, fabric and others. There are two kinds of human-annotated data in MINC: small RGB patches with a corresponding class labels(Fig.2.(a)) and images which have been partially pixel-wise labeled at the object level(Fig.2.(b)).

Unfortunately, neither of these annotations can be used in our applications. In RGB patches, there are many non-values (e.g. grey parts in Fig.2.(a)) because these patches extend beyond the limits of the image border. During the fine-tuning process, those non-values give a strong erroneous supervision to the VGG-16 network which lacks normalization layers. This prevents the VGG-16 network from converging. On the other hand, in the partially pixel-wise labeled images, all background pixels are masked, and only one foreground object is labeled, thus losing all context information for CRF-RNN training.

Some data preprocessing is required before network training. The original images from MINC are resized to 500×500 images so that semantic segmentation, based on CNN, can be performed in real-time. Next, 256×256 patches (which have one kind of material in the center) are extracted from the 500×500 images, as shown in Fig.2.(c). This ensures that there are no non-values in the extracted patches.

Next, all partially pixel-wise labeled images, belonging to a single original image, are combined together to generate a single fully pixel-wise labeled image, Fig.2.(d). Because not all objects are labeled in the original image, it is not possible to generate 100% pixel-wise labeled images. Therefore, any unlabeled pixels and repeated labeled pixels(object edges) are labeled as 255, which can then be ignored during the training process. Finally the pixel-wise labeled images are resized to 500×500.

82,1092 patches with the class labels for training, and 96,747 patches with the class labels for testing were generated. 1,498 pixel-wise labeled images for training and 300 pixel-wise labeled images for testing were generated.
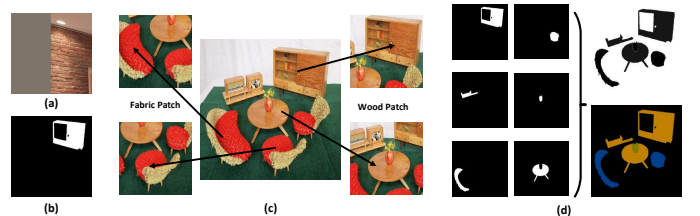


Fig. 2: **Data preprocessing:** (a) A patch with non-values in MINC. (b) A partially pixel-wise labeled image in MINC. (c) Extracting new patches from original images in MINC. (d) Combining the partially pixel-wise labeled images to generate a fully pixel-wise labeled image.

## B. Network training

We initialised our network weights using the publicly available weights of VGG-16[26] model, pre-trained on ImageNet. However, the VGG-16 network is designed for the ImageNet challenge, which can classify 1000 different class labels. In contrast, for the MINC database, only 23 kinds of material need to be classified. Therefore we change the output number of inner-product layer fc8-minc to 23 instead of 1000. Next we fine-tune this network by using the newly extracted $256\times256$ patches from MINC.

Because small patches have better spatial resolution, while large patch have more contextual information, the choice of patch size for fine-tuning is a trade-off between spatial resolution and contextual information. Four different patch sizes are tested in our fine-tuning experiments. As shown in Table I, the accuracy of fine-tuning initially increases, but then decreases, as patch size grows. The highest (optimal) value is achieved when the patch size occupies around 30-50% of the original image.

| Patch size | $56\times56$ | $156\times156$ | $256\times256$ | $356\times356$ |
|---|---|---|---|---|
| Accuracy | 69.20% | 81.06% | 80.18% | 73.40% |

TABLE I: **Accuracy of fine-tuning versus patch size.** Performance is best when the patch size occupies around 30-50% of the original image.

After fine-tuning, the fully connected VGG-16 network is transplanted to a fully convolutional VGG-16 network. The last three inner product layers (fc6, fc7 and fc8-minc) are transformed to the convolutional layers (fc6-conv, fc7-conv and fc8-conv). Convolutional layers inherit the weights of the inner product layers.

