

SceneCode: Monocular Dense Semantic Reconstruction using Learned Encoded Scene Representations

Shuaifeng Zhi, Michael Bloesch, Stefan Leutenegger, Andrew J. Davison
 Dyson Robotics Laboratory at Imperial College
 Department of Computing, Imperial College London, UK
 {s.zhi17, m.bloesch, s.leutenegger, a.davison}@imperial.ac.uk

Abstract

Systems which incrementally create 3D semantic maps from image sequences must store and update representations of both geometry and semantic entities. However, while there has been much work on the correct formulation for geometrical estimation, state-of-the-art systems usually rely on simple semantic representations which store and update independent label estimates for each surface element (depth pixels, surfels, or voxels). Spatial correlation is discarded, and fused label maps are incoherent and noisy.

We introduce a new compact and optimisable semantic representation by training a variational auto-encoder that is conditioned on a colour image. Using this learned latent space, we can tackle semantic label fusion by jointly optimising the low-dimensional codes associated with each of a set of overlapping images, producing consistent fused label maps which preserve spatial correlation. We also show how this approach can be used within a monocular keyframe based semantic mapping system where a similar code approach is used for geometry. The probabilistic formulation allows a flexible formulation where we can jointly estimate motion, geometry and semantics in a unified optimisation.

1. Introduction

Intelligent embodied devices such as robots need to build and maintain representations of their environments which permit inference of geometric and semantic properties, such as the traversability of rooms or the way to grasp objects. Crucially, if this inference is to be scalable in terms of computing resources, these representations must be *efficient*; and if devices are to operate *robustly*, the employed representations must cope with all of the variation present in the real world. However, current real-time scene understanding systems are still a long way from the performance needed for truly ground-breaking applications [5, 9].

An eventual token-like, composable scene understand-

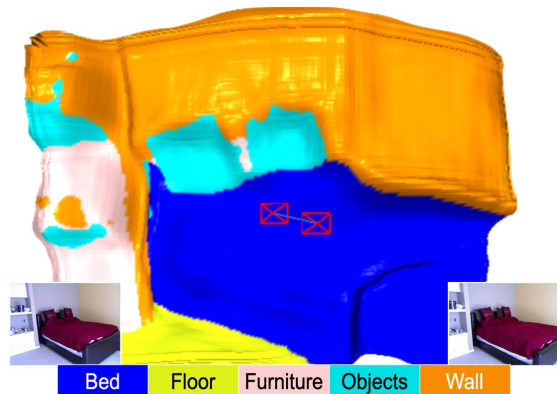


Figure 1: Dense semantic structure from motion on two frames from the NYUv2 dataset. Compact representations of semantics and geometry have been jointly optimised with camera pose to obtain smooth and consistent estimates.

ing may finally give artificial systems the capability to reason about space and shape in the intuitive manner of humans [37]. Bringing deep learning into traditional hand-designed estimation methods for SLAM has certainly led to big advances to representations which can capture both shape and semantics [6, 36], but so far these are problematic in various ways. The most straightforward approaches, such as [16, 14, 28, 38, 39, 26] which paint a dense geometric SLAM map with fused semantic labels predicted from views [25, 24, 13], are expensive in terms of representation size; label scenes in an incoherent way where each surface element independently stores its class; and do not benefit from semantic labelling improving motion or geometry estimation.

At the other end of the scale are approaches which explicitly recognise object instances and build scene models as 3D object graphs [27, 35, 30, 31]. These representations have the token-like character we are looking for, but are limited to mapping discrete ‘blob-like’ objects from known

classes and leave large fractions of scenes undescribed.

Looking for efficient representations of whole scenes, in this work we are inspired by CodeSLAM from Bloesch *et al.* [3] which used a learned encoding to represent the dense geometry of a scene with small codes which can be efficiently stored and jointly optimised in multi-view SLAM. While [3] encoded only geometry, here we show that we can extend the same conditional variational auto-encoder (CVAE) to represent the multimodal distributions of semantic segmentation. As in CodeSLAM, our learned low-dimensional semantic code is especially suitable for, but not limited to keyframe based semantic mapping systems, and allows for joint optimisation across multiple views to maximise semantic consistency. This joint optimisation alleviates the problems caused by the independence of surface elements assumed by most semantic fusion methods, and allows much higher quality multi-view labellings which preserve whole elements of natural scenes.

We show that compact representations of geometry and semantics can be jointly learned, resulting in the multitask CVAE used in this paper. This network makes it possible to build a monocular dense semantic SLAM system where geometry, poses and semantics can be jointly optimised.

To summarise, our paper has the following contributions:

- A compact and optimisable representation of semantic segmentation using an image-conditioned variational auto-encoder.
- A new multi-view semantic label fusion method optimising semantic consistency.
- A monocular dense semantic 3D reconstruction system, where geometry and semantics are tightly coupled into a joint optimisation framework.

2. Related Work

Structured semantic segmentations of the type we propose have been studied by several authors. Sohn *et al.* [33] proposed a CVAE to learn the distribution of object segmentation labels using Gaussian latent variables. Due to the learned distribution, the resulting object segmentation was more robust to noisy and partially observed data compared to discriminative CNN models. Pix2Pix from Isola *et al.* [15] used a conditional Generative Adversarial Network (GAN) to achieve image to image translation tasks in which the conditional distribution of semantic labels is implicitly learned. However, when used for semantic prediction from colour images, the GAN training process induces hallucinated objects. In addition, the distributions learned by GANs are not directly accessible and optimisable in the form we need for multi-view fusion.

