# iLabel: Interactive Neural Scene Labelling

Shuaifeng Zhi[1*]    Edgar Sucar[1*]    Andre Mouton[2]    Iain Haughton[2]
Tristan Laidlow[1]    Andrew J. Davison[1]
[1] Dyson Robotics Lab, Imperial College
[2] Dyson Ltd.

{s.zhi17,e.sucar18}@imperial.ac.uk

## Abstract

*Joint representation of geometry, colour and semantics using a 3D neural field enables accurate dense labelling from ultra-sparse interactions as a user reconstructs a scene in real-time using a handheld RGB-D sensor. Our iLabel system requires no training data, yet can densely label scenes more accurately than standard methods trained on large, expensively labelled image datasets. Furthermore, it works in an 'open set' manner, with semantic classes defined on the fly by the user.*

*iLabel's underlying model is a multilayer perceptron (MLP) trained from scratch in real-time to learn a joint neural scene representation. The scene model is updated and visualised in real-time, allowing the user to focus interactions to achieve efficient labelling. A room or similar scene can be accurately labelled into 10+ semantic categories with only a few tens of clicks. Quantitative labelling accuracy scales powerfully with the number of clicks, and rapidly surpasses standard pre-trained semantic segmentation methods. We also demonstrate a hierarchical labelling variant.*

## 1. Introduction

An intelligent agent must build an internal representation of its environment which goes beyond geometry and colour to include a semantic understanding of the scene. Research on neural field representations has shown that an MLP network can be trained from scratch in a single scene via automatic self-supervision to accurately and flexibly represent geometry and appearance [22,38]. In this paper we demonstrate that the internal scene structure learned by the network allows for efficient user-guided scene segmentation.

We introduce iLabel, the first online and interactive 3D scene capturing system with a unified neural field representation, which allows a user to achieve high-quality, dense scene reconstruction and multi-class semantic la-
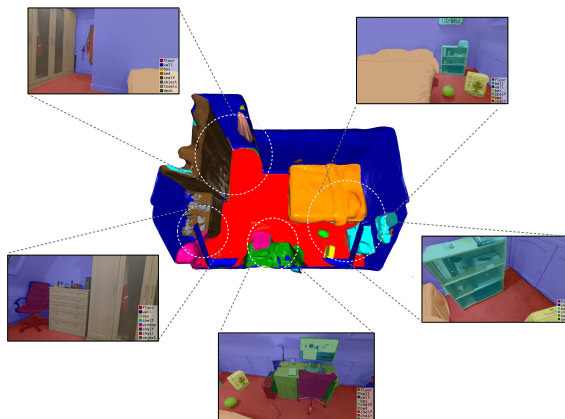


Figure 1. Whole-room semantic mesh labelled in real-time from only 140 interactive clicks and no prior training data. See https://youtu.be/bL7RZaMhRbk for a video demonstration.

belling from scratch with only minutes of scanning and a few tens of semantic click annotations. A real-time neural field SLAM forms the basis of our system. The user simultaneously scans a scene and provides sparse semantic annotations on selected keyframe images. By supervising the network on the sparse annotations, semantics are automatically propagated. The ability to render full predictions in real-time allows a user-in-the-loop to place annotations efficiently, fixing incorrect predictions or adding new classes.

Our approach requires no prior training on semantic datasets, and can therefore be applied in novel contexts, with categories defined on-the-fly by the user in an open-set manner. Standard methods for semantic scene segmentation use deep networks trained on datasets of thousands of images with dense, high-quality human annotations; even then they often have poor performance when the test scene is not a good match for the training set. We show that the quantitative labelling accuracy of iLabel scales powerfully with the number of clicks, and rapidly surpasses the accuracy of standard pre-trained semantic segmentation methods.

---

*Authors contributed equally to this work.

1

Alongside our core iLabel system for multi-class scene labelling via clicks, we present two variations. First, we show that hierarchical semantic labelling can be achieved by interpreting outputs as branches in a binary tree. Second, we demonstrate a 'hands free' labelling mode where an automatic uncertainty-guided framework selects a sequence of pixels for which to ask the user for label names without the need for clicks.

The only comparable interactive scene understanding system is SemanticPaint [41], which trains a classifier on top of a separate dense SLAM system. It requires alternating between training and prediction modes, making labelling cumbersome. We argue that the unified scene representation in iLabel is simpler and more user friendly, and also show qualitatively that iLabel obtains much more precise and complete segmentations.

We demonstrate iLabel in a wide variety of environments, from tabletop scenes to entire rooms and even outdoors. We believe iLabel to be a powerful and user-friendly tool, with much potential for interactive scene labelling for augmented reality or robotics, as well as providing intuitive insights into the ability of neural fields to jointly represent correlated quantities.

## 2. Related Work

### 2.1. iLabel System Overview

**Scene Representation for Visual SLAM**  Scene representation in visual SLAM has gradually progressed from sparse feature point sets [9, 10, 23] to dense geometric 3D maps (e.g. surfels, meshes and voxels) [8, 26, 28, 42] and more recently, to neural representations [3], increasingly involving semantics [20, 25, 32, 38, 40, 44, 45]. While classical dense scene representations (e.g. volumetric maps) have several advantages, a trade-off exists between computational cost and topological complexity. Several papers [3, 6, 39, 44] have shown that view-based code representations are able to learn rich prior information from off-line training, enabling joint optimisation of geometry, poses and semantics, to refine network predictions during inference. 3D neural field scene representations have recently gained popularity, owing to their ability to represent complex scene structures with a small memory footprint by exploiting 3D awareness and spatial continuity [22, 29, 34, 38]. More recently, iMAP [38] has been proposed as a real-time SLAM system built upon an efficient neural field representation and has demonstrated the ability to reconstruct high-quality, water-tight 3D meshes.