Using the fine-tuning model, the FCN-32s, FCN-16s and FCN-8s are trained step by step using 1,498 fully pixel-wise labeled images. Next, the CRF-RNN layer is inserted to form the part of the network following the FCN-8s. After inheriting the learned weights, this end-to-end FCN-8s with CRF-RNN network is trained again using 1,498 fully pixel-wise labeled images. The parameters of each trained network are shown in Table II. The number of mean-field iterations $T$ in the CRF-RNN is set to 5 during the training process. This helps avoid vanishing gradient problems and reduces the training time. During the test process, the number of mean-field iterations can be kept at 5 or be increased to 10 depending on the run-time required.

| | Learning rate | Momentum | Batch size | Weight decay | Training data |
|---|---|---|---|---|---|
| Fine-tuning | 1e-4 reduction with 0.1 | 0.95 | 50 | 0.0005 | 256*256 RGB patch |
| FCN-32s | 1e-10 | 0.99 | 1 | 0.0005 | 500*500 RGB image |
| FCN-16s | 1e-12 | 0.99 | 1 | 0.0005 | 500*500 RGB image |
| FCN-8s | 1e-14 | 0.99 | 1 | 0.0005 | 500*500 RGB image |
| FCN-8s with CRF-RNN | 1e-12 | 0.99 | 1 | 0.0005 | 500*500 RGB image |

TABLE II: The parameters of trained network.

## C. 2D material recognition

We evaluated our trained network using 300 fully pixel-wise labeled images from MINC for 2D pixel-wise material recognition.

*1) The qualitative analysis:* The semantic segmentation results of FCN-8s have non-sharp boundaries because of lacking neighbourhood consistency constraints. After inserting CRF-RNN into the network after FCN-8s, semantic label assignments are significantly improved. As shown in Fig.3, the first and second rows are the original and ground-truth images in MINC. The third and fourth rows are the 2D semantic segmentation results of FCN-8s, and FCN-8s with CRF-RNN, respectively. Clearly the semantic results of FCN-8s with CRF-RNN generate much clear shapes than FCN-8s alone, e.g. table leg in (m), person in (n), sofa in (o), and the chair back and vase in (p). In (l), a large section of "fabric" is erroneously recognised as "carpet". In contrast, this erroneous section is much smaller in (P) because of the neighbourhood consistency constraints of the fully connected CRF optimization.

*2) Quantitative analysis:* The standard parameters for scene understanding evaluation: *pixel accuracy, mean accuracy, mean IU* and *frequency weighed IU* are used for quantitative analysis, as shown in Table.III. End-to-end FCN-8s with CRF-RNN improve 3.53%, 5.16%, 4.62% and 3.92% for *pixel accuracy, mean accuracy, mean IU* and *frequency weighed IU* respectively, as compared to FCN-8s without CRF-RNN. The confusion matrices of material recognition are shown in Fig.4. The colour in the diagonal line is much darker than that in the other positions, suggesting good performance. After combining CRF-RNN with FCN-8s in a united framework, the recognition rate of each class increases by around 4-6%. We attribute this relatively small improvement to the small number (300) of only partially labeled images for this test.

| | Pixel acc. | Mean acc. | Mean IU | f.w. IU |
|---|---|---|---|---|
| FCN-8s | 78.41% | 71.91% | 56.51% | 66.07% |
| FCN-8s with CRF-RNN | 81.94% | 77.07% | 61.13% | 69.99% |

TABLE III: **Quantitative results of 2D material recognition in MINC.** End-to-end FCN-8s with CRF-RNN improve 3.53%, 5.16%, 4.62% and 3.92% for *pixel accuracy, mean accuracy, mean IU* and *frequency weighed IU* respectively, compared with FCN-8s alone.

