

Visiting the Invisible: Layer-by-Layer Completed Scene Decomposition

Chuanxia Zheng · Duy-Son Dao · Guoxian Song · Tat-Jen Cham · Jianfei Cai

arXiv:2104.05367v1 [cs.CV] 12 Apr 2021

Abstract Existing scene understanding systems mainly focus on recognizing the visible parts of a scene, ignoring the intact appearance of physical objects in the real-world. Currently, image completion has aimed to create plausible appearance for the invisible regions, but requires a manual mask as input. In this work, we propose a higher-level scene understanding system to tackle both visible and invisible parts of objects and backgrounds in a given scene. Particularly, we built a system to decompose a scene into individual objects, infer their underlying occlusion relationships, and even automatically learn which parts of the objects are occluded that need to be completed. In order to disentangle the occluded relationships of all objects in a complex scene, we use the fact that the front object without being occluded is easy to be identified, detected, and segmented. Our system interleaves the two tasks of instance segmentation and scene completion through multiple iterations, solving for objects layer-by-layer. We first provide a thorough experiment using a new realistically rendered dataset with ground-truths for all invisible regions. To bridge the domain gap to real imagery where ground-truths are unavailable, we then train another model with the pseudo-ground-truths generated from our trained synthesis model. We demonstrate results on a wide variety of datasets and show significant improvement over the state-of-the-art. The code will be available at <https://github.com/lyndonzheng/VINV>.

Keywords Layered scene decomposition · Scene completion · Amodal instance segmentation · Instance depth order · Scene recomposition.

Chuanxia Zheng · Guoxian Song · Tat-Jen Cham
School of Computer Science and Engineering, Nanyang Technological University, Singapore.
E-mail: chuanxia001@e.ntu.edu.sg, guoxian001@e.ntu.edu.sg
astjcham@ntu.edu.sg

Duy-Son Dao · Jianfei Cai
Department of Data Science & AI, Monash University, Australia.
E-mail: duy.dao@monash.edu, jianfei.cai@monash.edu

1 Introduction

The vision community has made rapid advances in scene understanding tasks, such as object classification and localization (Girshick et al., 2014; He et al., 2015; Ren et al., 2015), scene parsing (Badrinarayanan et al., 2017; Chen et al., 2017; Long et al., 2015), instance segmentation (Chen et al., 2019; He et al., 2017; Pinheiro et al., 2015), and layered scene decomposition (Gould et al., 2009; Yang et al., 2010; Zhang et al., 2015). Despite their impressive performance, these systems deal only with *visible* parts of scenes without trying to exploit *invisible* regions, which results in an uncompleted representation of real objects.

In parallel, significantly progress for the generation task has been made with the emergence of deep generative networks, such as GAN-based models (Goodfellow et al., 2014; Gulrajani et al., 2017; Karras et al., 2019), VAE-based models (Kingma and Welling, 2014; Vahdat and Kautz, 2020; Van Den Oord et al., 2017), and flow-based models (Dinh et al., 2014, 2017; Kingma and Dhariwal, 2018). Empowered by these techniques, image completion (Iizuka et al., 2017; Yu et al., 2018; Zheng et al., 2019) and object completion (Ehsani et al., 2018; Ling et al., 2020; Zhan et al., 2020) have made it possible to create the plausible appearances for occluded objects and backgrounds. However, these systems depend on manual masks or visible ground-truth masks as input, rather than automatically understand the full scene.

In this paper, we aim to build a system that has the ability to *decompose* a scene into individual objects, *infer* their underlying occlusion relationships, and moreover *imagine* what occluded objects may look like, *while using only an image as input*. This novel task involves the classical recognition task of instance segmentation to predict the geometry and category of all objects in a scene, and the generation task of image completion to reconstruct invisible parts of objects and backgrounds.

To decompose a scene into instances with completed appearances in one pass is extremely challenging. This is be-

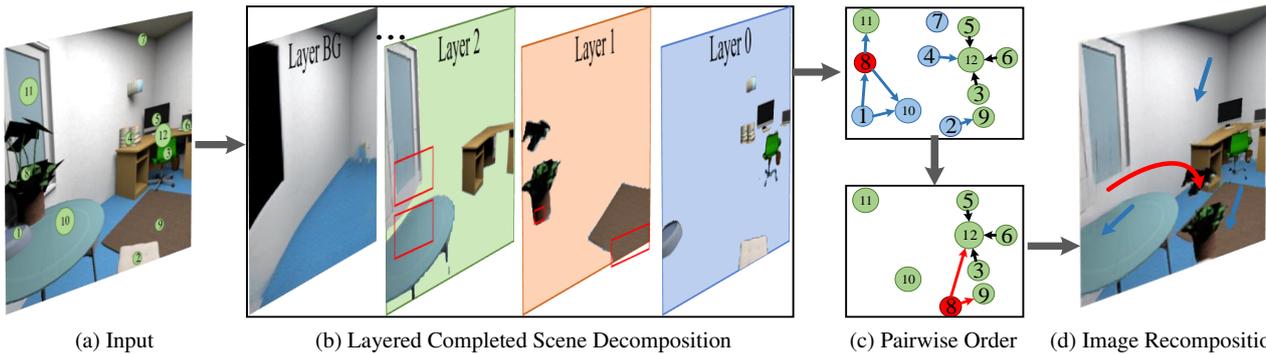


Fig. 1 Example results of scene decomposition and recomposition. (a) Input. (b) Our model structurally decomposes a scene into individual completed objects. Red rectangles highlight the original *invisible* parts. (c) The inferred pairwise order (top graph) and edited order (bottom graph) of the instances. Blue nodes indicate the deleted objects while the red node is the moved object. (d) The new recombined scene.

cause realistic natural scenes often consist of a vast collection of physical objects, with complex scene structure and occlusion relationships, especially when one object is occluded by multiple objects, or when instances have deep hierarchical occlusion relationships.