**Online Scene Understanding and Labelling**  Existing real-time, dense semantic mapping systems typically contain two parallel modules: 1) an RGB-D based geometric SLAM system, maintaining a dense 3D map of the scene, and 2) a semantic segmentation module that predicts dense

semantic labels of the scene [12, 19, 24, 37]. Multi-view semantic predictions are incrementally fused into the geometric model, yielding densely-labelled, coherent 3D scenes. While semantic segmentation has been performed using a variety of techniques [4, 14, 16, 27], it is an inherently user-dependent and subjective problem [18]. User-in-the-loop systems are therefore crucial in enabling full flexibility when defining semantic relations between entities in a scene. In this context, the works most closely related to ours are SemanticPaint [41] and Semantic Paintbrush [21].

SemanticPaint [41] is an online, user-in-the-loop system that allows the user to label a scene during capture. To this end, the user interacts with a 3D volumetric map, built from an RGB-D SLAM system, via voice and hand gestures [28]. A streaming random forest classifier, using hand-crafted features, learns continuously from the user gestures in 3D space. The forest predictions are used as unary terms in a conditional random field (CRF) to propagate the user annotations to unseen regions. As the CRFs are built upon the reconstructed data, there is an underlying assumption that this data is good enough to support label propagation. SemanticPaint is therefore restricted to comparably simple scenes and its efficacy in complex real-word scenarios is limited. A significant distinguishing factor between iLabel and SemanticPaint is ease-of-use. SemanticPaint has several distinct modes, requiring the user to switch between modes repeatedly and at well-time intervals to obtain optimal results. In contrast, iLabel offers a much simpler and intuitive user experience, such that high-quality segmentations are obtained with far fewer interactions and no expert knowledge/intuition.

Semantic Paintbrush [21] extends SemanticPaint to operate in outdoor scenes. Using a purely passive stereo setup for extended range and outdoor depth estimation, users visualise the reconstruction through a pair of optical see-through glasses and can draw directly onto it using a laser pointer to annotate the objects in the scene. The system learns in an online manner from the these annotations and is thus able to segment other regions in the 3D map.

In contrast to [21, 41], iLabel does not rely on hand-crafted features, benefiting instead from a powerful joint internal representation of shape and appearance.

**Hierarchical Semantic Segmentation**  Finding the hierarchical structure of complex scenes is a long-standing problem. Early attempts [1, 2] used image statistics to extract an ultrametric contour map (UCM), leading to further work on using convolutional neural networks (CNNs) for hierarchical image segmentation in a supervised manner [13, 17, 43]. We show that iLabel can build a user-defined hierarchical scene segmentation interactively and store it within the weights of an MLP.
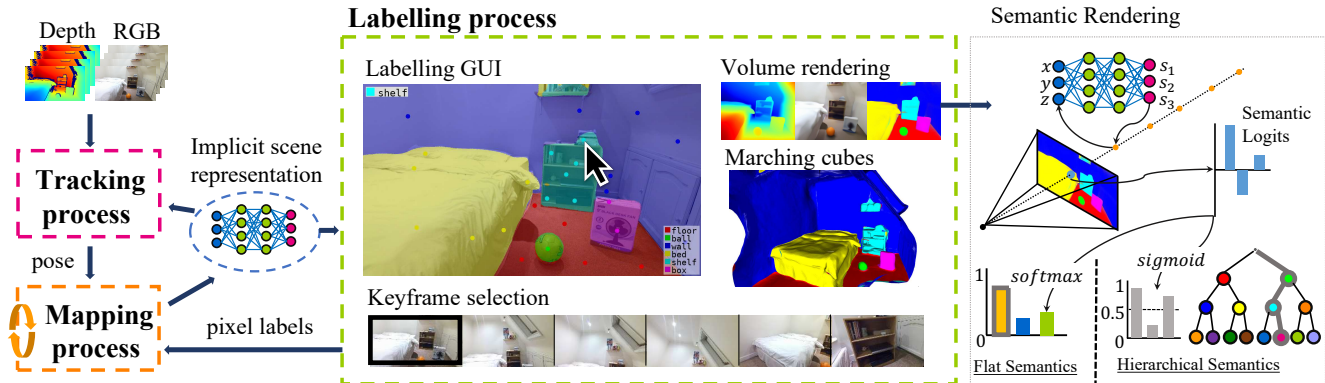
Figure 2. Overview of the iLabel system pipeline.

# 3. iLabel: Online, Interactive Open-Set Labelling and Learning

The core real-time SLAM elements of iLabel are similar to iMAP [38], which represents 3D scenes using a neural field MLP which maps a 3D coordinate to a colour and volume density. It jointly optimises the MLP and the poses of keyframes through differential volume rendering with actively sampled sparse pixels, while tracking the position of a moving RGB-D camera against the neural representation.

iLabel adds a semantic head to the MLP that predicts either a flat class distribution or a binary hierarchical tree (see Section 3.1). While SLAM continues, a user provides annotations via clicks in the keyframes. Scene semantics are then optimised through semantic rendering of these user-selected pixels. The smoothness and compactness priors present in the MLP mean that the user-supplied labels are automatically and densely propagated throughout the scene. iLabel is thus able to produce accurate, dense predictions from very sparse annotations and to often even auto-segment objects and other regions not labelled by the user. The ability to simultaneously reconstruct and label a scene in real-time allows for ultra-efficient labelling of new regions and for easy correction of errors in the current semantic predictions. Figure 2 gives an overview of the iLabel system.