*3) Run-time performance:* Our experiments were performed using an i7-6800k(3.4Hz) 8-cores CPU and NVIDIA TITAN X GPU (12G). For a $500\times500$ image, the 2D semantic segmentation based on the GPU version of FCN-8s costs 0.13s-0.15s, and that of FCN-8s with CRF-RNN costs 0.4s-0.6s (10 iterations) or 0.2s-0.3s (5 iterations). The run-time greatly decreases if smaller RGB images, e.g. $224\times224$, are used, enabling real-time, or near-to-real-time, pixel-wise material segmentation.

## D. 3D semantic reconstruction

We next evaluate our proposed method in a real-world application. Simultaneous 3D reconstruction and material recog-
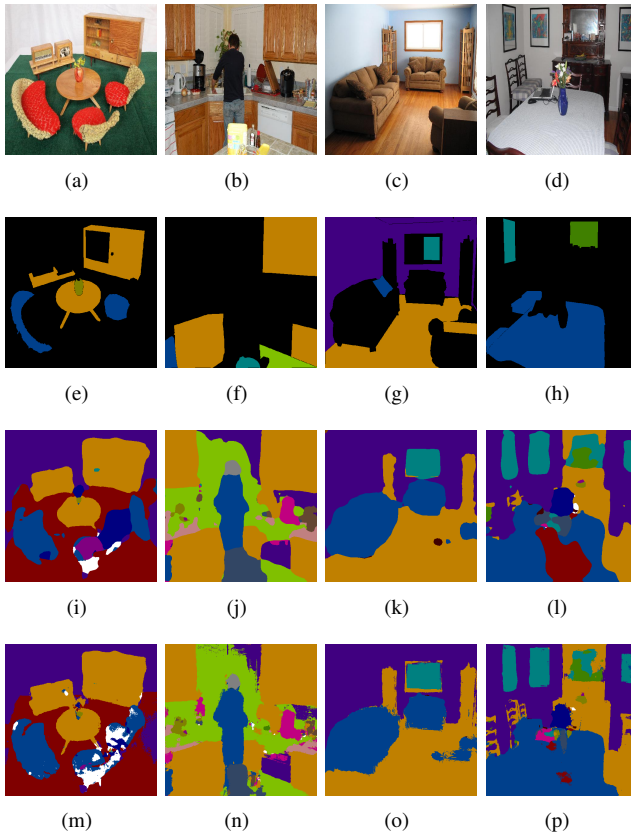
Fig. 3: **Qualitative results of 2D material recognition in MINC:** (a)(b)(c)(d) are original RGB images, (e)(f)(g)(h) are ground truth images, (i)(j)(k)(l) are semantic segmentation results of FCN-8s, (m)(n)(o)(p) are semantic segmentation results of FCN-8s with CRF-RNN. Clearly the semantic results of FCN-8s with CRF-RNN generate much clearer shapes than FCN-8s alone, e.g. table leg in (m), person in (n), sofa in (o), and chair back and vase in (p). In (l), a large part of fabric is erroneously labeled as carpet, while the size of this error greatly decreases in (P) due to the neighbourhood consistency constraints of the fully connected CRF optimization.
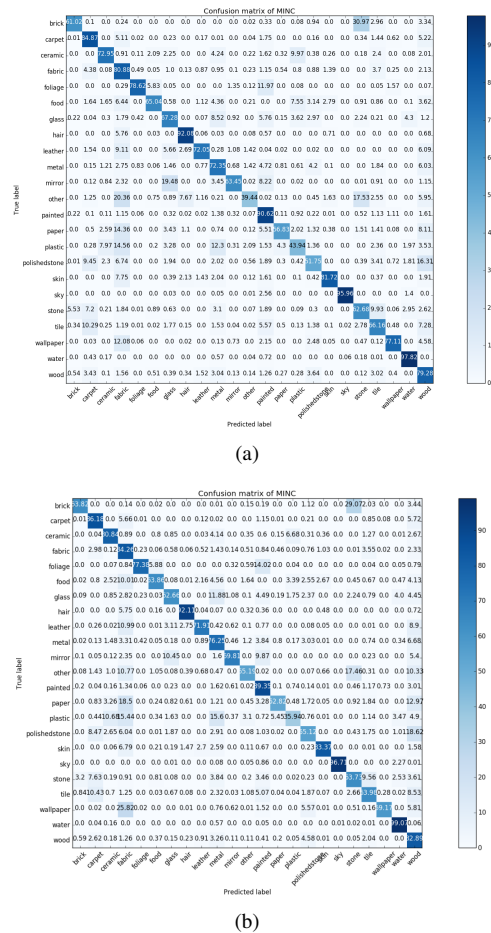