Kohl *et al.* recently proposed a probabilistic U-Net [23] to address the ambiguities of semantic segmentation due

to insufficient context information. A CVAE was designed to learn the multimodal distribution of segmentations given colour images through a low-dimensional latent space, and it was shown that ambiguities can be well modelled by a compact latent code. We build on this idea and show that we can use the learned latent space to integrate multi-view semantic labels, and build a monocular dense SLAM system capable of jointly optimising geometry and semantics.

3. Compact Geometry + Semantics Encoding

Our multitask CVAE (see Figure 2) learns the conditional probability densities for depths and semantic segmentations conditioned on colour images in a manner similar to the compact representation of geometry in CodeSLAM [3]. The network consists of three main parts: a U-shaped multitask network with skip connections and two variational auto-encoders for depth and semantic segmentation.

The U-shaped multitask network contains one shared encoder with a ResNet-50 backbone [12] and two separate decoders adopting RefineNet units [24]. Unlike the original implementation, batch normalisation is added after each convolution in the RefineNet unit to stabilise training. Each of the two variational auto-encoders consists of a VGG-like fully convolutional recognition model (encoder) followed by a linear generative model (decoder), which is coupled with the U-net and thus conditioned on colour images.

More specifically, in the linear decoder the latent code is first broadcast spatially to have the same height/width and then 1×1 convolved to have the same dimensionality as the image feature maps from the last RefineNet unit. A merged tensor is then computed by doing a three-fold concatenation of the broadcast/convolved code, the RefineNet unit, and an element-wise multiplication of the two. Finally, convolution (without nonlinear activation) and bilinear up-sampling is applied to obtain the prediction. The motivation for this procedure is to obtain a linear relationship between code and prediction which is conditioned on the input image in a nonlinear manner [3] — the linearity enabling pre-computation of Jacobians during inference at test time (see Section 4). The predicted depth and semantics (unscaled logits before softmax function) can thus be formulated as:

$$D(\mathbf{c}_d, I) = D_0(I) + J_d(I) \mathbf{c}_d, \quad (1)$$

$$S(\mathbf{c}_s, I) = S_0(I) + J_s(I) \mathbf{c}_s, \quad (2)$$

where $J_{s/d}$ represents the learned linear influence, and $D_0(I) = D(0, I)$ and $S_0(I) = S(0, I)$. Due to our variational setup, $D_0(I)$ and $S_0(I)$ can be interpreted as the most likely prediction given the input image alone. Note the generality of this framework, which could be combined with arbitrary network architectures.

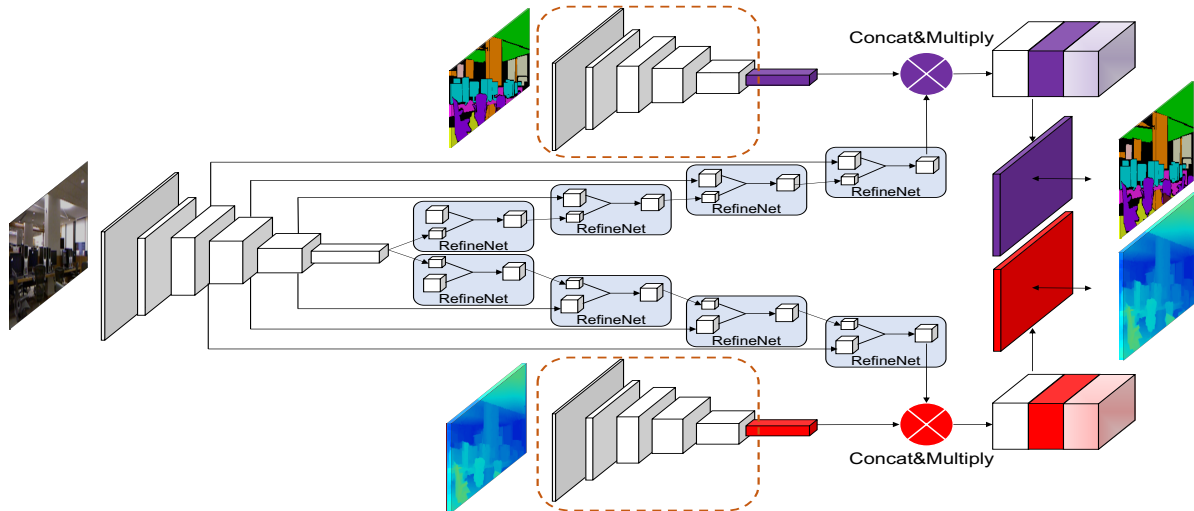


Figure 2: The proposed multitask conditional variational auto-encoder (multitask CVAE). Depth images and semantic labels (one-hot encoded) are encoded to two low-dimensional latent codes via VGG-like fully convolutional networks. These recognition models shown in the dashed boxes are not accessible during inference. The RGB images are processed by a U-shaped network with a ResNet-50 backbone. Finally, the sub-parts are combined by \otimes operations standing for a combination of broadcasting, concatenation, and element-wise multiplication.

3.1. Network Training Configuration

Both the depth and semantic predictions are jointly trained using groundtruth data. In addition to the reconstruction losses discussed in the following sections, the variational setup requires a KL-divergence based loss on the latent space [21]. In order to avoid a degrading latent space, we employ a KL annealing strategy [4, 34] where we gradually increase the weights of the KL terms from 0 after 2 training epochs. Finally, the weights of semantic vs. depth reconstruction losses are trained in an adaptive manner to account for task-dependent uncertainty [18]. In all of our experiments, we train the whole network in an end-to-end manner using the Adam optimiser [20] with an initial learning rate of 10^{-4} and a weight decay of 10^{-4} . The ResNet-50 is initialised using ImageNet pre-trained weights, and all other weights are initialised using He *et al.*'s method [11].