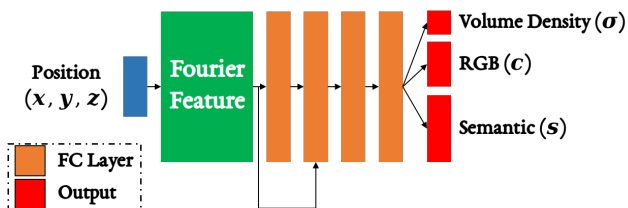## 3.1. Semantics Representation and Optimisation



Figure 3. We employ a 4-layer MLP with feature size of 256.

At the heart of iLabel is continuous optimisation of the underlying implicit scene representation, which follows the network design of iMAP with an additional semantic head (Figure 3):

$$F_\theta(\mathbf{p}) = (\mathbf{c}, \mathbf{s}, \rho), \qquad (1)$$

where $F_\theta$ is an MLP parameterised by $\theta$; $\mathbf{c}$, $\mathbf{s}$ and $\rho$ are the radiance, semantic logits and volume density at the 3D position $\mathbf{p} = (x, y, z)$, respectively. The scene representation is optimised with respect to volumetric renderings of depth, colour and semantics, computed by compositing the queried network values along the back-projected ray of pixel $[u, v]$:

$$\hat{D}[u,v] = \sum_{i=1}^{N} w_i d_i, \quad \hat{I}[u,v] = \sum_{i=1}^{N} w_i \mathbf{c}_i, \quad \hat{S}[u,v] = \sum_{i=1}^{N} w_i \mathbf{s}_i, \qquad (2)$$

where $w_i = o_i \prod_{j=1}^{i-1}(1 - o_j)$ is the ray-termination probability of sample $i$ at depth $d_i$ along the ray; $o_i = 1 - \exp(-\rho_i \delta_i)$ is the occupancy activation function; $\delta_i = d_{i+1} - d_i$ is inter-sample distance.

As in [38], geometry and keyframe camera poses $\{T_{WC}\}$ are optimised by minimising the discrepancy between the captured and rendered RGB-D images from sparsely sampled pixels. Semantics are optimised with respect to the user-labelled pixels, with two different activations and losses, corresponding to the two semantic modes described below. The right side of Figure 2 gives an overview of the semantic rendering process and the activation functions applied to the rendered logits.

**Flat Semantics** As in [45], the network outputs $\mathbf{s}_i$ are multi-class semantic logits which are converted into image space by differential volume rendering (Equation 2) followed by a *softmax* activation $\hat{\mathbf{S}}[u,v] = softmax(\hat{S}[u,v])$. Semantics are then optimised using the image cross-entropy loss between the provided class ID and the rendered predictions.

3

**Hierarchical Semantics**  We propose a novel hierarchical semantic representation through a binary tree, which allows for labelling and predicting semantics at different hierarchical levels. While the network output, $\mathbf{s}_i$, is still represented by an $n$-dimensional flat vector, $n$ now corresponds to the depth of the binary tree as opposed to the number of semantic classes. The semantic logits are rendered in the same manner, but the image activation and loss functions differ.

A *sigmoid* activation function is applied to the rendered logits, producing values in the range $[0, 1]$. The $j^{\text{th}}$ rendered output value, $\hat{\mathbf{S}}_j[u,v] = sigmoid(\hat{S}_j[u,v])$, corresponds to the branching factor at tree level $j$. To obtain a hierarchical semantic prediction, each value $\hat{\mathbf{S}}_j[u,v]$ is set to 0 or 1 by thresholding $\hat{\mathbf{S}}_j[u,v]$ at 0.5. In the hierarchical setting, the user-supplied label corresponds to selecting a specific node in the binary tree. This label is transformed into a binary branching representation, and a binary cross-entropy loss is computed for each rendered value. A label selecting a tree node at level $L$ only conditions the loss on the output values up to and including level $L$: $\hat{\mathbf{S}}_j[u,v], j \in \{1, ..., L\}$.

With reference to the top half of Figure 8, the network outputs three values corresponding to the three levels in the tree. First, the user separates the scene into *foreground* and *background* classes. A background label corresponds to the vector $[0, *, *]$ where $*$ indicates that no loss is calculated for the second and third rendered values. The user then divides the background class further into *wall* and *floor*, where the *wall* label corresponds to vector $[0, 1, *]$. The binary hierarchical representation allows the user to separate objects in stages. For example the user first separates a whole bookshelf from the rest of the scene, and later separates the books from the shelf without contradicting the initial labels, meaning that no labelling effort is wasted.

## 3.2. Semantic User Interaction Modes

Our system allows for two modes of interaction: 1) **manual interaction**, the usual interactive mode of iLabel, where users provide semantic labels in image space via clicks, and 2) **automatic query generation**, where the system generates automatic queries for the labels of informative pixels, driven by semantic prediction uncertainty (Figure 4). The latter mode eases the burden of manual annotation, and users could provide labels via text or voice.