Fig. 4: **Quantitative results of 2D material recognition in MINC.** Confusion matrices of (a)FCN-8s and (b)FCN-8s with CRF-RNN. The colour in the diagonal line is much darker than that in the other positions, suggesting good performance. After combining CRF-RNN with FCN-8s in a united framework, the recognition rate of each class increases by around 4-6%.

nition was performed in a real office which contains many different materials, such as brick, wood, metal, paper, carpet, painted surfaces, and others.

*1) Quantitative analysis:* The qualitative results of each step in our system in a multi-material office are shown in Fig.5. The local/global 3D map and local/global 3D semantic map are shown in Fig. 6 and 7 respectively. It can be seen that most of materials are correctly classified and segmented. However, some small objects cannot be recognised because they do not provide enough pixels in the RGB image. The pixel in the border between two different materials is easily assigned a wrong prediction label. In addition, some errors also result from illumination variances.

*2) Quantitative analysis: Pixel accuracy, mean accuracy, mean IU and frequency weighed IU* are used for quantitative

evaluation. First, 40 key frames of 3D reconstruction in the multi-material office were obtained according to visual odometry from RGBD SLAM. Next, we densely annotated all materials in all key frames using JS Segment Annotator[1]. Finally, pixel-wise true or false numbers were counted between the corresponding pixels from ground-truth and predicted images.

Table IV shows quantitative results. *Pixel accuracy* (80.10%) is satisfactory, while *mean accuracy* (58.75%) appears much lower than that reported in MINC evaluation (76.87%). However, these numbers are misleading, because we only use 40 test samples and there is large variance in material detection rates. Pixel-wise recognition accuracy of some materials e.g. mirror(0%) and paper(6.78%) is very low. However, the mirror only appears in one instance. So just one failure to recognise mirror generates a score of 0%, which

[1]http://kyamagu.github.io/js-segment-annotator/

(a)                             (b)

(c)                             (d)

(e)                             (f)

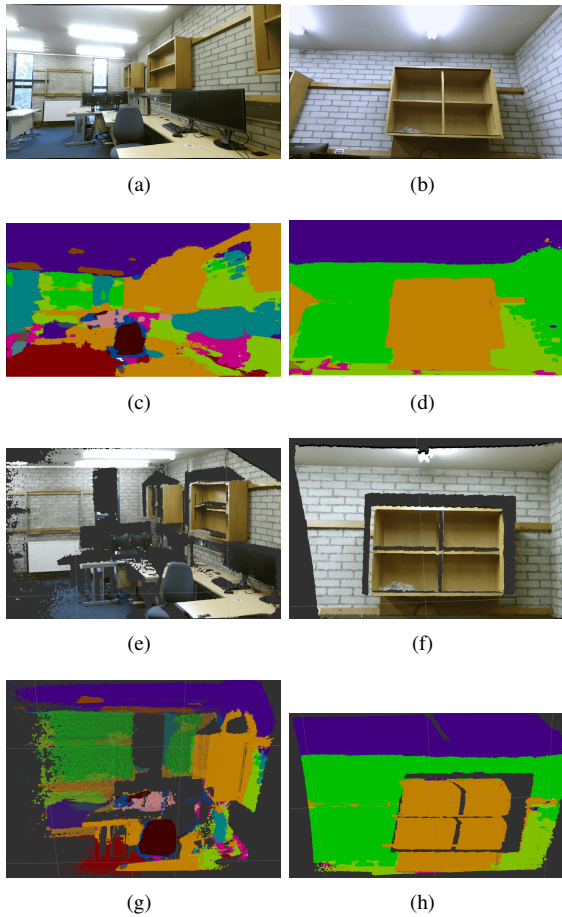(g)                             (h)