Our core idea is from the observation that *it is much easier to identify, detect and segment foreground objects than occluded objects*. Motivated by this, we propose a **Completed Scene Decomposition Network (CSDNet)** that learns to segment and complete each object in a scene layer-by-layer consecutively. As shown in Fig. 1, our layered scene decomposition network only segments the fully visible objects out in each layer (Fig. 1(b)). If the system is able to properly segment the foreground objects, it will automatically learn which parts of occluded objects are actually invisible that need to be filled in. The completed image is then passed back to the layered scene decomposition network, which can again focus purely on detecting and segmenting visible objects. As the interleaving proceeds, a structured instance depth order (Fig. 1(c)) is progressively derived by using the inferred absolute layer order. The thorough decomposition of a scene along with spatial relationships allows the system to freely recompose a new scene (Fig. 1(d)).

Another challenge in this novel task is the lack of data: there is no complex, realistic dataset that provides intact ground-truth appearance for originally occluded objects and backgrounds in a scene. While latest works (Li and Malik, 2016; Zhan et al., 2020) introduced a self-supervised way to tackle the amodal completion using only visible annotations, they can not do a fair quantitative comparison as no real ground-truths are available. To mitigate this issue, we constructed a high-quality rendered dataset, named **Completed Scene Decomposition (CSD)**, based on more than 2k indoor rooms. Unlike the datasets in (Dhamo et al., 2019; Ehsani et al., 2018), our dataset is designed to have more typical camera viewpoints, with near-realistic appearance.

As elaborated in Section 5.2, the proposed system performs well on this rendered dataset, both qualitatively and

quantitatively outperforming existing methods in completed scene decomposition, in terms of instance segmentation, depth ordering, and amodal mask and content completion. To further demonstrate the generalization of our system, we extend it to real datasets. As there is no ground truth annotations and appearance available for training, we created pseudo-ground-truths for real images using our model that is purely trained on **CSD**, and then fine-tuned this model accordingly. This model outperforms state-of-the-art methods (Qi et al., 2019; Zhan et al., 2020; Zhu et al., 2017) on amodal instance segmentation and depth ordering tasks, despite these methods being specialized to their respective tasks rather than our holistic completed scene decomposition task.

In summary, we propose a layer-by-layer scene decomposition network that jointly learns structural scene decomposition and completion, rather than treating them separately as the existing works (Dhamo et al., 2019; Ehsani et al., 2018; Zhan et al., 2020). To our knowledge, it is the first work that proposes to complete objects based on the global context, instead of tackling each object independently. To address this novel task, we render a high-quality rendered dataset with ground-truth for all instances. We then provide a thorough ablation study using this rendered dataset, in which we demonstrate that the method substantially outperforms existing methods that address the task in isolation. On real images, we improve the performance to the recent state-of-the-art methods by using pseudo-ground-truth as weakly-supervised labels. The experimental results show that our **CSDNet** is able to acquire a full decomposition of a scene, *with only an image as input*, which conduces to a lot of applications, *e.g.* object-level image editing.

The rest of the paper is organized as follows. We discuss the related work in Section 2, and describe our layer-by-layer CSDNet in detail in Section 3. In Section 4 we present our rendered dataset. We then show the experiment results on this synthetic dataset as well as the results on real-world images in Section 5, followed by a conclusion in Section 6.

Table 1 Comparison with related work based on three aspects: outputs, inputs and data. I: image, In: inmodal segmentation, O: occlusion order, SP: scene parsing, AB: amodal bounding box, AS: amodal surface, A: amodal segmentation, D: depth, IRGB: intact RGB object.

Paper	Outputs	Inputs	Data
	SP, O	I	LabelMe, PASVOC, others
(Yang et al., 2011)	In, O	I	PASVOC
(Tighe et al., 2014)	SP, O	I	LabelMe, SUN
(Zhang et al., 2015)	In, O	I	KITTI
(Guo and Hoiem, 2012)	AS	I	StreetScenes, SUN, others
(Kar et al., 2015)	AB	I	PASVOC, PAS3D
(Liu et al., 2016)	AS, O	I, D	NTUv2-D
(Li and Malik, 2016)	A	I, In	PASVOC
(Zhu et al., 2017)	A, O	I	COCOA (from COCO)
(Follmann et al., 2019)	A	I	COCOA, COCOA-cls, D2S
(Qi et al., 2019)	A	I	KINS (from KITTI)
(Hu et al., 2019)	A	I	Synthesis video
(Ehsani et al., 2018)	A, O, IRGB	I, In	DYCE, PAS3D
(Zhan et al., 2020)	A, O, IRGB	I, In	KINS, COCOA
(Ling et al., 2020)	A, IRGB	I, In	KINS
(Yan et al., 2019)	A, IRGB	I	Vehicle
(Burgess et al., 2019)	In, IRGB	I	Toy
(Dhamo et al., 2019)	A, D, IRGB	I	SUNCG, Stanford 2D-3D
Ours	A, O, IRGB	I	KINS, COCOA, SUNCG

2 Related Work

A variety of scene understanding tasks have previously been proposed, including layered scene decomposition (Yang et al., 2011), instance segmentation (He et al., 2017), amodal segmentation (Li and Malik, 2016), and scene parsing (Chen et al., 2017). In order to clarify the relationships of our work to the relevant literature, Table 1 gives a comparison based on three aspects: what the goals are, which information is used, and on which dataset is evaluated.