**Automatic Query Generation**  Uncertainty-based sampling is used in this work to actively propose pixel positions for label acquisition because it can integrate seamlessly with deep neural networks with little computational overhead [30, 33]. Several uncertainty measures are explored: softmax entropy, least confidence and margin sampling [33]. For example, the softmax entropy is defined as:

$$u_{entropy} = -\sum_{c=1}^{C} \hat{\mathbf{S}}^c[u,v]\log(\hat{\mathbf{S}}^c[u,v]), \qquad (3)$$

where $C$ is the number of semantic categories.

At system run-time, semantic labels and corresponding uncertainty maps of all registered keyframes are rendered. To decide which keyframe to allocate queries, we first compute frame-level entropy by accumulating pixel-wise entropy within frames and assign a higher probability to sampling the keyframe with higher frame-level entropy. Given a selected keyframe, we then randomly select the queried pixel coordinate from a pool of pixel positions with top-K highest entropy values. The frame-level and pixel-level uncertainty are updated every certainty mapping steps. K is set to 1% or 5% of pixel numbers to avoid repeated queries at nearby positions.



Figure 4. In hands-free mode with automatic query generation, semantic class uncertainty is used to actively select a pixel for which to request a label; in this case an unlabelled stool with ambiguous class prediction and high uncertainty is selected.

## 3.3. Implementation Details

iLabel operates in a multiprocessing, single or multi-GPU framework, running three concurrent processes: 1) tracking, 2) mapping, and 3) labelling (see Figure 2).

The mapping process encompasses optimising the MLP parameters with respect to a growing set of $W$ keyframes and associated RGB-D observations: $\{(I_i, D_i, T_i)\}_{i=1}^{W}$. As per [38], the photometric loss $L_p$ and geometric loss $L_g$ are minimised on sparse, information-guided pixels. iLabel performs an additional optimisation on $K$ user-selected pixels ($\xi_i$) in each keyframe and introduces a semantic loss $L_s$, minimising the following objective function:

$$\arg\min_{\theta} \frac{1}{K} \sum_{i=1}^{W} \sum_{(u,v)\in\xi_i} \underbrace{e_i^g[u,v]}_{L_g} + \alpha_p \underbrace{e_i^p[u,v]}_{L_p} + \alpha_s \underbrace{e_i^s[u,v]}_{L_s},$$
$$(4)$$

Figure 5. Segmentation results for challenging skeletal objects; left: pre-trained CNN on ScanNet (see Section 4.2), right: iLabel.



Figure 6. Catalog of object mesh assets separated with iLabel.



Figure 7. Precise segmentations can be obtained from just 1 or 2 interactive clicks per object. (Left: clicks; middle: dense labels rendered into a keyframe; right: full 3D mesh with labels.)

where:

$$e_i^p[u,v] = \left| I_i[u,v] - \hat{I}_i[u,v] \right|, e_i^s[u,v] = -\sum_{c=1}^{C} \mathbf{S}_i^c[u,v] \log(\hat{\mathbf{S}}_i^c[u,v]),$$

$$e_i^g[u,v] = \frac{\left| D_i[u,v] - \hat{D}_i[u,v] \right|}{\sqrt{\hat{D}_{var}[u,v]}}, \hat{D}_{var}[u,v] = \sum_{i=1}^{N} w_i(\hat{D}[u,v] - d_i)^2,$$

and in the hierarchical setting:

$$e_i^s[u,v] = \sum_{l=1}^{L} -\mathbf{S}_i^c[u,v] \log(\hat{\mathbf{S}}_i^c[u,v]) - (1 - \mathbf{S}_i^c[u,v]) \log(1 - \hat{\mathbf{S}}_i^c[u,v]).$$

The labelling process coordinates user interactions (clicks and labels) and controls the rendering of semantic images and meshes (via marching cubes on a dense voxel grid queried from the MLP). The ADAM optimiser is used with poses and map learning rates of 0.003 and 0.001. $\alpha_c$ and $\alpha_s$ are 5 and 8.

iLabel does not have an explicit/specific refinement process, and all user clicks are involved in the joint optimisation (Equation 4). The optimisation keeps working and growing with changing sparse samples for colour and geometry reconstruction, and increasing annotated pixels for semantics, colour and depth as well.

# 4. Experiments and Applications

iLabel is an interactive tool intended for real-time use and we therefore emphasise that its strengths are best illustrated *qualitatively*. We provide extensive examples to demonstrate iLabel in a variety of interesting scenes, and highly recommend that reviewers watch our **attached video** (Figure 1) which shows the full interactive labelling process. We show qualitative comparisons with the only comparable system SemanticPaint and clearly demonstrate better segmentation quality. Additionally we perform a quantitative evaluation to show how segmentation quality scales with additional click labels, using a state-of-the-art, fully-supervised RGB-D segmentation baseline [5].

## 4.1. Qualitative Evaluation

As the geometry, colour and semantic heads share a single MLP backbone, user annotations are naturally propagated to untouched regions of the scene without specifying an explicit propagation mechanism (e.g. the pairwise terms of a CRF used in [41]). This, together with a user-in-the-loop, enables ultra-efficient scene labelling with only a small number of well-placed clicks.

We have observed that the resulting embeddings are highly correlated for coherent 3D entities in the scene (e.g. objects, surfaces, etc.). Consequently, iLabel is able to segment these entities very efficiently, even with a single click. This is illustrated in Figures 7 and 9, where only a few clicks generate complete and precise segmentations for a wide range of objects and entities, ranging from small, coherent objects (e.g. fruit) to deformable and intricate entities (clothing and furniture). In Figure 10 we disable colour optimisation to further highlight that in iLabel geometry provides a strong signal for separating objects.