Fig. 5: **Qualitative results of 3D semantic reconstruction in a multi-material office**: RGB images from Kinect2 (a)(b), 2D semantic segmentation images (c)(d), 3D point clouds (e)(f), 3D semantic point clouds (g)(h)
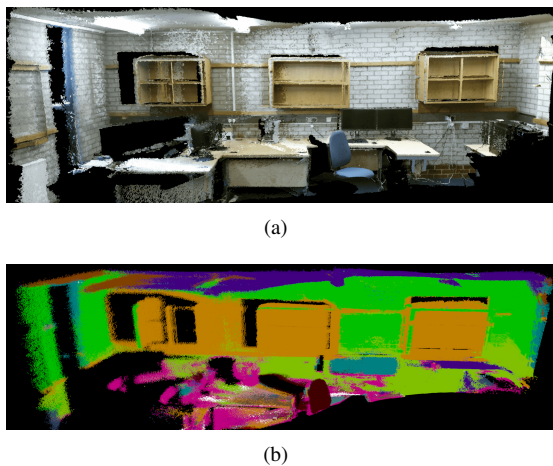


(a)

(b)

Fig. 6: **Qualitative results of 3D semantic reconstruction in a multi-material office:** (a) Local 3D map. (b) Local 3D semantic map.
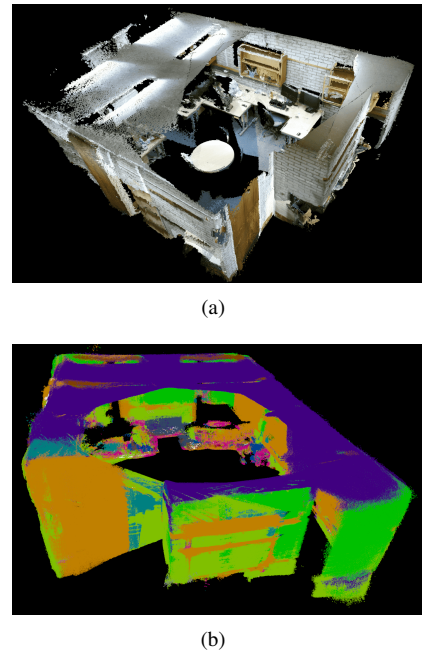


(a)

(b)

Fig. 7: **Qualitative results of 3D semantic reconstruction in a multi-material office:** (a) Global 3D map. (b) Global 3D semantic map.

greatly decreases the overall mean accuracy score. In future, we will obtain more samples, in different scenes, to perform more informative quantitative analysis.

|                         | Pixel acc. | Mean acc. | Mean IU | f.w. IU |
|-------------------------|------------|-----------|---------|---------|
| 3D semantic reconstruction | 80.10%   | 58.75%    | 39.45%  | 68.76%  |

TABLE IV: **Quantitative results of 3D semantic reconstruction in a multi-material office.** The *pixel accuracy*(80.10%) is good, while the *mean accuracy*(58.75%) appears much lower than that(76.87%) of MINC evaluation. However, this is an artifact of the small number of samples (40). Failure to detect just one instance of mirror, causes an accuracy of 0% for that category, which misleadingly skews the overall mean accuracy score to appear low.