Layered scene decomposition for inmodal perception. The layered scene decomposition for visible regions has been extensively studied in the literature. Shade *et al.* (Shade et al., 1998) first proposed a representation called a layered depth image (LDI), which contains multiple layers for a complex scene. Based on this image representation that requires occlusion reasoning, the early works focused on ordering the semantic map as occluded and visible regions. Winn and Shotton (Winn and Shotton, 2006) proposed a LayoutCRF to model several occlusions for segmenting partially occluded objects. Gould *et al.* (Gould et al., 2009) decomposed a scene into semantic regions together with their spatial relationships. Sun *et al.* (Sun et al., 2010) utilized an MRF to model the layered image motion with occlusion ordering. Yang *et al.* (Yang et al., 2010, 2011) formulated a layered object detection and segmentation model, in which occlusion ordering for all detected objects was derived. This inferred order for all objects has been used to improve scene parsing (Tighe et al., 2014) through a CRF. Zhang *et al.* (Zhang et al., 2015) combined CNN and MRF to predict instance segmentation

with depth ordering. While these methods evaluate occlusion ordering, their main goal is to improve the inmodal perception accuracy for object detection, image parsing, or instance segmentation using the spatial occlusion information. In contrast to these methods, our method *not* only focuses on visible regions with structural inmodal perception, but also tries to solve for amodal perception. *i.e.* to learn *what is behind the occlusion*.

Amodal image/instance perception. Some initial steps have been taken toward amodal perception, exploring the invisible regions. Guo and Hoiem (Guo and Hoiem, 2012) investigated background segmentation map completion by learning relationships between occluders and background. Subsequently, (Liu et al., 2016) introduced the Occlusion-CRF to handle occlusions and complete occluded surfaces. Kar *et al.* (Kar et al., 2015) focused on amodal bounding box completion, where the goal is to predict the intact extent of the bounding box. The common attribute in these earlier amodal perception works is using piecemeal representations of a scene, rather than a full decomposition that infers the amodal shapes for all objects.

The success of advanced deep networks trained on large-scale annotated datasets has recently led to the ability to get more comprehensive representations of a scene. Instance segmentation (Dai et al., 2016; Li et al., 2017; Pinheiro et al., 2015, 2016) deal with detecting, localizing and segmenting all objects of a scene into individual instances. This task combines the classical object detection (Girshick, 2015; Girshick et al., 2014; He et al., 2015; Ren et al., 2015) and semantic segmentation (Badrinarayanan et al., 2017; Chen et al., 2017; Long et al., 2015). However, these notable methods typically segment the scene only into visible regions, and do *not* have an explicit structural representation of a scene. We believe a key reason is the lack of large-scale datasets with corresponding annotations for amodal perception and occlusion ordering. The widely used datasets, such as Pascal VOC 2012 (Everingham et al., 2010), NYU Depth v2 (Silberman et al., 2012), COCO (Lin et al., 2014), KITTI (Geiger et al., 2012), and CityScapes (Cordts et al., 2016), contain only annotations for the visible instances, purely aiming for 2D inmodal perception.

To mitigate the lack of annotated datasets, Li *et al.* (Li and Malik, 2016) presented a self-supervised approach by pasting occluders into an image. Although reasonable amodal segmentation results are shown, a quantitative comparison is unavailable due to the lack of ground-truth annotations for invisible parts. In more recent works, the completed masks for occluded parts are provided in COCOA (Zhu et al., 2017) and KINS (Qi et al., 2019), which are respectively a subset of COCO (Lin et al., 2014) and KITTI (Geiger et al., 2012). However, their annotations for invisible parts are manually labeled, which is highly subjective (Ehsani et al., 2018; Zhan et al., 2020). Furthermore, these datasets are mainly used