The coordinate-based representation avoids quantisation and allows the network to be queried at arbitrary resolutions. This property allows reconstruction of detailed geometry and skeletal shapes that, when semantically labelled, render very precise segmentations. Figure 5 illustrates high-fidelity segmentations of objects which are challenging for a standard CNN.

iLabel can be used as an efficient tool for generating labelled scene datasets. For example, a scene of a complete room with 13 classes, can be fully segmented with high precision with only 140 user clicks (Figure 1). Alternatively, iLabel can be used to tag individual objects for generating object-asset catalogues (Figure 6) to aid robotic manipulation tasks, for example.

While iLabel is particularly powerful at segmenting coherent entities, Figure 11 also demonstrates its ability to propagate user-supplied labels to disjoint objects exhibiting similar properties. Each example shows label transfer
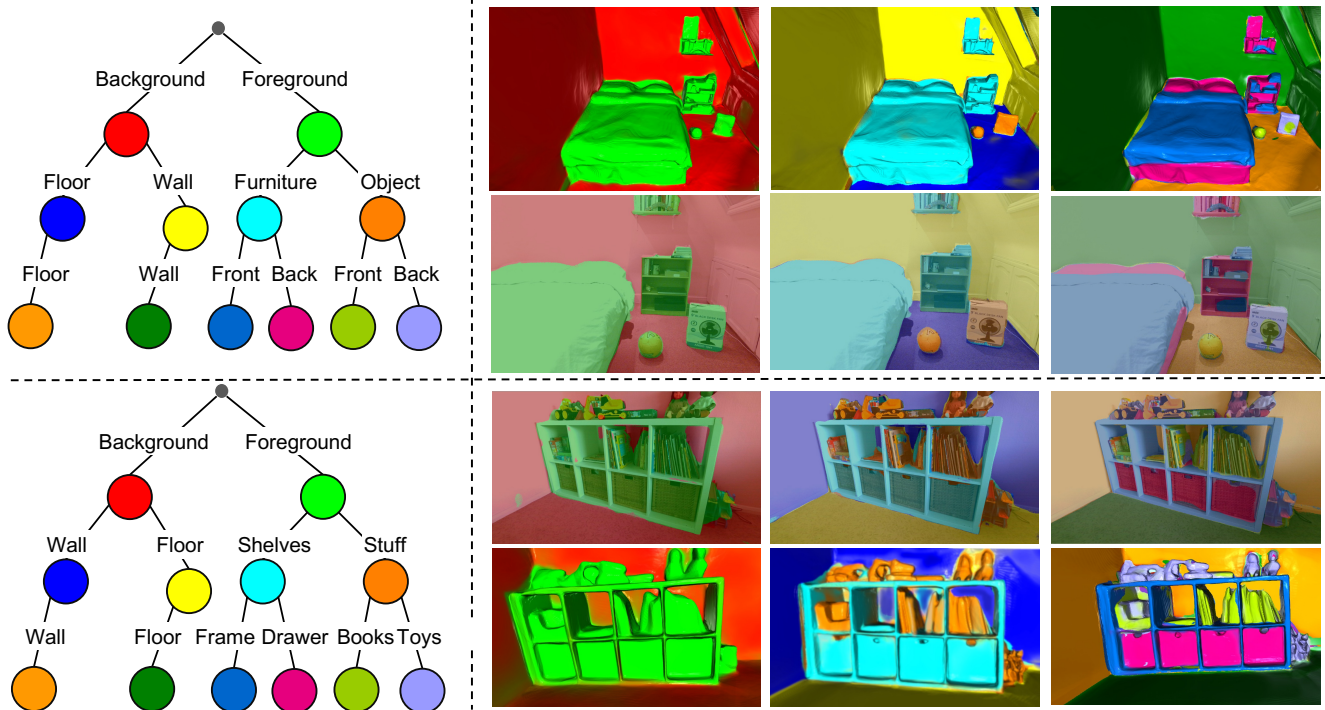
Figure 8. Binary tree as well as the segmentations at each level from the hierarchical mode of iLabel.
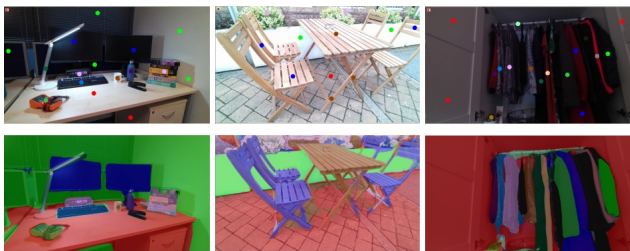


Figure 9. Ultra-efficient label propagation: iLabel produces high-quality segmentations of coherent 3D entities with very few user clicks, approximately 20–30 per scene.

between similar objects where only one has been labelled (e.g. (a) boxes on the bed, (b) food boxes and plastic cups and (c) toy dinosaurs). The table and chairs scene in Figure 11 (d) is especially interesting. Only four clicks are supplied: the label for the chair leg (blue) propagates to the leg of the table and the legs of the other chairs, while the table-top label (yellow) propagates to the seats of the chairs.

**Hierarchical scene segmentation** Figure 8 demonstrates iLabel's hierarchical mode. The colour-coded hierarchy (defined on-the-fly) is shown together with segmentations and scene reconstructions from each level. The results show the capacity of this representation to group objects at different levels, which has potential in applications where different tasks demand different groupings.



Figure 10. In removing the use of colour optimisation for scene reconstruction, only a few extra clicks are required to achieve a comparable quality of segmentation to that shown in Figure 7.