For depth images, as in [3], the raw depth values d are first transformed via a hybrid parametrisation called *proximity*, $p = a/(a + d)$, where a is the average depth value, which is set to 2m in all of our experiments. In this way, we can handle raw depth values ranging from 0 to $+\infty$ and assign more precision to regions closer to the camera. An L_1 loss function together with data dependent Homoscedastic uncertainty [17] is used as the reconstruction error:

$$L_{\phi, \theta}(d) = \sum_{i=1}^N \left[\frac{|\tilde{p}_i - p_i|}{b_i} + \log(b_i) \right], \quad (3)$$

where N is the number of pixels, \tilde{p}_i and p_i are the predicted

proximity and input proximity of the i -th pixel, and b_i is the predicted uncertainty of the i th pixel.

Semantic segmentation labels, which are discrete numbers, are one-hot encoded before being input to the network. Therefore, the multi-class cross-entropy function is a natural option for calculating the reconstruction loss using the predicted softmax probabilities and one-hot encoded labels:

$$L_{\phi, \theta}(s) = \frac{1}{N} \sum_{i=1}^N \sum_{c=1}^C k_c^{(i)} \log p_c^{(i)}, \quad (4)$$

where C is the number of classes, $k_c^{(i)}$ is the c -th element of the one-hot encoded labels for the i -th pixel and $p_c^{(i)}$ is the predicted softmax probability in the same position.

4. Fusion via Multi-View Code Optimisation

In a multi-view setup, depth, semantics, and motion estimates can be refined based on consistency in overlapping regions by making use of dense correspondences. While the use of photometric consistency is well-established, here we also introduce semantic consistency, i.e. any given part of our scene should have the same semantic label irrespective of viewpoint. The semantic consistency is less affected by disturbances such as non-Lambertian reflectance, but may be subject to quantisation errors and cannot be directly measured.

Given no additional information, an all-zero code is most likely code because of the multivariate Gaussian prior assumption during training (Section 3.1). This zero code

can thus be used, both as an initialisation value and as a prior during optimisation at test time (during which we have no access to depths or semantic labels). The probabilistic formulation of the system allows it to embed depth, semantics and motion into a unified probabilistic framework and thereby combine an arbitrary number of information sources including images, semantic constraints, priors, motion models or even measurements from other sensors.

4.1. Geometry Refinement

In analogy to [3], given an image I_A with its depth code c_d^A , and a second image I_B with estimated relative rigid body transformation $T_{BA} = (R_{BA}, t_{BA}) \in SO(3) \times \mathbb{R}^3$, the dense correspondence for each pixel u in view A is:

$$w(u_A, c_d^A, T_{BA}) = \pi(T_{BA} \pi^{-1}(u_A, D_A[u_A])), \quad (5)$$

where π and π^{-1} are the projection and inverse projection functions, respectively. D_A stands for $D_A = D(c_d^A, I_A)$, and the square bracket operation $[u]$ means a value look-up at pixel location u . We can then establish the photometric error r_i based on the photo-consistency assumption [19]:

$$r_i = I_A[u_A] - I_B[w(u_A, c_d^A, T_{BA})]. \quad (6)$$

Similarly, we can derive the geometric error term r_z as:

$$r_z = D_B[w(u_A, c_d^A, T_{BA})] - [T_{BA} \pi^{-1}(u_A, D_A[u_A])]_Z, \quad (7)$$

where $[\cdot]_Z$ refers to the depth value of a point.

Both photometric and geometric errors are differentiable w.r.t. the input camera poses and latent codes, so that Jacobians can be computed using the chain rule. Due to the designed linear relationship we can pre-compute the Jacobian of depth prediction w.r.t. the code which is computationally expensive due to the dense convolution operations.

4.2. Semantics Refinement

Given images I_A, I_B sharing a common field of view (FOV), and their pre-softmax predictions S_A and S_B generated from semantic codes c_s^A and c_s^B , we propose to establish a semantic error term via dense warping:

$$r_s = DS(S_A[u_A], S_B[w(u_A, c_d^A, T_{BA})]), \quad (8)$$

where DS can be an arbitrary function measuring distance/dissimilarity [7]. In the scope of this paper, DS is chosen to be the Euclidean distance after applying softmax on the logits. Establishing the semantic error on top of semantic labels is not adopted here due to the loss of information and the induced non-differentiability.

The underlying intuition of Equation 8 is that corresponding pixels must have the same semantic label, and thus similar (but not necessary the same) softmax categorical probabilities. Unlike the photo-consistency assumption,

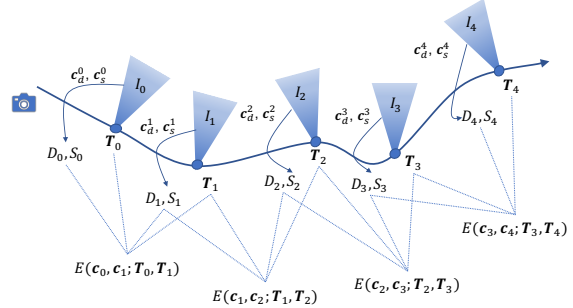


Figure 3: Semantic mapping formulation. Each keyframe has a colour image I , depth code c_d and semantic code c_s . Second order optimisation can be applied to jointly or separately optimise camera motion, geometry and semantics.

the semantic consistency assumption is comparatively weak since it is not anchored to any actual measurement, though this is somewhat alleviated by the zero-code prior described above. Nevertheless, as the viewpoint varies, different semantic cues may become available and a previously semantically ambiguous region may become more distinctive. Instead of fusing this information element-wise [28], the estimates are propagated all the way back to the semantic code, allowing spatial information fusion.

