

SemanticFusion: Dense 3D Semantic Mapping with Convolutional Neural Networks

John McCormac, Ankur Handa, Andrew Davison, and Stefan Leutenegger
Dyson Robotics Lab, Imperial College London

Abstract—Ever more robust, accurate and detailed mapping using visual sensing has proven to be an enabling factor for mobile robots across a wide variety of applications. For the next level of robot intelligence and intuitive user interaction, maps need to extend beyond geometry and appearance — they need to contain semantics. We address this challenge by combining Convolutional Neural Networks (CNNs) and a state-of-the-art dense Simultaneous Localisation and Mapping (SLAM) system, ElasticFusion, which provides long-term dense correspondences between frames of indoor RGB-D video even during loopy scanning trajectories. These correspondences allow the CNN's semantic predictions from multiple view points to be probabilistically fused into a map. This not only produces a useful semantic 3D map, but we also show on the NYUv2 dataset that fusing multiple predictions leads to an improvement even in the 2D semantic labelling over baseline single frame predictions. We also show that for a smaller reconstruction dataset with larger variation in prediction viewpoint, the improvement over single frame segmentation increases. Our system is efficient enough to allow real-time interactive use at frame-rates of $\approx 25\text{Hz}$.

I. INTRODUCTION

The inclusion of rich semantic information within a dense map enables a much greater range of functionality than geometry alone. For instance, in domestic robotics, a simple fetching task requires knowledge of both what something is, as well as where it is located. As a specific example, a user communicating with a robot with a shared spatial and semantic understanding may issue commands such as ‘fetch the coffee mug from the nearest table on your right.’ Similarly, the ability to query semantic information within a map is useful for humans directly, providing a database for answering spoken queries about the semantics of a previously made map; ‘How many chairs do we have in the conference room? What is the distance between the lectern and its nearest chair?’ In this work, we combine the geometric information from a state-of-the-art SLAM system ElasticFusion [26] with recent advances in semantic segmentation using Convolutional Neural Networks (CNNs).

Our approach is to use the SLAM system to provide correspondences from the 2D frame into a globally consistent 3D map. This allows the CNN's semantic predictions from multiple viewpoints to be probabilistically fused into a dense semantically annotated map, as shown in Figure 1. ElasticFusion is particularly suitable for fusing semantic labels because its surfel-based surface representation is automatically deformed to remain consistent after the small and large loop closures which would frequently occur during typical interactive use by an agent (whether human or robot). As the surface representation is deformed and corrected, individual

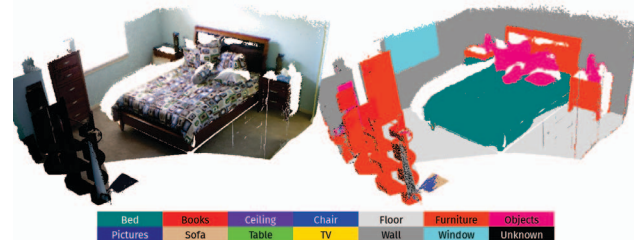


Fig. 1: **The output of our system:** On the left, a dense surfel based reconstruction from a video sequence in the NYUv2 test set. On the right the same map, semantically annotated with the classes given in the legend below.

surfels remain persistently associated with real-world entities and this enables long-term fusion of per-frame semantic predictions over wide changes in viewpoint. The geometry of the map itself can also provide useful information which can be used to regularise the final predictions.

Our pipeline is designed to work online, and although we have not focused on performance, the efficiency of each component leads to a real-time capable ($\approx 25\text{Hz}$) interactive system. The resulting map could also be used as a basis for more expensive offline processing to further improve both the geometry and the semantics; however that has not been explored in the current work.

We evaluate the accuracy of our system on the NYUv2 dataset, and show that by using information from the unlabelled raw video footage we can improve upon baseline approaches performing segmentation using only a single frame. This suggests the inclusion of SLAM not only provides an immediately useful semantic 3D map, but also that many state-of-the-art 2D single frame semantic segmentation approaches may see a boost in performance when combined with SLAM.

The NYUv2 dataset was not taken with full room reconstruction in mind, and often does not provide significant variations in viewpoints for a given scene. To explore the benefits of SemanticFusion within a more thorough reconstruction, we developed a small dataset of a reconstructed office room, annotated with the NYUv2 semantic classes. Within this dataset we witness a more significant improvement in segmentation accuracy over single frame 2D segmentation. This indicates that the system is particularly well suited to longer duration scans with wide viewpoint variation aiding to disambiguate the single-view 2D semantics.

II. RELATED WORK

The works most closely related are Stücker *et al.* [23] and Hermans *et al.* [7]; both aim towards a dense semantically annotated 3D map of indoor scenes. They both obtain per-pixel label predictions for incoming frames using **Random Decision Forests**, whereas ours exploits recent advances in Convolutional Neural Networks that provide state-of-the-art accuracy with a real-time capable run-time performance. They both fuse predictions from different viewpoints in a classic Bayesian framework. Stücker *et al.* [23] used a Multi-Resolution Surfel Map-based SLAM system capable of operating at 12.8Hz, however unlike our system they do not maintain a single global semantic map as local key frames store aggregated semantic information and these are subject to graph optimisation in each frame. Hermans *et al.* [7] did not use the capability of a full SLAM system with explicit loop closure: they registered the predictions in the reference frames using only camera tracking. Their run-time performance was 4.6Hz, which would prohibit processing a live video feed, whereas our system is capable of operating online and interactively. As here, they explore regularising their predictions using Krähenbühl and Koltun's [12] fully-connected CRF inference scheme to obtain a final semantic map.