Figure 11. Generalisation: iLabel is able to transfer user labels to objects exhibiting similar properties. It is worth highlighting that the segmentation in (d) was achieved with only 4 clicks.

**Comparison to SemanticPaint** SemanticPaint (SPaint) [41] is currently the only comparable online interactive scene understanding system. With several distinct modes (labelling, propagation, training, predicting, correcting, smoothing), which do not operate simultaneously, users

6

(a) Input annotations     (b) **SPaint:** Initial strokes     (c) **SPaint:** Additional strokes     (d) **iLabel:** Initial strokes

Figure 12. Comparison results between iLabel and SemanticPaint for user annotations in (a). (b) SPaint results for initial strokes; (c) SPaint results after corrections; (d) iLabel segmentations obtained using only the input strokes in (a).

have to switch between modes repeatedly (with careful consideration given to the duration spent in each mode) to obtain optimal results. In contrast, iLabel presents a unified interface for scene reconstruction, whereby user interaction, label propagation, learning and prediction occur simultaneously. The more intuitive and simpler interface presented by iLabel means that high-quality segmentations are obtained with far fewer interactions and no expert knowledge/intuition.

Qualitative comparisons between iLabel and SPaint are given in Figures 12 and 13. Scenes with varying degrees of complexity were chosen to demonstrate the superiority of iLabel even in scenes well-suited to SPaint (e.g. final row Figure 12). For each scene in Figure 12, users annotated objects/regions with the strokes shown in (a). From these initial annotations only, iLabel was able to generate high-quality segmentations (Figure 12 (d)). In contrast, SPaint produced comparatively noisy and incomplete initial segmentations (Figure 12 (b)). Multiple mode switches and additional corrective strokes were required to generate the final SPaint results (Figure 12 (c)). We argue that the results produced by iLabel with only the initial user inputs ($< 10$ strokes), surpass those of SPaint after the additional user interactions. Figure 13 additionally illustrates the quality of

the 3D meshes generated by each technique, further highlighting the superiority of iLabel.
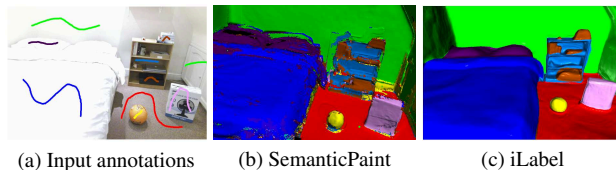


(a) Input annotations     (b) SemanticPaint     (c) iLabel

Figure 13. Qualitative comparison between iLabel and Semantic-Paint showing generated meshes.

## 4.2. Quantitative evaluation

We evaluate iLabel's 2D semantic segmentation performance in both user-interaction and automatic query generation modes, with varying numbers of clicks per scene, on the public datasets Replica [36] and ScanNet [7]. Both datasets are publicly available for research purposes under their licence. We report the mean Intersection Over Union (mIOU), averaged over ground truth labels remapped to NYU-13 class definitions.

**Baseline** While pre-trained segmentation models serve a different purpose than an interactive scene-specific system (to generalise to unseen scenes) we use them as a baseline

to demonstrate the labelling efficiency of our system. iLabel scales rapidly with the number of clicks and rapidly surpasses the pretrained model, even when this has been trained on very similar scenes
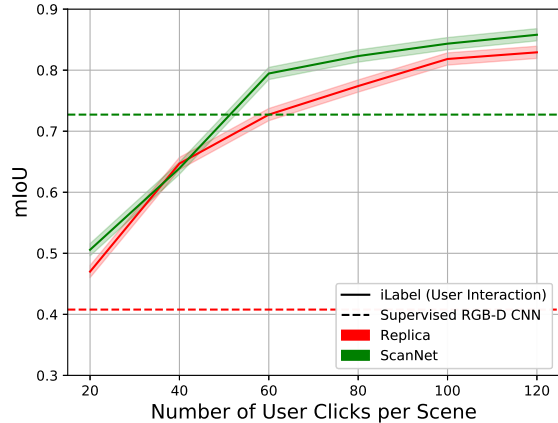
Performance is evaluated against SA-Gate [5] with a ResNet-101 DeepLabV3+ backbone [4], which is the current state-of-the-art in RGB-D segmentation. For Replica, we pre-train SA-Gate using the SUN-RGBD dataset [35] and fine-tune on our generated Replica sequences to avoid over-fitting. We adopt a leave-one-out strategy, whereby fine-tuning is performed independently for each test scene using the remaining Replica scenes. For ScanNet, we train SA-Gate directly on the official training sets, achieving 63.98% mIOU on the validation sets of 13 classes. Approximately 11k (9860 and 475 images for our SUN-RGBD training and validation splits, 900 images for Replica fine-tuning) and 25k training images were used for baseline CNN training on each Replica and ScanNet experiment, respectively. The ResNet-101 backbone is initialised with ImageNet pre-trained weights [31] through all the experiments. As per [5], depth maps use HHA encoding [11], before which fast depth completion [15] is used for hole-filling in ScanNet.