The semantic error term in Equation 8 is differentiable not only w.r.t. the semantic codes c_s^A and c_s^B , but also w.r.t. the camera pose and depth of the reference keyframe. This naturally enables semantic information to influence motion and structure estimation, i.e., the framework will for instance attempt to align chairs with chairs and walls with walls. Again, the Jacobians of the semantic logits w.r.t. the semantic code can be pre-computed.

Although our semantics refinement approach targets a monocular keyframe based SLAM system, it can be adopted as a semantic label fusion module in arbitrary SLAM system such as stereo or RGB-D SLAM systems.

5. Monocular Dense Semantic SLAM

We can integrate the geometry and semantics refinement processes into a preliminary keyframe based monocular SLAM system. The map is represented by a collection of keyframes, each with a camera pose and two latent codes, one for geometry and one for semantics, as shown in Figure 3. We follow the standard paradigm of dividing the system into tracking (front-end) and mapping (back-end) and alternate between them [22]. In the present paper, for efficiency reasons, the tracking module estimates the relative 3D motion between the current frame and the last keyframe using the photometric residual only [2].

The mapping module relies on dense N-frame structure

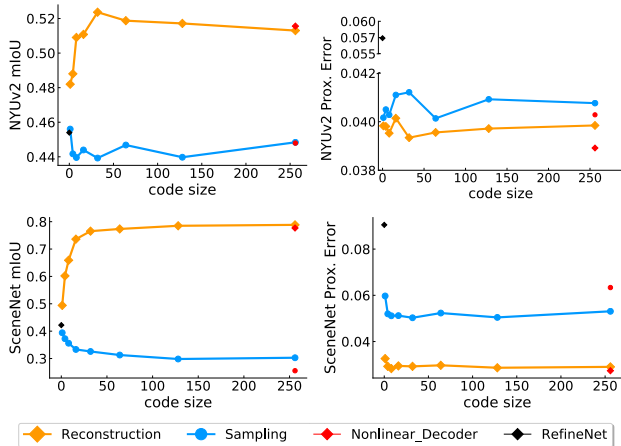


Figure 4: Reconstruction and zero-code prediction performance of different set-ups on the NYUv2 and SceneNet RGB-D test sets. Reconstruction performance increases with code size. The quality of zero-code predictions is comparable to a discriminative RefineNet model for semantic segmentation, and better on depth prediction. Using a non-linear decoder leads to little improvement.

from motion, by minimising photometric, geometric and semantic residuals with a zero-code prior between any two overlapping frames, which can be formulated as a non-linear least-square problem. As in CodeSLAM [3], we employ loss functions that (i) remove invalid correspondences, (ii) perform relative weighting for different residuals, (iii) include robust Huber weighting, (iv) down-weight strongly slanted and potentially occluded pixels. The differentiable residuals are minimised by a damped Gauss-Newton. In addition, the linear decoder allows us to pre-compute the Jacobians of the network prediction w.r.t. the code for each keyframe. Because the semantic residual relies not only on the semantic code but also on data association, during mapping we first jointly optimise the geometry and poses, then optimise the semantic residual, and lastly jointly optimise both geometry and semantics. In this way, we tightly couple geometry and semantics into a single optimisation framework.

6. Experiments

Please also see our submitted video which includes further demonstrations: <https://youtu.be/MCgbgW3WA1M>.

To test our method, we use three indoor datasets: the synthetic SceneNet RGB-D [29] dataset, and the real-world NYUv2 [32] and Stanford 2D-3D-Semantic datasets [1]. Compared to outdoor road scenes [10, 8], indoor scenes have different challenges with large variations in spatial arrangement and object sizes, and full 6-D motion.

6.1. Datasets

NYUv2 has 1,449 pre-aligned and annotated images (795 in the training set and 654 in the test set). We cropped all the available images from 640×480 to valid regions of 560×425 before further processing. The 13 class semantic segmentation task is evaluated in our experiments.

Stanford 2D-3D-Semantic is a large scale real world dataset with a different set of 13 semantic class definitions. 70,496 images with random camera parameters are split into a training set of 66,792 images (areas 1, 2, 4, 5, 6) and a test set of 3,704 images (area 3). We rectified all images to a unified camera model.

The synthetic **SceneNet RGB-D** dataset provides perfect ground truth annotations for 5M images. We use a subset: our training set consists of 110,000 images by sampling every 30th frame of each sequence from the first 11 original training splits. Our test dataset consists of 3,000 images by sampling every 100th frame from the original validation set.

All input images are resized to 256×192 . During training, we use data augmentation including random horizontal flipping and jittering of brightness and contrast. At test time, only single scale semantic prediction is evaluated.

6.2. Image Conditioned Scene Representation

We first quantitatively inspect the influence of code size on the NYUv2 and SceneNet RGB-D datasets by measuring reconstruction performance. We use the same latent code size for depth images and semantic labels for simplicity. We also train a discriminative RefineNet for semantic segmentation and depth estimation separately as a single task prediction-only baseline models (i.e. code size of 0). Figure 4 shows results for depth and semantic encoding with different code size and setups. The reconstruction performance indicates the capacity of the latent encoding for variational auto-encoders. Due to the encoded information, the reconstruction is consistently better than single view monocular prediction. Furthermore, we do not benefit from a non-linear decoder and we observe diminishing returns when the code size is larger than 32, and therefore choose this for later experiments.

The qualitative effects of our image conditioned auto-encoding of size 32 are shown in Figure 5. The zero-code predictions are usually similar to the encoded predictions, though errors in ambiguous regions are corrected given the additional encoded information. Figure 6 displays the learned image dependent Jacobians of the semantic logits w.r.t. entries in the code. We see how each code entry is responsible for certain semantically meaningful regions (e.g. examine the sofa Jacobians). Additionally, each code entry also has a tendency to decrease the probability of other ambiguous classes. For two images from different viewpoints, the image dependent Jacobians show high consistency.