Previous work by Salas-Moreno *et al.* aimed to create a fully capable SLAM system, SLAM++ [19], which maps indoor scenes at the level of semantically defined objects. However, their method is limited to mapping objects that are present in a pre-defined database. It also does not provide the dense labelling of entire scenes that we aim for in this work, which also includes walls, floors, doors, and windows which are equally important to describe the extent of the room. Additionally, the features they use to match template models are hand-crafted unlike our CNN features that are learned in an end-to-end fashion with large training datasets.

The majority of other approaches to indoor semantic labelling either focuses on offline batch mapping methods [24], [11] or on single-frame 2D segmentations which do not aim to produce a semantically annotated 3D map [3], [20], [15], [22]. Valentin *et al.* [24] used a CRF and a per-pixel labelling from a variant of TextonBoost to reconstruct semantic maps of both indoor and outdoor scenes. This produces a globally consistent 3D map, however inference is performed on the whole mesh once instead of incrementally fusing the predictions online. Koppula *et al.* [11] also tackle the problem on a completed 3D map, forming segments of the map into nodes of a graphical model and using hand-crafted geometric and visual features as edge potentials to infer the final semantic labelling. In outdoor semantic labelling, Vineet *et al.* [25] produced an incremental 3D reconstruction from stereo pairs and used a Random Forest with a CRF for semantic labelling. However they also did not have a full SLAM system capable of loop closure while maintaining a globally consistent map structure.

Our semantic mapping pipeline is inspired by the recent success of Convolution Neural Networks in semantic

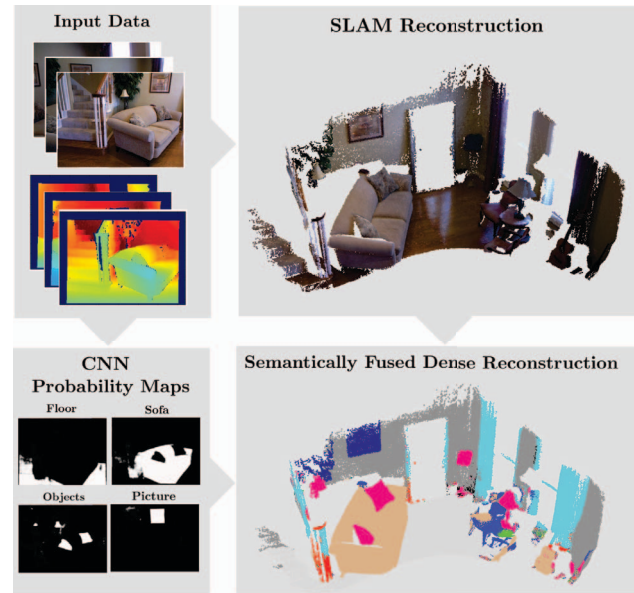


Fig. 2: **An overview of our pipeline:** Input images are used to produce a SLAM map, and a set of probability prediction maps (here only four are shown). These maps are fused into the final dense semantic map via Bayesian updates.

labelling and segmentation tasks [13], [16], [17]. CNNs have proven capable of both state-of-the-art accuracy and efficient test-time performance. They have exhibited these capabilities on numerous datasets and a variety of data modalities, in particular RGB [17], [16], Depth [1], [6] and Normals [2], [4], [5]. In this work we build on the CNN model proposed by Noh *et al.* [17], but we modify it to take advantage of the directly available depth data in a manner that does not require significant additional pre-processing.

III. METHOD

Our SemanticFusion pipeline is composed of three separate units; a real-time SLAM system ElasticFusion, a Convolutional Neural Network, and a Bayesian update scheme, as illustrated in Figure 2. The role of the SLAM system is to provide correspondences between frames, and a globally consistent map of fused surfels. Separately, the CNN receives a 2D image (for our architecture this is RGB or RGBD, for Eigen *et al.* [2] it also includes normals), and returns a set of per-pixel class probabilities. Finally, a Bayesian update scheme keeps track of the class probability distribution for each surfel, and uses the correspondences provided by the SLAM system to update those probabilities based on the CNN's predictions. Finally, we also experiment with a CRF regularisation scheme to use the geometry of the map itself to improve the semantic predictions [7], [12]. The following section outlines each of these components in more detail.

A. SLAM Mapping

We choose ElasticFusion as our SLAM system.¹ For each arriving frame, k , ElasticFusion tracks the camera pose