*3) Implementation and run-time performance:* We implemented our system on an i7-6800k(3.4Hz) 8-cores CPU and NVIDIA TITAN X GPU (12G). IAI Kinect2 package[2] is employed to interface with ROS and calibrate the Kinect's RGB and depth cameras. The FCN-8s with CRF-RNN is implemented using Caffe toolbox[3]. The overall system is implemented using C++ and GPU programming within a ROS framework.

Run-time performance of our system is around 2Hz (10 iterations) or 4Hz (5 iterations) using the QHD RGB and depth images from Kinect2. The 540×960 RGB image is reduced to 500×500 RGB image for material recognition, and then

[2]https://github.com/code-iai/iai_kinect2/
[3]http://caffe.berkeleyvision.org/

increased to 540×960 RGB image for semantic reconstruction. The run-time performance can be boosted to around 10Hz if the QHD RGB image is decreased to 224×224 RGB image, using 5 CRF iterations for material recognition. In contrast, the run-time performance of SemanticFusion[8] claims up to 25.3Hz using 224×224 RGB image. However, this does not include the significant time needed for CRF post-processing. On average, SemanticFusion takes 20.3s to perform 10 CRF iterations. In contrast, our system is a fully end-to-end system and our run-time includes the CRF optimization which is embedded within our network. For real-time 3D reconstruction, most of the frames are abandoned and only a few key frames are used. So 5Hz-10Hz run-time performance is enough to ensure a real-time semantic reconstruction assuming a 30fps RGB-D camera.

A video demo can be found https://www.youtube.com/watch?v=bVbrb_aE6uw.

## V. Conclusions

In this paper, we report the first system for simultaneous 3D reconstruction and material recognition. It is a real-time, fully end-to-end system, which does not require hand-crafted features or post-processing CRF optimization. Its run-time performance can be boosted to around 10Hz, enabling real-time 3D semantic reconstruction with a 30fps camera. We presented both quantitative and qualitative experimental results, which support the effectiveness of our method.

## Acknowledgment

## References

[1] F. Endres, J. Hess, J. Sturm, D. Cremers, and W. Burgard, "3-D Mapping with an RGB-D camera," *IEEE Transactions on Robotics*, vol. 30, no. 1, pp. 177–187, 2014.

[2] J. Engel, T. Schöps, and D. Cremers, "LSD-SLAM: Large-Scale Direct Monocular SLAM," in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2014, vol. 8690 LNCS, no. PART 2, pp. 834–849.

[3] R. a. Newcombe, A. J. Davison, S. Izadi, P. Kohli, O. Hilliges, J. Shotton, D. Molyneaux, S. Hodges, D. Kim, and A. Fitzgibbon, "KinectFusion: Real-time dense surface mapping and tracking," in *2011 10th IEEE International Symposium on Mixed and Augmented Reality*, 2011, pp. 127–136.

[4] R. F. Salas-Moreno, R. a. Newcombe, H. Strasdat, P. H. J. Kelly, and A. J. Davison, "SLAM++: Simultaneous localisation and mapping at the level of objects," in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. IEEE, jun 2013, pp. 1352–1359.

[5] K. Tateno, F. Tombari, and N. Navab, "When 2.5D is not enough: Simultaneous reconstruction, segmentation and recognition on dense SLAM," in *2016 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, may 2016, pp. 2295–2302.

[6] A. Hermans, G. Floros, and B. Leibe, "Dense 3D semantic mapping of indoor scenes from RGB-D images," in *2014 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, may 2014, pp. 2631–2638.

[7] V. Prisacariu, M. Lidegaard, D. Murray, P. Torr, S. Golodetz, and P. Patrick, "Incremental Dense Semantic Stereo Fusion for Large-Scale Semantic Scene Reconstruction," *Icra*, pp. Prisacariu, V. et al., 2015. Incremental Dense Sem, 2015.

[8] J. McCormac, A. Handa, A. Davison, and S. Leutenegger, "SemanticFusion: Dense 3D Semantic Mapping with Convolutional Neural Networks," *arXiv preprint*, 2016.