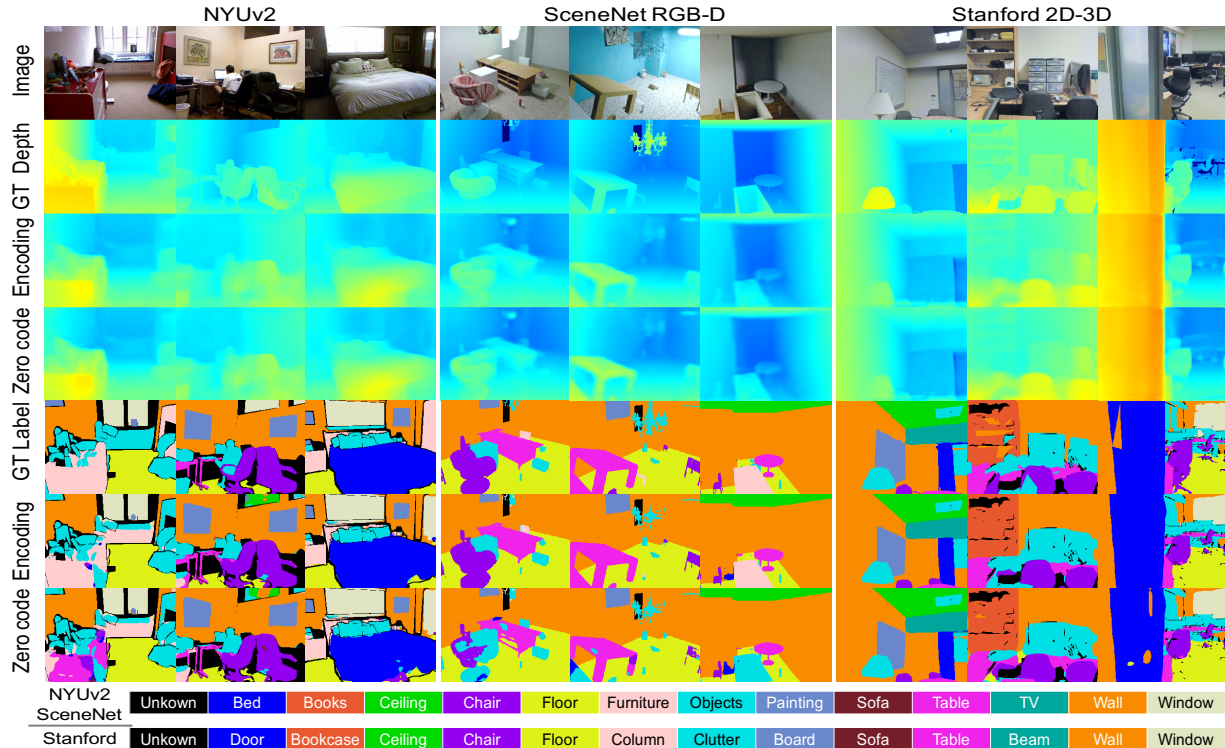


Figure 5: Qualitative results on the NYUv2 (left), SceneNetRGB-D (middle) and Stanford (right) datasets. Input colour images are at the top. We show ground truth, encoded predictions (code from encoder) and zero-code predictions (monocular predictions) for depth and semantic labels. Incorrect semantic predictions in regions which are ambiguous for monocular predictions are corrected by optimising the compact latent codes. Black regions are masked unknown classes.

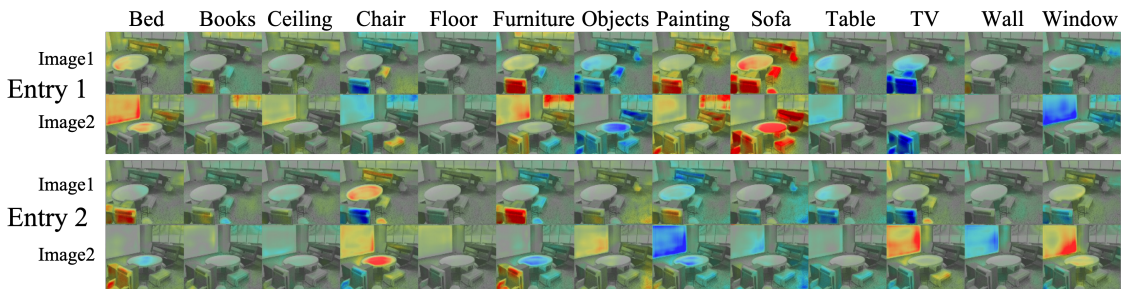


Figure 6: The Jacobians of semantic logits w.r.t. two code entries for a pair of wide baseline views. The columns represent the influence of the code entry across each semantic classes. Red and blue mean positive and negative influence, respectively. Semantically meaningful regions can be refined coherently during optimisation, leading to smooth and complete segmentation, and this property is automatically carried over into the semantic fusion process.

6.3. Semantic Label Fusion using Learned Codes

Our semantic refinement process can be regarded as a label fusion scheme for multi-view semantic mapping. An important advantage of code based fusion compared to the usual element-wise update methods for label fusion is its ability to naturally obtain spatially and temporally consis-

tent semantic labels by performing joint estimation in the latent space. This means that pixels are not assumed i.i.d when their semantic probabilities are updated, leading to smoother and more complete label regions.