¹Available on <https://github.com/mp3guy/ElasticFusion>

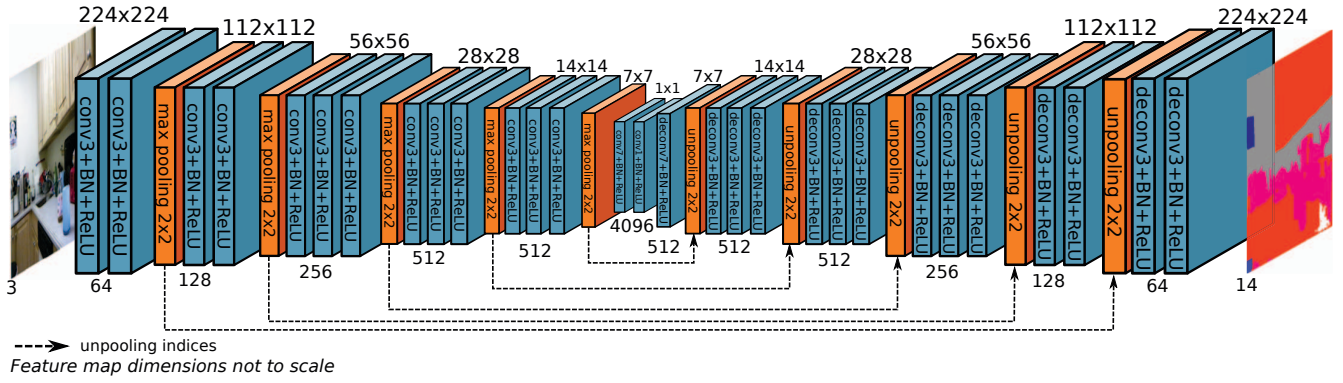


Fig. 3: **CNN Architecture:** RGB-CNN of Noh *et al.* [17] used in our experiments. “conv3” denotes an 3×3 kernel size, with unit border padding. BN denotes batch normalisation. For RGBD experiments, the input actually has 4 channels.

via a combined ICP and RGB alignment, to yield a new pose T_{WC} , where W denotes the World frame and C the camera frame. New surfels are added into our map using this camera pose, and existing surfel information is combined with new evidence to refine their positions, normals, and colour information. Additional checks for a loop closure event run in parallel and the map is optimised immediately upon a loop closure detection. loop closure : refine the map

The deformation graph and surfel based representation of ElasticFusion lend themselves naturally to the task at hand, allow probability distributions to be ‘carried along’ with the surfels during loop closure, and also fusing new depth readings to update the surfel’s depth and normal information, without destroying the surfel, or its underlying probability distribution. It operates at real-time frame-rates at VGA resolution and so can be used both interactively by a human or in robotic applications. We maintained many of the default parameters in the public implementation, The depth cutoff was extended from 3m to 8m to allow reconstruction on sequences with geometry outside of the 3m range, and we disabled the RGB component of tracking on the NYUv2 dataset due to a white border produced from preprocessing.

B. CNN Architecture

Our CNN is implemented in *caffe* [10] and adopts the Deconvolutional Semantic Segmentation network architecture proposed by Noh *et al.* [17], depicted in Figure 3. Their architecture is itself based on the VGG 16-layer network [21], but with the addition of max unpooling and deconvolutional layers which are trained to output a dense pixel-wise semantic probability map. This CNN was originally trained for RGB input, and in the following sections when using a network with this setup we describe it RGB-CNN.

Given the availability of depth data, we modified the original network architecture to accept depth information as a fourth channel. Unfortunately, the depth modality lacks the large scale training datasets of its RGB counterpart. The NYUv2 dataset only consists of 795 labelled training images. To effectively use depth, we initialized the depth filters with

the average intensity of the other three inputs, which had already been trained on a large dataset, and converted it from the 0–255 colour range to the 0–8m depth range by increasing the weights by a factor of $\approx 32 \times$.

We rescale incoming images to the native 224×224 resolution for our CNNs, using bilinear interpolation for RGB, and nearest neighbour for depth. In our experiments with the Eigen *et al.* implementation we rescale the inputs in the same manner to 320×240 resolution. We upsample the network output probabilities to full 640×480 image resolution using nearest neighbour when applying the update to surfels, described in the section below.

C. Incremental Semantic Label Fusion

In addition to normal and location information, each surfel (index s) in our map, \mathcal{M} , stores a discrete probability distribution, $P(L_s = l_i)$ over the set of class labels, $l_i \in \mathcal{L}$. Each newly generated surfel is initialised with a uniform distribution over the semantic classes, as we begin with no *a priori* evidence as to its latent classification.

After a prespecified number of frames, we perform a forward pass of the CNN with the image I_k coming directly from the camera. Depending on the CNN architecture, this image can include any combination of RGB, depth, or normals. Given the data I_k of the k^{th} image, the output of the CNN is interpreted in a simplified manner as a per-pixel independent probability distribution over the class labels $P(O_u = l_i | I_k)$, with u denoting pixel coordinates.

Using the tracked camera pose T_{WC} , we associate every surfel at a given 3D location ${}_W x(s)$ in the map, with pixel coordinates u via the camera projection $u(s, k) = \pi(T_{CW}(k) {}_W x(s))$, employing the homogeneous transformation matrix $T_{CW}(k) = T_{WC}^{-1}(k)$ and using homogeneous 3D coordinates. This enables us to update all the surfels in the visible set $\mathcal{V}_k \subseteq \mathcal{M}$ with the corresponding probability distribution by means of a recursive Bayesian update

$$P(l_i | I_{1,\dots,k}) = \frac{1}{Z} P(l_i | I_{1,\dots,k-1}) P(O_{u(s,k)} = l_i | I_k), \quad (1)$$

which is applied to all label probabilities per surfel, finally normalising with constant Z to yield a proper distribution.

It is the SLAM correspondences that allow us to accurately associate label hypotheses from multiple images and combine evidence in a Bayesian way. The following section discusses how the naïve independence approximation employed so far can be mitigated, allowing semantic information to be propagated spatially when semantics are fused from different viewpoints.

D. Map Regularisation

We explore the benefits of using map geometry to regularise predictions by applying a fully-connected CRF with Gaussian edge potentials to surfels in the 3D world frame, as in the work of Hermans *et al.* [7], [12]. We do not use the CRF to arrive at a final prediction for each surfel, but instead use it incrementally to update the probability distributions. In our work, we treat each surfel as a node in the graph. The algorithm uses the mean-field approximation and a message passing scheme to efficiently infer the latent variables that approximately minimise the Gibbs energy E of a labelling, \mathbf{x} , in a fully-connected graph, where $x_s \in \{l_i\}$ denotes a given labelling for the surfel with index s .