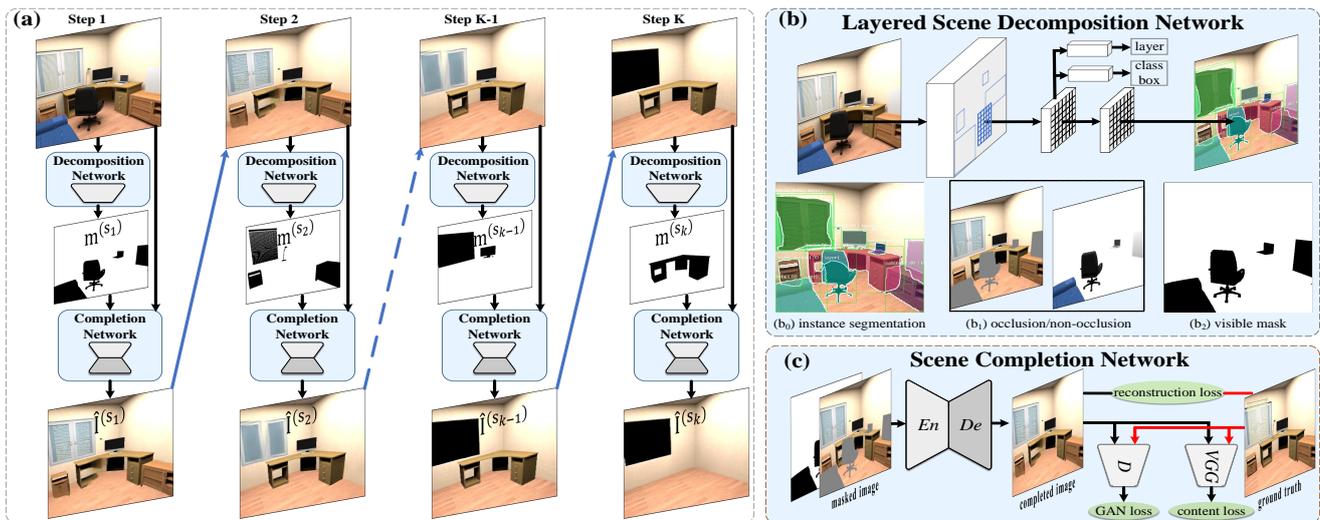


Fig. 2 An illustration of the CSDNet framework. (a) Overall layer-by-layer completed scene decomposition pipeline. In each step, the layered decomposition network selects out the fully visible objects. The completion network will complete the resultant holes with appropriate imagery. The next step starts again with the completed image. (b) The layered scene decomposition network estimates instance masks and binary occlusion relationships. (c) The completion network generates realistic content for the original *invisible* regions.

for the task of inferring amodal semantic maps and are not suitable for the task of RGB appearance completion, since the ground truth RGB appearance for occluded parts are not available. In contrast, we jointly address these two amodal perception tasks using our constructed CSD dataset.

Amodal perception for both mask and appearance. The problem of generating the amodal RGB appearance for the occluded parts is highly related to semantic image completion. The latest methods (Iizuka et al., 2017; Nazeri et al., 2019; Pathak et al., 2016; Yang et al., 2017; Yu et al., 2018; Zheng et al., 2019) extend GANs (Goodfellow et al., 2014) and CGANs (Mirza and Osindero, 2014) to address this task, generating new imagery to fill in partially erased image regions. However, they mainly focus on object removal, needing users to interactively annotate the objects to be removed.

SeGAN (Ehsani et al., 2018) involved an end-to-end network that sequentially infers amodal masks and generates complete RGB appearances for instances. The instance depth order is estimated by comparing the areas of the full and visible masks. PCNet (Zhan et al., 2020) used a self-supervised learning approach to recover masks and content using only visible annotations. However, these works mainly present results in which the ground truth visible mask is used as input, and are sensitive to errors in this visible mask. As stated in (Zhan et al., 2020), their focus is on amodal completion, rather than a scene understanding for amodal perception. While the recent work of Yan *et al.* (Yan et al., 2019) tried to visualize the invisible from a single input image, it only tackles the occluded “vehicle” category, for which there is much less variation in amodal shape and RGB appearance, and thus easier to model.

There are two recent works that attempt to learn structural scene decomposition with amodal perception. MONet (Burgess et al., 2019) combined an attention network and a CVAE (Kingma and Welling, 2014) for jointly modeling objects in a scene. While it is nominally able to do object appearance completion, this unsupervised method has only been shown to work on simple toy examples with minimal occlusions. Dhama *et al.* (Dhama et al., 2019) utilized Mask-RCNN (He et al., 2017) to obtain visible masks, and conducted RGBA-D completion for each object. However, depth values are hard to accurately estimate from a single image, especially in real images without paired depth ground-truths. Besides, they still considered the decomposition and completion separately. In practice, even if we use domain transfer learning for depth estimation, the pixel-level depth value for all objects are unlikely to be consistent in a real scene. Therefore, our method uses an instance-level occlusion order, called the “2.1D” model (Yang et al., 2011), to represent the structural information of a scene, which is easier to be inferred and manipulated.

3 Method

In this work, we aim to derive a higher-level structural decomposition of a scene. When given a single RGB image I , our goal is to decompose all objects in it and infer their fully completed RGB appearances, together with their underlying occlusion relationships (As depicted in Fig. 1). Our system is designed to carry out inmodal perception for *visible* structured instance segmentation, and also solve the amodal perception task of completing shapes and appearances for original *invisible* parts.

Instead of directly predicting the invisible content and decoupling the occlusion relationships of all objects at one pass, we use the fact that foreground objects are more easily identified, detected and segmented without occlusion. Our CSDNet decomposes the scene layer-by-layer. As shown in Fig. 2(a), in each step s_{k-1} , given an image $\mathbf{I}^{(s_{k-1})}$, the layered segmentation network creates masks as well as occlusion labels for all detected objects. Those instances classified as fully visible are extracted out and the scene completion network generates appropriate appearances for the invisible regions. The completed image $\hat{\mathbf{I}}^{(s_{k-1})}$ will then be re-submitted for layered instance segmentation in the next step s_k . This differs significantly from previous works (Burgess et al., 2019; Dhamao et al., 2019; Ehsani et al., 2018; Ling et al., 2020; Zhan et al., 2020), which do not adapt the segmentation process based on completion results.