To isolate only label estimation, our experiments use the SceneNet RGB-D dataset where precise ground truth depth and camera poses are available to enable perfect data as-

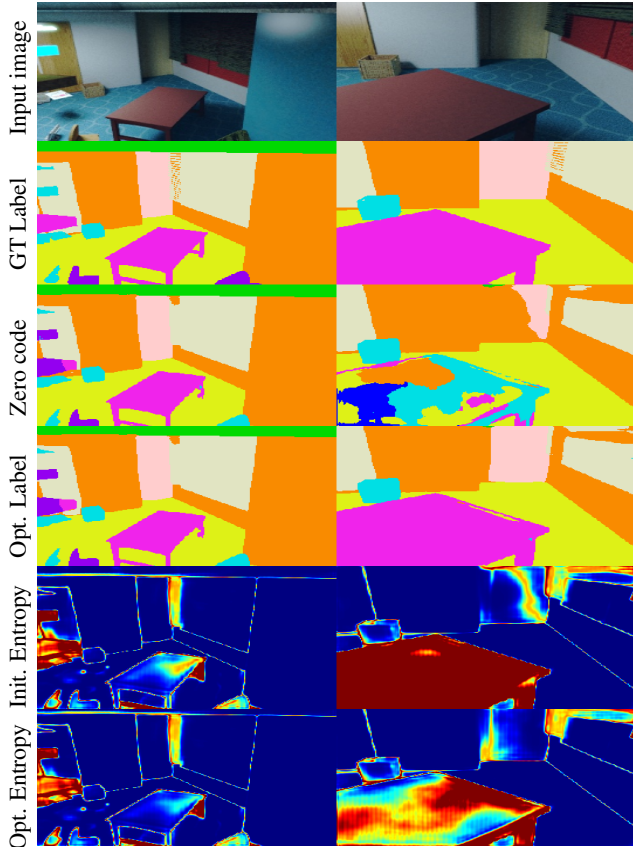


Figure 7: An example of two-view semantic label fusion of our method. From top to bottom rows: input colour image, ground truth semantic label, zero-code prediction, optimised label (minimising semantic cost), information entropy of the zero-code softmax probabilities, information entropy of the optimised softmax probabilities.

Statistics	Mean	Std	Max	Min
Rotation (degree)	5.950	9.982	163.382	0.028
Translation (meter)	0.149	0.087	0.701	0.001

Table 1: The statistics of the relative 3D motion between consecutive frames extracted from SceneNet RGB-D.

sociation. We also mask out and ignore occluded regions. We use the zero-code monocular predictions as the initial semantic predictions for all fusion methods.

In Figure 7 we show the result of semantic label fusion given two views taken with a large baseline. The RHS zero-code prediction struggles to recognise the table given the ambiguous context. The high entropy indicates that semantic labels are likely to change during optimisation. In contrast, the LHS zero-code prediction is able to accurately segment the table with relatively low entropy. By minimis-

ing the semantic cost between two views, the optimised semantic representations are able to generate a consistent predictions, successfully leading to the disambiguation of the RHS into a well segmented and smooth prediction. The entropy of both views are reduced as well. Similar improvements can also be observed in other regions. In addition, it is interesting to observe that the entropy map exhibits consistency with the scene structure, showing that the network can recognise the spatial extent of an object but struggles with the precise semantic class.

Qualitative results of different label fusion methods are shown in Figure 8. The results of both element-wise fusion approaches are obtained by integrating the probabilities of the other images into each current frame, while our result simply comes from pairing all the later frames to the first frame. For a sequence of 5 frames with small baselines, the zero-code predictions are all similar. As a result, when there is a difficult, ambiguous region (indicated by low quality zero-code predictions and high entropy), the element-wise label fusion methods lead to results which are only marginally better. However, the representation power in the learned compact code enables much smoother predictions with correct labels to be obtained through optimisation. After optimisation, the reduced entropy for these regions indicates that the network is much more confident.

Next we provide a quantitative comparison. 2000 images sampled from 1000 sequences (2 images per sequence) from the SceneNet RGB-D validation set are used to evaluate the performance. We augment every extracted image with a variable number of subsequent images in the sequence to obtain short multi-view sequences (1-4 frames). Since the trajectories of SceneNet RGB-D are randomly generated, a good variety of relative transformations and baselines are included in this set (Table 1).

Table 2 shows the effectiveness of three multi-view label fusion methods given various number of views. Our label fusion approach using code optimisation outperforms others methods. The improvement in total pixel accuracy is not significant because of the large area of walls and floors in the dataset. However, the large improvement in the mIoU metric shows that our method is able to consider more on high-order statistics, indicating smoother predictions and better results on other small objects.

Effect of Code Prior during Semantic Optimisation

During semantic optimisation we use a zero-code regularisation term. Without this term, the optimisation may be drawn to locally consistent but incorrect semantic labels. Table 2 shows that the accuracy of two-view label fusion without a zero-code prior is even lower than single view prediction, underlining the importance of this prior.

#Views	Method	Pix. Acc.	Cls. Acc.	mIoU
1	-	75.167	63.330	41.713
2	Multiplication	75.424	63.629	42.326
	Average	75.374	63.549	42.220
	Ours	75.725	63.750	43.842
	Ours (w/o prior)	74.498	60.646	39.600
3	Multiplication	75.542	63.815	42.692
	Average	75.451	63.754	42.213
	Ours	75.815	63.827	44.231
4	Multiplication	75.578	63.950	42.795
	Average	75.358	63.767	42.102
	Ours	75.668	63.720	44.263

Table 2: The effectiveness of different label fusion methods on 2000 images sampled from SceneNet RGB-D. The large improvement on the metric of intersection over union shows that our label fusion lead to smoother predictions.