The energy $E(\mathbf{x})$ consists of two parts, the unary data term $\psi_u(x_s)$ is a function of a given label, and is parameterised by the internal probability distribution of the surfel from fusing multiple CNN predictions as described above. The pairwise smoothness term, $\psi_p(x_s, x_{s'})$ is a function of the labelling of two connected surfels in the graph, and is parameterised by the geometry of the map:

$$E(\mathbf{x}) = \sum_s \psi_u(x_s) + \sum_{s < s'} \psi_p(x_s, x_{s'}). \quad (2)$$

For the data term we simply use the negative logarithm of the chosen labelling's probability for a given surfel,

$$\psi_u(x_s) = -\log(P(L_s = x_s | \mathbf{I}_{1,\dots,k})). \quad (3)$$

In the scheme proposed by Krähenbühl and Koltun [12] the smoothness term is constrained to be a linear combination of K Gaussian edge potential kernels, where \mathbf{f}_s denotes some feature vector for surfel, s , and in our case $\mu(x_s, x_{s'})$ is given by the Potts model, $\mu(x_s, x_{s'}) = [x_s \neq x_{s'}]$:

$$\psi_p(x_s, x_{s'}) = \mu(x_s, x_{s'}) \left(\sum_{m=1}^K w^{(m)} k^{(m)}(\mathbf{f}_s, \mathbf{f}_{s'}) \right). \quad (4)$$

Following previous work [7] we use two pairwise potentials; a bilateral appearance potential seeking to closely tie together surfels with both a similar position and appearance, and a spatial smoothing potential which enforces smooth predictions in areas with similar surface normals:

$$k^1(\mathbf{f}_s, \mathbf{f}_{s'}) = \exp \left(-\frac{|\mathbf{p}_s - \mathbf{p}_{s'}|^2}{2\theta_\alpha^2} - \frac{|\mathbf{c}_s - \mathbf{c}_{s'}|^2}{2\theta_\beta^2} \right), \quad (5)$$

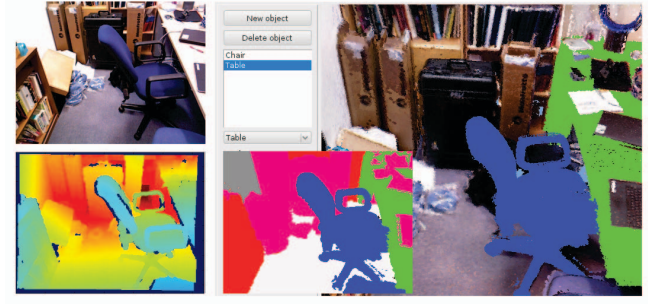


Fig. 4: **Our office reconstruction dataset:** On the left are the captured RGB and Depth images. On the right, is our 3D reconstruction and annotation. Inset into that is the final ground truth rendered labelling we use for testing.

$$k^2(\mathbf{f}_s, \mathbf{f}_{s'}) = \exp \left(-\frac{|\mathbf{p}_s - \mathbf{p}_{s'}|^2}{2\theta_\alpha^2} - \frac{|\mathbf{n}_s - \mathbf{n}_{s'}|^2}{2\theta_\gamma^2} \right). \quad (6)$$

The gaussian edge potentials allow for an efficient mean field approximation algorithm for inference even in a fully connected CRF. The computational cost of this algorithm is linear in the number of surfels, which is particularly useful in our case as the SLAM system can potentially enable long trajectories and millions of surfels.

We chose unit standard deviations of $\theta_\alpha = 0.05\text{m}$ in the spatial domain, $\theta_\beta = 20$ in the RGB colour domain, and $\theta_\gamma = 0.1$ radians in the angular domain. We experimented with varying these parameters on the reconstruction dataset below, but this did not lead to any noticeable improvement. We also maintained w^1 of 10 and w^2 of 3 for all experiments. These were the default settings in Krähenbühl and Koltun's public implementation² [12].

IV. EXPERIMENTS

A. Network Training

We initialise our CNNs with weights from Noh *et al.* [17] trained for segmentation on the PASCAL VOC 2012 segmentation dataset [3]. For depth input we initialise the fourth channel as described in Section III-B, above. We finetuned this network on the training set of the NYUv2 dataset for the 13 semantic classes defined by Couprie *et al.* [1]. The NYUv2 training images contain depth inpainted with the colorization scheme of Levin *et al.* [14] and we continue to preprocess depth in this manner in the experiments below.

For optimisation we used standard stochastic gradient descent, with a learning rate of 0.01, momentum of 0.9, and weight decay of 5×10^{-4} . After 10k iterations we reduced the learning rate to 1×10^{-3} . We use a mini-batch size of 64, and trained the networks for a total of 20k iterations over the course of 2 days on an Nvidia GTX Titan X.

²Available from: <http://www.philkr.net/home/densecrf>

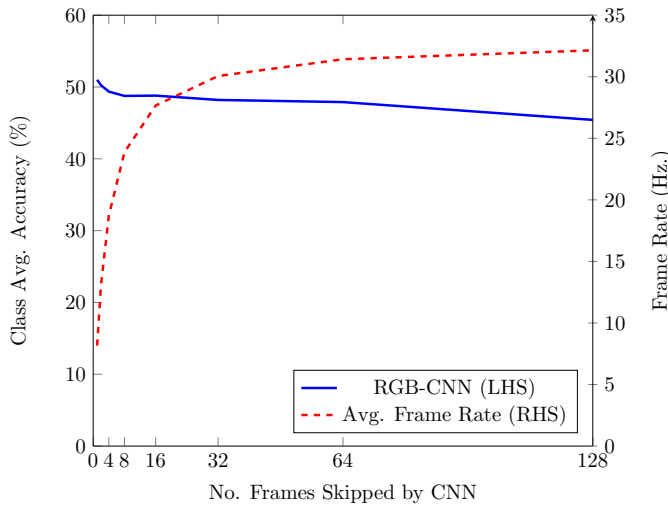


Fig. 5: The class average accuracy of our RGB-CNN on the office reconstruction dataset against the number of frames skipped between fusing semantic predictions. We perform this evaluation without CRF smoothing. The right hand axis shows the estimated run-time performance in terms of FPS.