**Results** Figure 14a shows the performance of iLabel compared against the supervised RGB-D CNN baseline (dashed horizontal line) on 5 Replica scenes and 6 ScanNet scenes from the validation set. The Replica dataset is a low data regime with only 7 scenes used for fine tuning, which makes generalisation specially hard. iLabel is specially suited for this settings, and surpasses the baseline with only 20 clicks per scene. In the ScanNet dataset where much more data is available, iLabel reaches similar accuracy to the baseline with around 50 clicks, and continues to improve surpassing the baseline by 20% at 120 clicks.
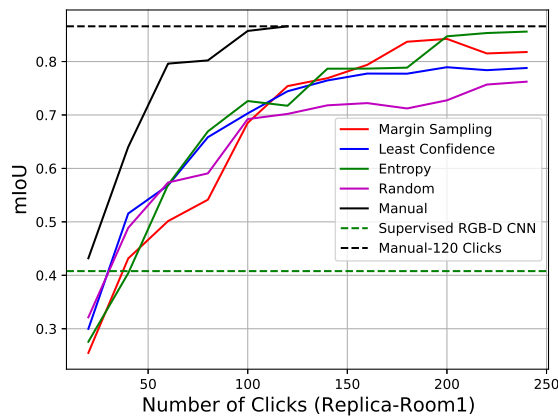
Figure 14b shows the effectiveness of automatic query generation, which opens the possibility for hands-free scene labelling, e.g., by voice command. As expected, this mode is less efficient than manual clicks and takes around 240 clicks to reach similar performance. We show how random uniform pixel sampling achieves a lower performance, specially when more labels have been added, highlighting the importance of uncertainty guided pixel selection.

## 5. Potential Negative Societal Impacts

As a visual perception module, iLabel can enable intelligent robots to label novel environments in an open-set manner with only minimal human input. As with any system designed to capture data, user privacy can be negatively impacted. Privacy concerns may be particularly important for iLabel as the scene representations it creates are compact ($\approx 1$ MB) making the process both portable and scalable. However, these same characteristics may also enable pos-



(a) Manual Interaction.



(b) Automatic Query Generation.

Figure 14. Quantitative evaluation of 2D semantic segmentation on the Replica and ScanNet datasets. Both interaction modes are evaluated and outperform supervised baselines with a small annotation budget.

itive technologies such as assistive robotics or inspection platforms that require semantic scene understanding.

## 6. Conclusion

We have shown that online, scene-specific training of a compact MLP model which encodes scene geometry, appearance and semantics allows ultra-sparse interactive labelling to produce accurate dense semantic segmentation, far surpassing the performance of standard pre-trained approaches. Despite promising results, our system's label propagation mechanism works well mainly for proximal regions and/or those sharing similar geometry or texture. A deeper understanding of this mechanism is necessary to enable better control of this process and to improve generalisation performance. As architectures and methods for neural implicit representation of scenes continue to improve, we expect these gains to be passed on to our labelling approach, and for tools like iLabel to become highly practical

for applications where users are able to teach AI systems efficiently about useful scene properties.

## 7. Acknowledgements

## References

[1] Pablo Arbelaez, Michael Maire, Charless Fowlkes, and Jitendra Malik. Contour detection and hierarchical image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 33(5):898–916, 2010. 2

[2] P. Arbeláez, J. Pont-Tuset, J. Barron, F. Marques, and J. Malik. Multiscale combinatorial grouping. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014. 2

[3] M. Bloesch, J. Czarnowski, R. Clark, S. Leutenegger, and A. 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

[4] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018. 2, 8

[5] Xiaokang Chen, Kwan-Yee Lin, Jingbo Wang, Wayne Wu, Chen Qian, Hongsheng Li, and Gang Zeng. Bi-directional cross-modality feature propagation with separation-and-aggregation gate for rgb-d semantic segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020. 5, 8

[6] Jan Czarnowski, Tristan Laidlow, Ronald Clark, and Andrew J Davison. Deepfactors: Real-time probabilistic dense monocular slam. *IEEE Robotics and Automation Letters*, 5(2):721–728, 2020. 2

[7] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. ScanNet: Richly-annotated 3d reconstructions of indoor scene. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 7

[8] Angela Dai, Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Christian Theobalt. BundleFusion: Real-time Globally Consistent 3D Reconstruction using On-the-fly Surface Re-integration. *ACM Transactions on Graphics (TOG)*, 36(3):24:1–24:18, 2017. 2

[9] A. J. Davison. Real-Time Simultaneous Localisation and Mapping with a Single Camera. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2003. 2

[10] Jakob Engel, Thomas Schoeps, and Daniel Cremers. LSD-SLAM: Large-scale direct monocular SLAM. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2014. 2

[11] 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. 8

[12] 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. 2

[13] Kunihito Kato Hiroaki Aizawa, Yukihiro Domae. Hierarchical pyramid representations for semantic segmentation. *arXiv preprint arXiv 2104.01792*, 2021. 2

[14] P Krähenbühl and Vladlen Koltun. Efficient Inference in Fully Connected CRFs with Gaussian Edge Potentials. In *Neural Information Processing Systems (NIPS)*, 2011. 2

[15] Jason Ku, Ali Harakeh, and Steven L Waslander. In defense of classical image processing: Fast depth completion on the cpu. In *Proceedings of the Canadian Conference on Computer and Robot Vision (CRV)*, 2018. 8

[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. 2

[17] Kevis-Kokitsi Maninis, Jordi Pont-Tuset, Pablo Arbeláez, and Luc Van Gool. Convolutional oriented boundaries. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016. 2

[18] David Martin, Charless Fowlkes, Doron Tal, and Jitendra Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2001. 2

[19] J. McCormac, R. Clark, M. Bloesch, A. J. Davison, and S. Leutenegger. Fusion++:volumetric object-level slam. In *Proceedings of the International Conference on 3D Vision (3DV)*, 2018. 2