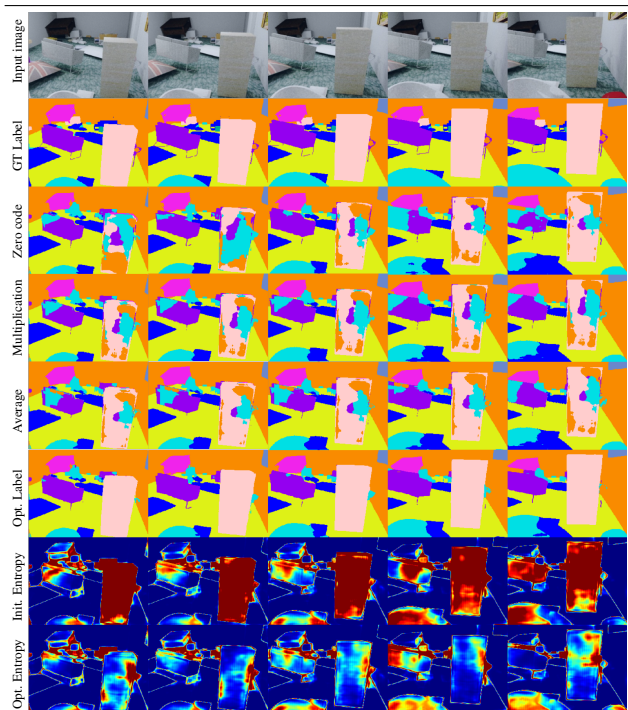


Figure 8: Qualitative comparison of different label fusion methods. 5 consecutive frames with a small baseline are chosen. Our method can effectively fuse multi-view semantic labels to generate smoother semantic predictions.

6.4. Monocular Dense Semantic SLAM

We present example results from our preliminary full monocular dense semantic SLAM system. Due to the prior information of geometry encoded in the system, the system is very robust during initialisation and can manage pure rotational motion. The system currently runs in a sliding window manner. Figures 1 and 9 show examples of two-view dense semantic structure from motion from different

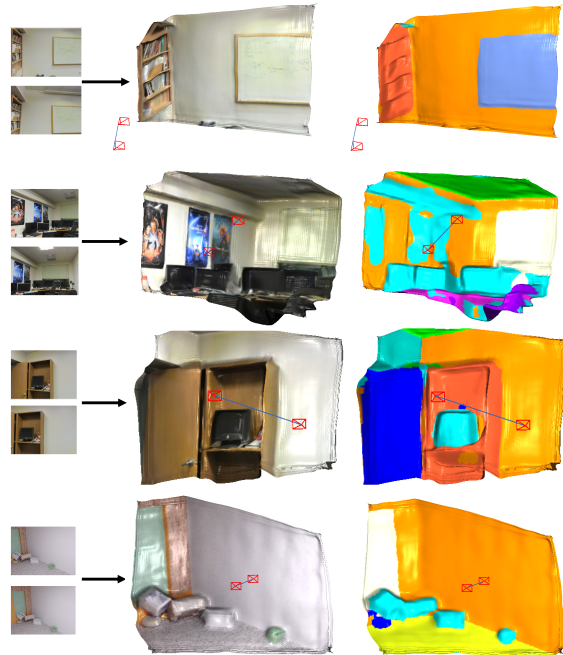


Figure 9: Qualitative results of two-view structure from motion on two selected frames from Stanford dataset (first 3 rows) and SceneNet RGB-D dataset (last row). The compact representations of both semantics and geometry are (jointly) optimised with came pose to obtain a dense map with consistent semantic labels and relative camera motion.

datasets.

7. Conclusion and Future Work

We have shown that an image-conditioned learned compact representation can coherently and efficiently represent semantic labels. This code can be optimised across multiple overlapping views to implement semantic fusion with many advantages over the usual methods which operate in a per-surface-element independent manner. As well as proving this fusion capability experimentally, we have built and demonstrated a prototype full dense, semantic monocular SLAM system based on learned codes where geometry, poses and semantics can all be jointly optimised.

In future work, we aim to unify learned geometric and semantic representations still further as we continue to reach towards scene models with optimal representational efficiency for truly useful real-time Spatial AI systems.

8. Acknowledgements

Research presented in this paper has been supported by Dyson Technology Ltd. Shuaifeng Zhi holds a China Scholarship Council-Imperial Scholarship. We are very grateful to Jan Czarowski for research and software collaboration on this project.