B. Reconstruction Dataset

We produced a small experimental RGB-D reconstruction dataset, which aimed for a relatively complete reconstruction of an office room. The trajectory used is notably more loopy, both locally and globally, than the NYUv2 dataset which typically consists of a single back and forth sweep. We believe the trajectory in our dataset is more representative of the scanning motion an active agent may perform when inspecting a scene.

We also took a different approach to manual annotation of this data, by using a 3D tool we developed to annotate the surfels of the final 3D reconstruction with the 13 NYUv2 semantic classes under consideration (only 9 were present). We generated 2D projections of our 3D annotations using the camera pose trajectory from the SLAM system, enabling us to render ground truth labels for any frame in the input video sequence. This approach was much more efficient than producing manual 2D single-frame annotations, and produces a more temporally consistent ground truth. The tool, and the resulting annotations are depicted in Figure 4. Every 100th frame of the sequence was used as a test sample to validate our predictions against the annotated ground truth, resulting in 49 test frames.

C. CNN and CRF Update Frequency Experiments

We used the dataset to evaluate the accuracy of our system when only performing a CNN prediction on a subset of the incoming video frames. We used the RGB-CNN described above, and evaluated the accuracy of our system when performing a prediction on every 2^n frames, where $n \in \{0..7\}$. We calculate the average frame-rate based upon the run-time analysis discussed in Section IV-F. As shown in Figure 5, the accuracy is highest (50.8%) when every

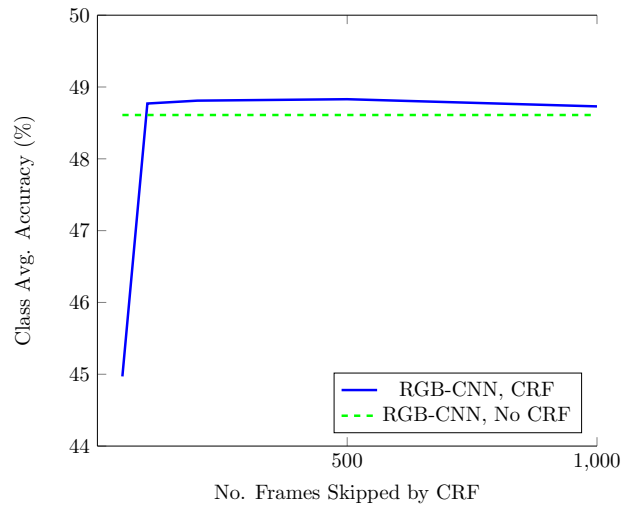


Fig. 6: The average class accuracy processing every 10th frame with a CNN, with a variable number of frames between CRF updates. If applied too frequently the CRF was detrimental to performance, and the performance improvement from the CRF was not significant for this CNN.

frame is processed by the network, however this leads to a significant drop in frame-rate to 8.2Hz. Processing every 10th frame results in a slightly reduced accuracy (48.6%), but over three times the frame-rate of 25.3Hz. This is the approach taken in all of our subsequent evaluations.

We also evaluated the effect of varying the number of frames between CRF updates (Figure 6). We found that when applied too frequently the CRF resulted in a significant reduction in accuracy. Performing an update every 500 frames results in a slight improvement, and so we use that as the default update rate in all subsequent experiments.

D. Accuracy Evaluation

We evaluate the accuracy of our SemanticFusion pipeline against the accuracy achieved by a single frame CNN segmentation. Here we give two accuracy metrics; pixel average accuracy, the proportion of correctly classified pixels out of all ground truth labelled pixels, and class average accuracy, which is the average of the diagonal of the prediction's normalised confusion matrix. The results of this evaluation on the reconstruction dataset are summarised in Table I. We observe that in all cases semantically fusing additional viewpoints improved the accuracy of the segmentation over a single frame system. For the RGB-CNN, performance improved from 39.4% for a single frame to 48.6% when projecting the predictions from the 3D SemanticFusion map. The RGBD-CNN, also saw a marked improvement, from 43.6% for a single frame to 48.3% with SemanticFusion.

We also evaluate our system on the office dataset when using predictions from the state-of-the-art CNN developed by Eigen *et al.*³ based on the VGG architecture. To maintain

³We use the publicly available network weights and implementation from: <http://www.cs.nyu.edu/~deigen/dnl/>.