Our *key novel insight* is that scene completion generates completed shapes for originally occluded objects by leveraging the global scene context, so that they are subsequently easier to be detected and segmented without occlusion. Conversely, better segmented masks are the cornerstones to complete individual objects by precisely predicting which regions are occluded. Furthermore, this interleaved process enables extensive *information sharing between these two networks*, to holistically solve for multiple objects, and produces a structured representation for a scene. This contrasts with existing one-pass methods (Burgess et al., 2019; Dhamao et al., 2019; Ehsani et al., 2018; Ling et al., 2020; Zhan et al., 2020), where the segmentation and completion are processed separately and instances are handled independently. Together with the benefit of occlusion reasoning, our system is able to explicitly learn *which parts of the objects and background are occluded that need to be completed*, instead of freely extending to arbitrary shapes.

3.1 Layered Scene Decomposition

As shown in Fig. 2, our layered scene decomposition network comprehensively detect objects in a scene. For each candidate instance, it outputs a class label, a bounding-box offset and an instance mask. The system is an extension of Mask-RCNN (He et al., 2017), which consists of two main stages. In the first stage, the image is passed to a *backbone network* (e.g. ResNet-50-FPN (Lin et al., 2017)) and next to a *region proposal network* (RPN (Ren et al., 2015)) to get object proposals. In the second stage, the network extracts features using *RoIAlign* from each candidate box, for passing to object classification and mask prediction. We refer readers to (Chen et al., 2019; He et al., 2017) for details.

To determine if an object is fully visible *or* partially occluded, a new parallel branch for this binary occlusion classification is added, as shown in Fig. 2(b). This decomposition is done consecutively layer-by-layer, where at each

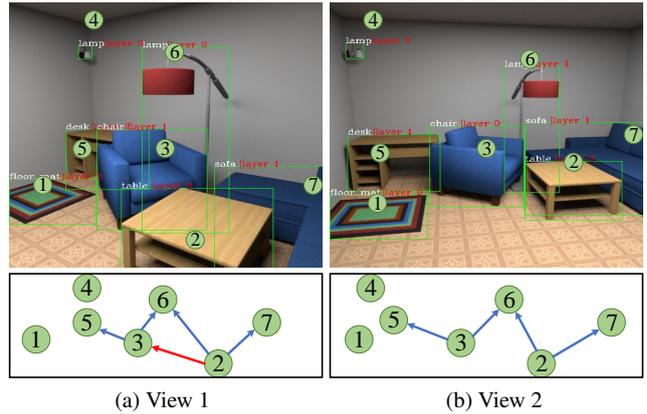


Fig. 3 Instance depth order representation. Top images show absolute layer order (Qi et al., 2019) in different views. Bottom directed graphs give the pairwise order between objects.

step it is applied to a single RGB derived from the counterpart scene completion step. While only binary decisions are made in each step, after a full set of iterations, a comprehensive layered occlusion ordering is obtained. The following parts describe how this is done, looking first at the instance depth order representation, followed by how occlusion head is designed.

Instance depth order representation. Absolute layer order and pairwise occlusion order are two standard representations for occlusion reasoning in a scene (Sun et al., 2010; Tighe et al., 2014). As shown in Fig. 3, the definition for our *absolute layer order* \mathcal{L} follows (Qi et al., 2019), where fully visible objects are labeled as 0, while other occluded objects have 1 order higher than the highest-order instance occluding them (see top images in Fig. 3). We interpret the *pairwise occlusion order matrix* as a directed graph $G = (\Omega, W)$ (see bottom graphs in Fig. 3), where Ω is a discrete set of all instances with number N , and W is a $N \times N$ matrix. $W_{i,j}$ is the occlusion relationship of instance i to instance j . We use three numbers to encode the order — $\{-1$: occluded, 0 : no relationship, 1 : front}. For example, the chair (instance #3) is occluded by the table (instance #2), so the pairwise order for the chair is $W_{3,2} = -1$, while the pairwise order for the table is inversely labeled as $W_{2,3} = 1$.

Occlusion head. In practice we find it hard to directly predict these instance depth orders. The absolute layer order index $l \in \mathcal{L}$ cannot be predicted purely from local features in a bounding box, since it depends on the global layout of all objects in a scene. Furthermore, this index is very sensitive to viewpoints, e.g. in Fig. 3, the desk (instance #5) is occluded by only one object (instance #3) in both views, but the absolute layer order indices of the desk are different: “2” vs “1”. In contrast, pairwise ordering G captures the occlusion relationships between pairs of objects, but all pairs have to be analyzed, leading to scalability issues in the current instance segmentation network. As R-CNN-based system cre-

ates 2,000 candidate objects, this pairwise analysis requires building an unwieldy $2k \times 2k$ features.

We circumvent these problems as our occlusion classifier only predicts a *binary occlusion label*: $\{0, 1\}$ in each step, where 0 is fully visible, and 1 is occluded, following the setting of absolute layer order. During training, each ground-truth binary occlusion label is determined from the pairwise order of the actual objects present in the scene (see details in the Appendix). The occlusion head in our layered scene decomposition network is a *fc* layer, which receives aligned features from each RoI and predicts the binary occlusion label.

Decomposition Loss. The multi-task loss function for layered scene decomposition is defined as follows:

$$L_{\text{decomp}} = \sum_{t=1}^T \alpha_t (L_{\text{cls}}^t + L_{\text{bbox}}^t + L_{\text{mask}}^t + L_{\text{occ}}^t) + \beta L_{\text{seg}} \quad (1)$$

where classification loss L_{cls}^t , bounding-box loss L_{bbox}^t , mask loss L_{mask}^t and semantic segmentation loss L_{seg} are identical to those defined in HTC (Chen et al., 2019), and L_{occ}^t is the occlusion loss at the cascade refined stage t (three cascade refined blocks in HTC (Chen et al., 2019)), using binary cross-entropy loss (Long et al., 2015) for each RoI.