[20] J. McCormac, A. Handa, A. J. Davison, and S. Leutenegger. SemanticFusion: Dense 3D semantic mapping with convolutional neural networks. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2017. 2

[21] Ondrej Miksik, Vibhav Vineet, Morten Lidegaard, Ram Prasaath, Matthias Nießner, Stuart Golodetz, Stephen L Hicks, Patrick Pérez, Shahram Izadi, and Philip HS Torr. The semantic paintbrush: Interactive 3d mapping and recognition in large outdoor spaces. In *ACM Conference on Human Factors in Computing Systems (CHI)*, 2015. 2

[22] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020. 1, 2

[23] R. Mur-Artal and J. D. Tardós. ORB-SLAM2: An Open-Source SLAM System for Monocular, Stereo, and RGB-D Cameras. *IEEE Transactions on Robotics (T-RO)*, 33(5):1255–1262, 2017. 2

[24] Yoshikatsu Nakajima, Byeongkeun Kang, Hideo Saito, and Kris Kitani. Incremental class discovery for semantic segmentation with rgbd sensing. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2019. 2

[25] Gaku Narita, Takashi Seno, Tomoya Ishikawa, and Yohsuke Kaji. Panopticfusion: Online volumetric semantic mapping at the level of stuff and things. In *Proceedings of the IEEE/RSJ Conference on Intelligent Robots and Systems (IROS)*, 2019. 2

[26] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohli, J. Shotton, S. Hodges, and A. Fitzgibbon. KinectFusion: Real-Time Dense Surface Mapping and Tracking. In *Proceedings of the International Symposium on Mixed and Augmented Reality (ISMAR)*, 2011. 2

[27] Duc Thanh Nguyen, Binh-Son Hua, Lap-Fai Yu, and Sai-Kit Yeung. A robust 3d-2d interactive tool for scene segmentation and annotation. *IEEE Transactions on Visualization and Computer Graphics*, 2017. 2

[28] M. Nießner, M. Zollhöfer, S. Izadi, and M. Stamminger. Real-time 3D Reconstruction at Scale using Voxel Hashing. In *Proceedings of SIGGRAPH*, 2013. 2

[29] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2

[30] Pengzhen Ren, Yun Xiao, Xiaojun Chang, Po-Yao Huang, Zhihui Li, Xiaojiang Chen, and Xin Wang. A survey of deep active learning. *arXiv preprint arXiv:2009.00236*, 2020. 4

[31] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. 8

[32] 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. 2

[33] Burr Settles. Active learning literature survey. 2009. 4

[34] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. In *Neural Information Processing Systems (NIPS)*, 2019. 2

[35] Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. SUN RGB-D: A RGB-D scene understanding benchmark suite. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 567–576, 2015. 8

[36] Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik Wijmans, Simon Green, Jakob J. Engel, Raul Mur-Artal, Carl Ren, Shobhit Verma, Anton Clarkson, Mingfei Yan, Brian Budge, Yajie Yan, Xiaqing Pan, June Yon, Yuyang Zou, Kimberly Leon, Nigel Carter, Jesus Briales, Tyler Gillingham, Elias Mueggler, Luis Pesqueira, Manolis Savva, Dhruv Batra, Hauke M. Strasdat, Renzo De Nardi, Michael Goesele, Steven Lovegrove, and Richard Newcombe. The Replica dataset: A digital replica of indoor spaces. *arXiv preprint arXiv:1906.05797*, 2019. 7

[37] J. Stückler and S. Behnke. Multi-resolution surfel maps for efficient dense 3d modeling and tracking. *Journal of Visual Communication and Image Representation*, 25(1):137–147, 2014. 2

[38] E. Sucar, S. Liu, J. Ortiz, and A. J. Davison. iMAP: Implicit Mapping and Positioning in Real-Time. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2021. 1, 2, 3, 4

[39] Edgar Sucar, Kentaro Wada, and Andrew Davison. NodeSLAM: Neural object descriptors for multi-view shape reconstruction. In *Proceedings of the International Conference on 3D Vision (3DV)*, 2020. 2

[40] N. Sünderhauf, T. T. Pham, Y. Latif, M. Milford, and I. Reid. Meaningful maps with object-oriented semantic mapping. In *Proceedings of the IEEE/RSJ Conference on Intelligent Robots and Systems (IROS)*, 2017. 2

[41] Julien Valentin, Vibhav Vineet, Ming-Ming Cheng, David Kim, Jamie Shotton, Pushmeet Kohli, Matthias Nießner, Antonio Criminisi, Shahram Izadi, and Philip Torr. Semanticpaint: Interactive 3d labeling and learning at your fingertips. *ACM Transactions on Graphics*, 34(5), November 2015. 2, 5, 6

[42] T. Whelan, R. F. Salas-Moreno, B. Glocker, A. J. Davison, and S. Leutenegger. ElasticFusion: Real-time dense SLAM and light source estimation. *International Journal of Robotics Research (IJRR)*, 35(14):1697–1716, 2016. 2

[43] Saining Xie and Zhuowen Tu. Holistically-nested edge detection. In *Proceedings of the International Conference on Computer Vision (ICCV)*, pages 1395–1403, 2015. 2

[44] S. Zhi, M. Bloesch, S. Leutenegger, and A. J. Davison. SceneCode: Monocular dense semantic reconstruction using learned encoded scene representations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2

[45] S. Zhi, T. Laidlow, S. Leutenegger, and A. J. Davison. In-Place Scene Labelling and Understanding with Implicit Scene Representation. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2021. 2, 3