Office Reconstruction: 13 Class Semantic Segmentation

Method	books	ceiling	chair	floor	objects	painting	table	wall	window	class avg.	pixel avg.
RGB	37.7	55.0	17.5	61.7	48.0	26.6	9.0	76.4	22.4	39.4	49.0
RGB-SF	60.3	81.4	30.1	65.2	43.9	29.5	9.8	81.0	36.3	48.6	55.0
RGB-SF-CRF	60.5	81.5	30.0	65.1	44.4	30.0	9.9	81.4	36.0	48.8	55.3
RGBD	61.8	48.2	28.6	63.9	41.8	39.5	9.1	80.6	18.9	43.6	47.0
RGBD-SF	66.4	78.7	36.8	63.4	41.9	26.2	12.1	84.2	25.3	48.3	54.7
RGBD-SF-CRF	66.4	78.1	37.2	64.2	40.8	27.5	10.6	85.1	22.7	48.1	54.8
Eigen [2]	57.8	54.3	57.8	72.8	49.4	77.5	24.1	81.6	38.9	57.1	62.5
Eigen-SF	60.8	58.0	62.8	74.9	53.3	80.3	24.6	86.3	38.8	60.0	65.8
Eigen-SF-CRF	65.9	53.3	65.1	76.8	53.1	79.6	22.0	87.7	41.4	60.5	67.0

TABLE I: **Reconstruction dataset results:** SF denotes that the labels were produced by SemanticFusion, and the results were captured immediately if a frame with ground truth labelling was present. When no reconstruction is present for a pixel, we fall back to the predictions of the baseline single frame network. All accuracy evaluations were performed at 320×240 resolution.

consistency with the rest of the system, we perform only a single forward pass of the network to calculate the output probabilities. The network requires surface normal information, and so to ensure the input pipeline is the same as in Eigen *et al.* [2], we preprocess the sequence with the MATLAB script linked to in the project page to produce normals from the current depth frame. The requirement of this preprocessing step prohibits using the Eigen *et al.* CNN live in real-time, but we include it here to show that even a state-of-the-art CNN with additional useful input channels such as normal information, can still benefit from fusing predictions from multiple viewpoints. With this setup we see an improvement of 2.9% over the single frame implementation with SemanticFusion, from 57.1% to 60.0%.

The performance benefit of the CRF was less clear. It provided a very small improvement of +0.5% for the Eigen network and +0.2% for the RGB-CNN, but a slight detriment to the RGBD-CNN of -0.2%.

E. NYU Dataset

We choose to validate our approach on the NYUv2 dataset [20], as it is one of the few datasets which provides all of the information required to evaluate semantic RGB-D reconstruction. The SUN RGB-D [22], although an order of magnitude larger than NYUv2 in terms of labelled images, does not provide the raw RGB-D videos and therefore could not be used in our evaluation.

The NYUv2 dataset itself is still not ideally suited to the role. Many of the 206 test set video sequences exhibit significant drops in frame-rate and thus prove unsuitable for tracking and reconstruction. In our evaluations we excluded any sequence which experienced a frame-rate under 2Hz. The remaining 140 test sequences result in 360 labelled test images of the original 654 image test set in NYUv2. The results of our evaluation are presented in Table II and some qualitative results are shown in Figure 7.

Overall, fusing semantic predictions resulted in a notable improvement over single frame predictions. However, the

total relative gains of 2.3% for the RGBD-CNN was approximately half of the 4.7% improvement witnessed in the office reconstruction dataset. We believe this is largely a result of the style of capturing NYUv2 datasets. The primarily rotational scanning pattern often used in test trajectories does not provide as many useful different viewpoints from which to fuse independent predictions. Despite this, there is still a significant accuracy improvement over the single frame predictions. The RGBD-CNN performed better than the RGB-CNN both in the baseline (+3.3%), and after SemanticFusion (+3.7%).

We also improved upon the state-of-the-art Eigen *et al.* [2] CNN, with the class average accuracy going from 59.9% to 63.2% (+3.3%). This result clearly shows, even on this challenging dataset, the capacity of SemanticFusion to not only provide a useful semantically annotated 3D map, but also to improve the predictions of state-of-the-art 2D semantic segmentation systems. It is also interesting to note that although better overall, the Eigen Multi-Scale CNN produced quite different accuracies for individual classes when compared against our CNNs. It performed particularly well on tv, table, and ceiling while under performing in bed, books, and objects. The differences in the Eigen CNN architecture, input channels, and training procedure make it difficult to speculate on why this may be the case without further detailed experiments, however it does suggest that the networks may complement each other and a hybrid approach could improve performance further.

We also give the accuracy results of the similar work of Hermans *et al.* [7] which used Random Decision Forests. It can be seen that the accuracy of even the baseline CNN approach is superior and the difference when using SemanticFusion is greater still. It is however difficult to draw precise comparisons as we exclude sequences with less than a 2Hz frame-rate, while they did not.

The improvement as a result of the CRF was not particularly significant, but slightly positive for all CNNs. Eigen's CNN saw +0.4% improvement and the RGBD-CNN saw

NYUv2 Test Set: 13 Class Semantic Segmentation

Method	bed	books	ceiling	chair	floor	furniture	objects	painting	sofa	table	tv	wall	window	class avg.	pixel avg.
RGB	61.1	52.0	27.9	44.3	93.9	59.7	60.8	68.0	30.2	21.0	14.7	86.0	60.4	52.3	62.2
RGB-SF	60.4	49.8	32.5	50.8	92.8	65.4	60.6	62.5	43.6	30.5	18.8	85.8	63.9	55.2	67.2
RGB-SF-CRF	60.5	49.7	32.7	50.9	93.0	65.5	60.5	62.6	43.6	30.0	18.8	85.9	64.0	55.3	67.2
RGBD	62.5	60.5	35.0	51.7	92.1	54.5	61.3	72.1	34.7	26.1	32.4	86.5	53.5	55.6	62.0
RGBD-SF	61.7	58.5	43.4	58.4	92.6	63.7	59.1	66.4	47.3	34.0	33.9	86.0	60.5	58.9	67.5
RGBD-SF-CRF	62.0	58.4	43.3	59.5	92.7	64.4	58.3	65.8	48.7	34.3	34.3	86.3	62.3	59.2	67.9
Eigen [2]	42.3	49.1	73.1	72.4	85.7	60.8	46.5	57.3	38.9	42.1	68.5	85.5	55.8	59.9	66.5
Eigen-SF	47.8	50.8	79.0	73.3	90.5	62.8	46.7	64.5	45.8	46.0	70.7	88.5	55.2	63.2	69.3
Eigen-SF-CRF	48.3	51.5	79.0	74.7	90.8	63.5	46.9	63.6	46.5	45.9	71.5	89.4	55.6	63.6	69.9
Hermans <i>et al.</i> [7]	68.4	45.4	83.4	41.9	91.5	37.1	8.6	35.8	28.5	27.7	38.4	71.8	46.1	48.0	54.3