3.2 Visiting the Invisible by Exploring Global Context

In our solution, we treat visiting the invisible as a *semantic image completion* (Pathak et al., 2016) problem. As illustrated in Fig. 2, in step s_{k-1} , after removing the front visible instances, the given image $\mathbf{I}^{(s_{k-1})}$ is degraded to become $\mathbf{I}_m^{(s_{k-1})}$. Our goal is to generate appropriate content to complete these previously *invisible* regions (being occluded) for the next layer $\mathbf{I}^{(s_k)}$. Unlike existing methods that complete each object independently (Burgess et al., 2019; Dhano et al., 2019; Ehsani et al., 2018; Ling et al., 2020; Yan et al., 2019; Zhan et al., 2020), our model completes multiple objects in each step layer-by-layer, such that the information from earlier scene completions propagate to later ones. The global scene context is utilized in each step.

To visit the invisible, it is critical to know which parts are invisible that need to be completed. The general image completion methods use manually interactive masks as input, which differs from our goal. Recent related works (Ehsani et al., 2018; Ling et al., 2020; Zhan et al., 2020) depend on the ground-truth visible masks as input to indicate which parts are occluded. In contrast, our system selects out fully visible objects and automatically learns which parts are occluded in each step. The holes left behind explicitly define the occluded regions for remaining objects, and thus the completed shapes for remaining objects must be deliberately *restricted to these regions*, instead of being allowed to grow freely using only the predicted visible masks.

We use the PICNet (Zheng et al., 2019) framework to train our completion network. While the original PICNet was designed for diversity, here we only want to obtain the best result closest to the ground-truth. Therefore, we only use the encoder-decoder structure, and eschew the random sampling aspect.

Completion Loss. The overall scene completion loss function is given by

$$L_{\text{comp}} = \alpha_{\text{rec}} L_{\text{rec}} + \alpha_{\text{ad}} L_{\text{ad}} + \alpha_{\text{per}} L_{\text{per}} \quad (2)$$

where reconstruction loss L_{rec} and adversarial loss L_{ad} are identical to those in PICNet (Zheng et al., 2019). The perceptual loss $L_{\text{per}} = |\mathbf{F}^{(l)}(\hat{\mathbf{I}}^{(s_k)}) - \mathbf{F}^{(l)}(\mathbf{I}^{(s_k)})|$ (Johnson et al., 2016), based on a pretrained VGG-19 (Simonyan and Zisserman, 2014), is the l_1 distance of features \mathbf{F} in l -th layer between the generated image $\hat{\mathbf{I}}^{(s_k)}$ and ground-truth $\mathbf{I}^{(s_k)}$.

3.3 Inferring Instance Pairwise Occlusion Order

As discussed in Section 3.1, absolute layer order \mathcal{L} is sensitive to errors. If one object is incorrectly selected as a front object in an earlier step, objects behind it will have their absolute layer order incorrectly shifted. Hence in keeping with prior works (Ehsani et al., 2018; Zhan et al., 2020), we use the pairwise occlusion order $G = (\Omega, W)$ to represent our final instance occlusion relationships for evaluation.

Given a single image \mathbf{I} , our model decomposes it into instances with completed RGB appearances $A_{\Omega}^{S_K}$. Here, A denotes the amodal perception instance (inclusive of both mask and appearance), Ω specifies instances in the scene, and S_K indicates which layers are the instances in (selected out in step s_k). When two segmented amodal masks $A_{\omega_i}^{s_i}$ and $A_{\omega_j}^{s_j}$ overlap, we infer their occlusion relationship based on the order of the object-removing process, formally:

$$W_{\omega_i, \omega_j} = \begin{cases} 0 & \text{if } O(A_{\omega_i}^{s_i}, A_{\omega_j}^{s_j}) = 0 \\ 1 & \text{if } O(A_{\omega_i}^{s_i}, A_{\omega_j}^{s_j}) > 0 \text{ and } s_i < s_j \\ -1 & \text{if } O(A_{\omega_i}^{s_i}, A_{\omega_j}^{s_j}) > 0 \text{ and } s_i \geq s_j \end{cases} \quad (3)$$

where $O(A_{\omega_i}^{s_i}, A_{\omega_j}^{s_j})$ is the area of overlap between instances ω_i and ω_j . If they do not overlap, they share no pairwise depth-order relationship in a scene. If there is an overlap and the instance ω_i is first selected out with a smaller layer order, the inferred pairwise order is $W_{\omega_i, \omega_j} = 1$; otherwise it is labeled as $W_{\omega_i, \omega_j} = -1$. Hence the instance occlusion order only depends on the order (selected out step) of removal between the two instances, and do not suffer from shift errors.