[9] J. Degol, M. Golparvar-Fard, and D. Hoiem, "Geometry-Informed Material Recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, jun 2016, pp. 1554–1562.

[10] S. Bell, P. Upchurch, N. Snavely, and K. Bala, "Material recognition in the wild with the Materials in Context Database," in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 07-12-June. IEEE, jun 2015, pp. 3479–3487.

[11] G. Schwartz and K. Nishino, "Material Recognition from Local Appearance in Global Context," nov 2016.

[12] N. Rusk, "A 4D Light-Field Dataset and CNN Architectures for Material Recognition," in *European Conference on Computer Vision*, ser. Lecture Notes in Computer Science, B. Leibe, J. Matas, N. Sebe, and M. Welling, Eds. Cham: Springer International Publishing, 2016, pp. 121–138.

[13] R. F. Salas-Moreno, B. Glocken, P. H. J. Kelly, and A. J. Davison, "Dense Planar SLAM," in *Proc. of ISMAR*. IEEE, sep 2014, pp. 157–164.

[14] H. Noh, S. Hong, and B. Han, "Learning Deconvolution Network for Semantic Segmentation," in *2015 IEEE International Conference on Computer Vision (ICCV)*, vol. 11-18-Dece. IEEE, dec 2015, pp. 1520–1528.

[15] P. Krähenbühl and V. Koltun, "Efficient inference in fully connected crfs with gaussian edge potentials," *Adv. Neural Inf. Process. Syst*, 2011.

[16] C. Liu, L. Sharan, E. H. Adelson, and R. Rosenholtz, "Exploring features in a Bayesian framework for material recognition," in *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. IEEE, jun 2010, pp. 239–246.

[17] D. Hu, L. Bo, and X. Ren, "Toward Robust Material Recognition for Everyday Objects," in *Procedings of the British Machine Vision Conference 2011*. British Machine Vision Association, 2011, pp. 48.1–48.11.

[18] X. Qi, R. Xiao, C.-G. Li, Y. Qiao, J. Guo, and X. Tang, "Pairwise Rotation Invariant Co-Occurrence Local Binary Pattern," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, no. 11, pp. 2199–2213, nov 2014.

[19] G. Schwartz and K. Nishino, "Visual Material Traits: Recognizing Per-Pixel Material Context," in *2013 IEEE International Conference on Computer Vision Workshops*. IEEE, dec 2013, pp. 883–890.

[20] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, and A. Vedaldi, "Describing Textures in the Wild," in *2014 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, jun 2014, pp. 3606–3613.

[21] K. Mcdonald-maier, "Evaluating Deep Convolutional Neural Networks for Material Classification," no. November, 2016.

[22] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, jun 2015, pp. 3431–3440.

[23] S. Zheng, S. Jayasumana, B. Romera-Paredes, V. Vineet, Z. Su, D. Du, C. Huang, and P. H. S. Torr, "Conditional Random Fields as Recurrent Neural Networks," *Proceedings of the IEEE International Conference on Computer Vision*, no. [1] S. Zheng, S. Jayasumana, B. Romera-Paredes, V. Vineet, Z. Su, D. Du, C. Huang, and P. H. S. Torr, Conditional Random Fields as Recurrent Neural Networks, Proc. IEEE Int. Conf. Comput. Vis., pp. 15291537, 2015., pp. 1529–1537, 2015.

[24] R. Kümmerle, G. Grisetti, H. Strasdat, K. Konolige, and W. Burgard, "G2o: A general framework for graph optimization," in *Proceedings - IEEE International Conference on Robotics and Automation*, 2011, pp. 3607–3613.

[25] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe," *Proceedings of the ACM International Conference on Multimedia - MM '14*, pp. 675–678, 2014.

[26] K. Simonyan and A. Zisserman, "Very Deep Convolutional Networks for Large-Scale Image Recognition," *International Conference on Learning Representations*, pp. 1–14, 2015.