TABLE II: NYUv2 test set results: SF denotes that the labels were produced by SemanticFusion, and the results were captured immediately if a keyframe was present. When no reconstruction is present for a pixel, we fall back to the predictions of the baseline single frame network. Note that we calculated the accuracies of [2] using their publicly available implementation. Our results are not directly comparable with Hermans *et al.* [7] as we only evaluate on a subset of the test set, and their annotations are not available. However, we include their results for reference. Following previous work [7] we exclude pixels without a corresponding depth measurement. All accuracy evaluations were performed at 320×240 resolution.

+0.3%, while the RGB-CNN only saw +0.1%. An interesting avenue for future work would be further experiments to improve the performance of both this and other kinds of map-based semantic regularisation schemes.

F. Run-time Performance

We benchmark the performance of our system on a random sample of 30 sequences from the NYUv2 test set. All tests were performed on an Intel Core i7-5820K 3.30GHz CPU and an Nvidia Titan Black GPU. Our SLAM system requires 29.3ms on average to process each frame and update the map. For every frame we also update our stored surfel probability table to account for any surfels removed by the SLAM system. This process requires an additional 1.0ms. As discussed above, the other components in our system do not need to be applied for every frame. A forward pass of our CNN requires 51.2ms and our Bayesian update scheme requires a further 41.1ms. Our standard scheme performs this every 10 frames, resulting in an average frame-rate of 25.3Hz.

Our experimental CRF implementation was developed only for the CPU in C++, but the message passing algorithm adopted could lend itself to an optimised GPU implementation. The overhead of copying data from the GPU and performing inference on a single threaded CPU implementation is significant, and we initialise the entire CRF lattice from scratch each time it is run. Therefore on average, it takes 20.3s to setup and perform 10 CRF iterations. In the evaluation above, we perform a CRF update once every 500 frames, but for online use it can be disabled entirely or applied once at the conclusion of a sequence.

V. CONCLUSIONS

Our results confirm the strong expectation that using a SLAM system to provide pixel-wise correspondences be-

分割为什么是基于2D的?

tween frames allows the fusion of per-frame 2D segmentations into a coherent 3D semantic map. It is the first time that this has been demonstrated with a real-time, loop-closure capable approach suitable for interactive room scanning. Not only that, the incorporation of such a map led to a significant increase in the corresponding 2D segmentation accuracy.

We exploited the flexibility of CNNs to improve the accuracy of a pretrained RGB network by incorporating an additional depth channel. In this work we opted for the simplest feasible solution to allow this new modality. Some recent work has explored other ways to incorporate depth information [8], but such an approach requires duplication of the lower network parameters and was infeasible in our system due to GPU memory limitations. Future research could incorporate CNN compression [9], which would not only enable the incorporation of other modalities, but also offer exciting new directions in real-time semantic segmentation on low memory and power mobile devices.

We believe that this is just the start of how knowledge from SLAM and machine-learned labelling can be brought together to enable powerful semantic and object-aware mapping. Our own reconstruction-focused dataset shows a much larger improvement in labelling accuracy via fusion than the NYU dataset with less varied trajectories, this underlines the importance of viewpoint variation. It also hints at the improvements possible with significantly longer trajectories, such as those of an autonomous robot in the field making direct use of the semantically annotated 3D map.

Going further, it is readily apparent, as demonstrated in a so far relatively simple manner in systems like SLAM++ [19] that not just should reconstruction be used to provide correspondence to help labelling, but that labelling/recognition can make reconstruction and SLAM much more accurate and efficient. A loop-closure capable surfel map as in Elas-

重建和识别是一个相互补充的过程

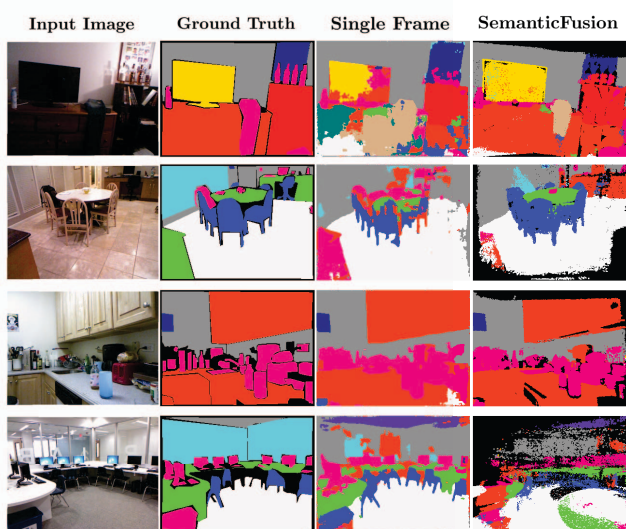


Fig. 7: **Qualitative NYUv2 test set results:** The results of SemanticFusion are using the RGBD-CNN with CRF against the same networks single frame predictions. For evaluation, the black regions of SemanticFusion denoting areas without a reconstruction, are replaced with the baseline CNN predictions. The first two rows show instances where SemanticFusion has clearly improved the accuracy of the 2D annotations. The third row shows an example of a very rotational trajectory, where there is little difference as a result of fusing predictions. The final row shows an example where the trajectory was clearly not taken with reconstruction in mind, and the distant geometry leads to tracking and mapping problems even within our subset requiring 2Hz frame-rate. Cases such as this provide an advantage to the accuracy of the single frame network.