3.4 Training on Real Data with Pseudo Ground-truth

Real-world data appropriate for a completed scene decomposition task is difficult to acquire, because ground truth

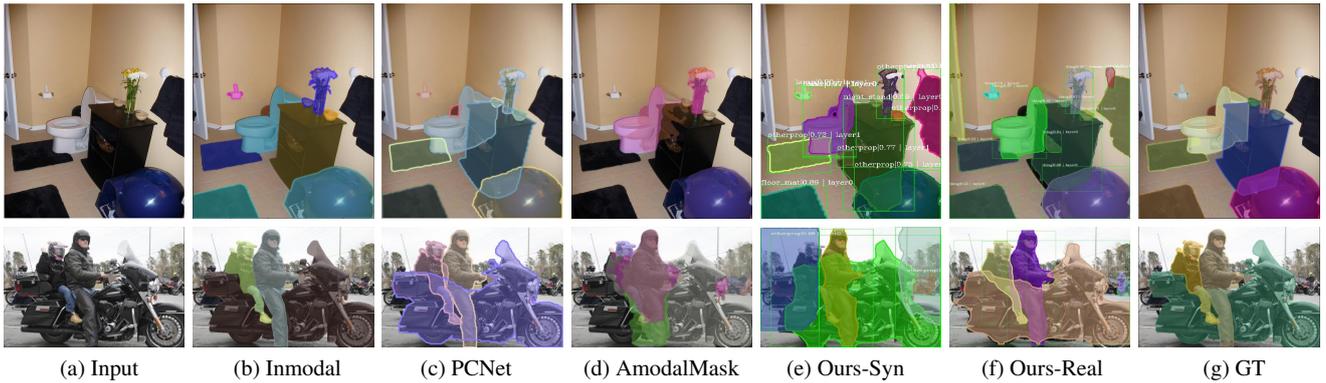


Fig. 4 Amodal instance segmentation results on the COCOA validation set. Our model trained on synthetic dataset (Ours-syn) achieves visually reasonable results in the similar real indoor scenes (example in top row), but it fails in some dissimilar real scenes (example in bottom row). After training on the real data with “pseudo ground-truths”, the model (Ours-Real) performs much better. Note that, unlike the PCNet (Zhan et al., 2020) that need visible inmodal ground-truth masks as input, our system decomposes a scene using only a RGB image.

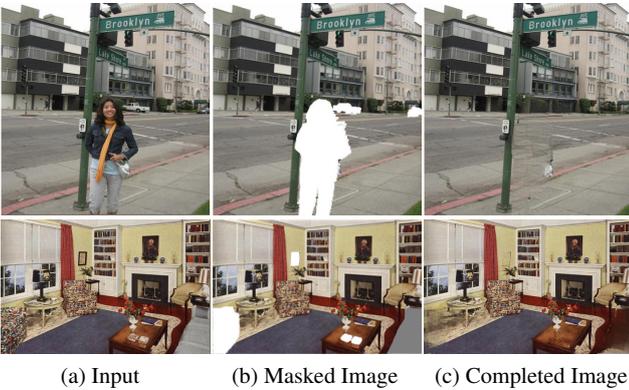


Fig. 5 Pseudo RGB ground-truth. (a) Input. (b) Masked image by selecting out the fully visible objects. (c) Pseudo ground-truth generated from our model trained on synthetic data.

shapes and RGB appearance for occluded parts are hard to collect without very extensive manual interventions, *e.g.* deliberate physical placement and removal of objects. Although our proposed model trained on the high-quality rendered data achieves visually reasonable results in some real scenes that share similarities to the rendered dataset (*e.g.* indoor scene in top row of Fig. 4), it does not generalize well to dissimilar real scenes (*e.g.* outdoor scene in bottom row of Fig. 4). These are caused by: 1) differences in labeled categories between synthetic and real datasets, especially between indoor and outdoor scenes; and 2) inconsistencies between synthetically trained image completion of masked regions and fully visible real pixels.

One alternative is to simply use an image completion network trained only on real images. From our experience, this performs poorly in a scene decomposition task. The reason is that while real-trained image completion methods are able to create perceptually-pleasing regions and textures for a single completion, they do not appear to have the ability to adhere to consistent object geometry and boundaries when de-occluding, which is crucial for object shape completion.

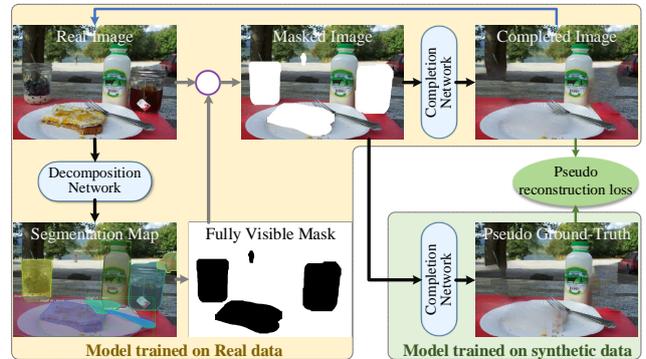


Fig. 6 Training pipeline for real images. We introduce a semi-supervised learning method for real data by providing *pseudo* RGB ground-truth for originally invisible regions.

As a result, errors accumulate even more dramatically as the decomposition progresses.

Our *key motivating insight* is this: instead of training the model entirely without ground-truth in the completion task, we train it in a semi-supervised learning manner, exploiting the scene structure and object shape knowledge that has been gained in our synthetically-trained CSDNet. As shown in Fig. 5, this synthetic completion model is able to generate visually adequate appearance, but more importantly it is better able to retain appropriate geometry and shapes. We can use this to guide the main image completion process in real-world data, while allowing a GAN-based loss to increase the realism of the output.