References

- [1] Iro Armeni, Alexander Sax, Amir R. Zamir, and Silvio Savarese. Joint 2D-3D-Semantic Data for Indoor Scene Understanding. *arXiv preprint arXiv:1702.01105*, 2017. 5
- [2] Simon Baker and Iain Matthews. Lucas-Kanade 20 years on: A Unifying Framework: Part 1. *International Journal of Computer Vision (IJCV)*, 56(3):221–255, 2004. 4
- [3] Michael Bloesch, Jan Czarnowski, Ronald Clark, Stefan Leutenegger, and Andrew J. Davison. CodeSLAM — Learning a Compact, Optimisable Representation for Dense Visual SLAM. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2, 3, 4, 5
- [4] Samuel R. Bowman, Luke Vilnis, Oriol Vinyals, Andrew M Dai, Rafal Jozefowicz, and Samy Bengio. Generating Sentences from a Continuous Space. *arXiv preprint arXiv:1511.06349*, 2015. 3
- [5] Cesar Cadena, Luca Carlone, Henry Carrillo, Yasir Latif, Davide Scaramuzza, José Neira, Ian Reid, and John J. Leonard. Past, Present, and Future of Simultaneous Localization and Mapping: Toward the Robust-Perception Age. *IEEE Transactions on Robotics (T-RO)*, 32(6):1309–1332, Dec 2016. 1
- [6] Cesar Cadena, Anthony R. Dick, and Ian Reid. Multi-modal Auto-Encoders as Joint Estimators for Robotics Scene Understanding. In *Proceedings of Robotics: Science and Systems (RSS)*, 2016. 1
- [7] Sung-Hyuk Cha and Sargur N. Srihari. On Measuring the Distance between Histograms. *Pattern Recognition*, 35(6):1355–1370, 2002. 4
- [8] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The Cityscapes Dataset for Semantic Urban Scene Understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 5
- [9] Andrew J. Davison. Futuremapping: The Computational Structure of Spatial AI Systems. *arXiv preprint arXiv:1803.11288*, 2018. 1
- [10] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for Autonomous Driving? The KITTI Vision Benchmark Suite. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012. 5
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2015. 3
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2
- [13] Yang He, Wei-Chen Chiu, Margret Keuper, and Mario Fritz. STD2P: RGBD Semantic Segmentation Using Spatio-Temporal Data-Driven Pooling. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 1
- [14] Alexander Hermans, Georgios Floros, and Bastian 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. 1
- [15] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2
- [16] Olaf Kahler and Ian Reid. Efficient 3D Scene Labelling Using Fields of Trees. In *Proceedings of the IEEE International Conference on Computer Vision*, 2013. 1
- [17] Alex Kendall and Yarin Gal. What Uncertainties Do We Need in Bayesian Deep Learning for Computer Vision? In *Neural Information Processing Systems (NIPS)*, 2017. 3
- [18] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-Task Learning Using Uncertainty to Weigh Losses for Scene Geometry and Semantics. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 3
- [19] Christian Kerl, Jürgen Sturm, and Daniel Cremers. Robust Odometry Estimation for RGB-D Cameras. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2013. 4
- [20] Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2015. 3
- [21] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2014. 3
- [22] Georg Klein and David W. Murray. Parallel Tracking and Mapping for Small AR Workspaces. In *Proceedings of the International Symposium on Mixed and Augmented Reality (ISMAR)*, 2007. 4
- [23] Simon A. A. Kohl, Bernardino Romera-Paredes, Clemens Meyer, Jeffrey De Fauw, Joseph R. Ledsam, Klaus H. Maier-Hein, S. M. Eslami, Danilo Jimenez Rezende, and Olaf Ronneberger. A Probabilistic U-Net for Segmentation of Ambiguous Images. In *Neural Information Processing Systems (NIPS)*, 2018. 2
- [24] Guosheng Lin, Anton Milan, Chunhua Shen, and Ian Reid. RefineNet: Multi-Path Refinement Networks for High-Resolution Semantic Segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 1, 2
- [25] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully Convolutional Networks for Semantic Segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015. 1
- [26] Lingni Ma, Jörg Stückler, Christian Kerl, and Daniel Cremers. Multi-View Deep Learning for Consistent Semantic Mapping with RGB-D Cameras. In *Proceedings of the IEEE/RSJ Conference on Intelligent Robots and Systems (IROS)*, 2017. 1
- [27] John McCormac, Ronald Clark, Michael Bloesch, Andrew J. Davison, and Stefan Leutenegger. Fusion++: Volumetric Object-Level SLAM. In *Proceedings of the International Conference on 3D Vision (3DV)*, 2018. 1

- [28] John McCormac, Ankur Handa, Andrew J. Davison, and Stefan Leutenegger. SemanticFusion: Dense 3D Semantic Mapping with Convolutional Neural Networks. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2017. 1, 4
- [29] John McCormac, Ankur Handa, Stefan Leutenegger, and Andrew J. Davison. SceneNet RGB-D: Can 5M Synthetic Images Beat Generic ImageNet Pre-training on Indoor Segmentation? In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2017. 5
- [30] Lachlan Nicholson, Michael Milford, and Niko Sünderhauf. QuadricSLAM: Constrained Dual Quadrics from Object Detections as Landmarks in Object-oriented SLAM. *IEEE Robotics and Automation Letters*, 2018. 1
- [31] Martin Rünz, Maud Buffier, and Lourdes Agapito. MaskFusion: Real-time Recognition, Tracking and Reconstruction of Multiple Moving Objects. In *Proceedings of the International Symposium on Mixed and Augmented Reality (ISMAR)*, 2018. 1
- [32] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor Segmentation and Support Inference from RGBD Images. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2012. 5
- [33] Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning Structured Output Representation using Deep Conditional Generative Models. In *Neural Information Processing Systems (NIPS)*, 2015. 2
- [34] Casper Kaae Sønderby, Tapani Raiko, Lars Maaløe, Søren Kaae Sønderby, and Ole Winther. How to Train Deep Variational Autoencoders and Probabilistic Ladder Networks. *arXiv preprint arXiv:1602.02282*, 2016. 3
- [35] Niko Sünderhauf, Trung T. Pham, Yasir Latif, Michael Milford, and Ian Reid. Meaningful Maps With Object-Oriented Semantic Mapping. In *Proceedings of the IEEE/RSJ Conference on Intelligent Robots and Systems (IROS)*, 2017. 1
- [36] Chamara Saroj Weerasekera, Yasir Latif, Ravi Garg, and Ian Reid. Dense Monocular Reconstruction using Surface Normals. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2017. 1
- [37] Jiajun Wu, Ilker Yildirim, Joseph J. Lim, William T. Freeman, and Joshua B. Tenenbaum. Galileo: Perceiving Physical Object Properties by Integrating a Physics Engine with Deep Learning. In *Neural Information Processing Systems (NIPS)*, 2015. 1
- [38] Yu Xiang and Dieter Fox. DA-RNN: Semantic Mapping with Data Associated Recurrent Neural Networks. *arXiv preprint arXiv:1703.03098*, 2017. 1
- [39] Jianxiang Xiao and Long Quan. Multiple View Semantic Segmentation for Street View Images. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2009. 1