ticFusion is highly suitable for applying operations such as class-specific smoothing (as in the extreme case of planar region recognition and fitting [18]), and this will be an interesting direction. More powerful still will be to interface with explicit object instance recognition and replace elements of the surfel map directly with 3D object models once confidence reaches a suitable threshold.

ACKNOWLEDGMENT

Research presented in this paper has been supported by Dyson Technology Ltd.

REFERENCES

- [1] C. Couprie, C. Farabet, L. Najman, and Y. LeCun, "Indoor semantic segmentation using depth information," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2013.
- [2] D. Eigen and R. Fergus, "Predicting Depth, Surface Normals and Semantic Labels with a Common Multi-Scale Convolutional Architecture," in *Proceedings of the International Conference on Computer Vision (ICCV)*, 2015.
- [3] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (VOC) challenge," *International Journal of Computer Vision (IJCV)*, no. 2, pp. 303–338, 2010.
- [4] S. Gupta, R. Girshick, P. Arbelaez, and J. Malik, "Learning Rich Features from RGB-D Images for Object Detection and Segmentation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2014.
- [5] S. Gupta, P. A. Arbeláez, R. B. Girshick, and J. Malik, "Aligning 3D models to RGB-D images of cluttered scenes," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [6] A. Handa, V. Pătrăucean, V. Badrinarayanan, S. Stent, and R. Cipolla, "SceneNet: Understanding Real World Indoor Scenes With Synthetic Data," *arXiv preprint arXiv:1511.07041*, 2015.
- [7] A. Hermans, G. Floros, and B. Leibe, "Dense 3d semantic mapping of indoor scenes from rgb-d images," in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2014.
- [8] J. Hoffman, S. Gupta, J. Leong, G. S., and T. Darrell, "Cross-Modal Adaptation for RGB-D Detection," in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2016.
- [9] F. N. Iandola, M. W. Moskewicz, K. Ashraf, S. Han, W. J. Dally, and K. Keutzer, "Squeezenet: Alexnet-level accuracy with 50x fewer parameters and <1mb model size," *CoRR*, 2016.
- [10] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," *arXiv preprint arXiv:1408.5093*, 2014.
- [11] H. S. Koppula, A. Anand, T. Joachims, and A. Saxena, "Semantic Labeling of 3D Point Clouds for Indoor Scenes," in *Neural Information Processing Systems (NIPS)*, 2011.
- [12] P. Krähenbühl and V. Koltun, "Efficient Inference in Fully Connected CRFs with Gaussian Edge Potentials," in *Neural Information Processing Systems (NIPS)*, 2011.
- [13] A. Krizhevsky, I. Sutskever, and G. Hinton, "ImageNet classification with deep convolutional neural networks," in *Neural Information Processing Systems (NIPS)*, 2012.
- [14] A. Levin, D. Lischinski, and Y. Weiss, "Colorization using Optimization," in *Proceedings of SIGGRAPH*, 2004.
- [15] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: Common objects in context," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2014, pp. 740–755.
- [16] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [17] H. Noh, S. Hong, and B. Han, "Learning deconvolution network for semantic segmentation," *arXiv preprint arXiv:1505.04366*, 2015.
- [18] R. F. Salas-Moreno, B. Glocker, P. H. J. Kelly, and A. J. Davison, "Dense Planar SLAM," in *Proceedings of the International Symposium on Mixed and Augmented Reality (ISMAR)*, 2014.
- [19] 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 Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013. [Online]. Available: <http://dx.doi.org/10.1109/CVPR.2013.178>
- [20] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus, "Indoor segmentation and support inference from RGBD images," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2012.
- [21] K. Simonyan and A. Zisserman, "Very Deep Convolutional Networks for Large-Scale Image Recognition," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2015.
- [22] S. Song, S. P. Lichtenberg, and J. Xiao, "SUN RGB-D: A RGB-D scene understanding benchmark suite," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 567–576.
- [23] J. Stückler, B. Waldvogel, H. Schulz, and S. Behnke, "Multi-resolution surfel maps for efficient dense 3d modeling and tracking," *Journal of Real-Time Image Processing JRTIP*, vol. 10, no. 4, pp. 599–609, 2015.
- [24] J. Valentin, S. Sengupta, J. Warrell, A. Shahrokni, and P. Torr, "Mesh Based Semantic Modelling for Indoor and Outdoor Scenes," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013.
- [25] V. Vineet, M. Ondrej, M. Lidegaard, M. Niener, S. Golodetz, V. A. Prisacariu, O. Kahler, D. W. Murray, S. Izadi, P. Perez, and P. H. S. Torr, "Incremental Dense Semantic Stereo Fusion for Large-Scale Semantic Scene Reconstruction," in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2015.
- [26] T. Whelan, S. Leutenegger, R. F. Salas-Moreno, B. Glocker, and A. J. Davison, "ElasticFusion: Dense SLAM without a pose graph," in *Proceedings of Robotics: Science and Systems (RSS)*, 2015.