Specifically, for a real image \mathbf{I} , we first train the layered decomposition network using the manual annotated amodal labels. In a step, after segmenting and selecting out the foreground objects, we obtain $\hat{\mathbf{I}}_{syn}^{(s_k)} = G(\mathbf{I}_m^{(s_k)}; \theta_{syn})$ to serve as “pseudo ground-truth” (green box in Fig. 6) through the completion model trained on synthetic data. We then train the completion network $G(\mathbf{I}_m^{(s_k)}; \theta_{real})$ using the loss function of equation (2) by comparing the output $\hat{\mathbf{I}}_{real}^{(s_k)}$ to “pseudo ground-

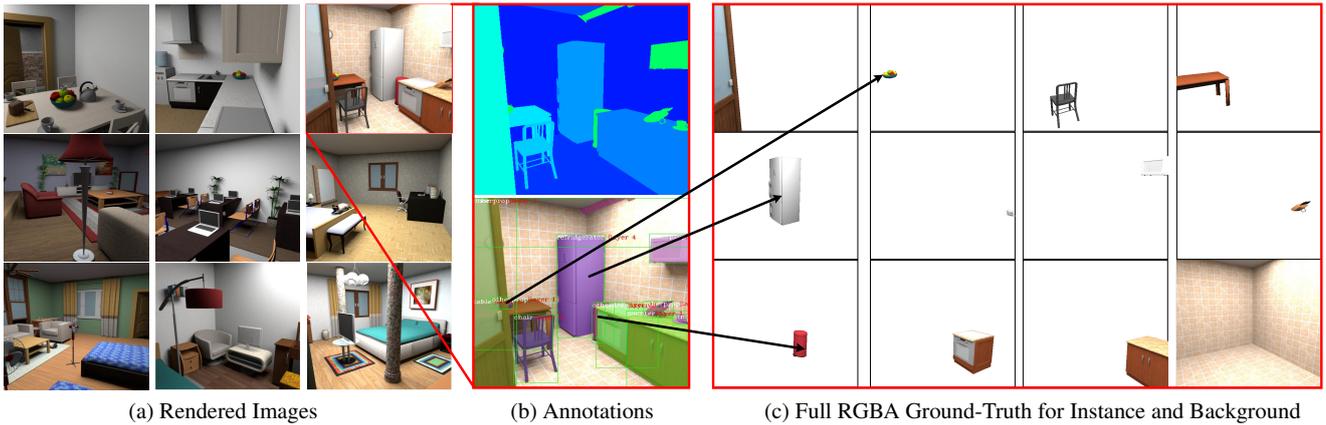


Fig. 7 Our rendered dataset. (a) High quality rendered RGB images. (b) Semantic map and instance annotations with bbox, category, ordering and segmentation map. (c) Intact RGBA ground-truth for instances and background.

truth” $\hat{\mathbf{I}}_{syn}^{(s_k)}$. Like (Sengupta et al., 2020), we also reduce the weights of reconstruction loss L_{rec} and perceptual loss L_{per} to encourage the output to be biased towards the real image distribution via the discriminator loss L_{ad} . It is worth noticing that the completed image is *passed back* to the layered decomposition network in the next layer, where the decomposition loss L_{decomp} in equation (1) will be backpropagated to the completion network. This connection allows the completion network to *learn to complete real world objects that might not be learned through the synthetic data*.

4 Synthetic Data Creation

Large datasets with complete ground-truth appearances for all objects are very limited. Burgess *et al.* (Burgess et al., 2019) created the *Objects Room dataset*, but only with toy objects. Ehsani *et al.* (Ehsani et al., 2018) and Dhano *et al.* (Dhano et al., 2019) rendered synthetic datasets. However, the former only includes 11 rooms, with most viewpoints being atypical of real indoor images. The latter’s OpenGL-rendered dataset appears to have more typical viewpoints with rich annotations, but the OpenGL-rendered images have low realism. Recently, Zhan *et al.* (Zhan et al., 2020) explored the *amodal completion* task through self-supervised learning without the need of amodal ground-truth. However, a fair quantitative comparison is not possible as no appearance ground-truth is available for invisible parts.

To mitigate this issue, we rendered a realistic dataset with Maya (Autodesk Maya, 2019). We can train the supervised model and test the unsupervised model on this synthetic data with masks and RGB appearance ground-truths for all occluded parts.

Data Rendering. Our rendered data is based on a total of 10.2k views inside over 2k rooms (CAD models from SUNCG (Song et al., 2017)) with various room types and lighting environments (see Fig. 7(a)). To select the typical viewpoints,

we first sampled many random positions, orientations and heights for the camera. Only when a view contains at least 5 objects will we render the image and the corresponding ground-truth of each instance. To avoid an excessive rendering workload, we separately rendered each isolated object, as well as the empty room, as shown in Fig. 7(c). This allows us to then freely create the ground-truths of each layer by compositing these individual objects and background using the instance occlusion order. The rendering details and statistics of the dataset can be found in the Appendix.

Data Annotation. Each rendered scene is accompanied by a global semantic map and dense annotations for all objects. As shown in Fig. 7(b) and Fig. 7(c), the intact RGB appearances are given, as well as categories (the 40 classes in NYUDv2 (Nathan Silberman and Fergus, 2012)), bounding boxes and masks for complete objects, as well as for only the visible regions. Furthermore, the absolute layer order and pairwise occlusion order shown in Fig. 3 are also defined in our rendered dataset.

5 Experiments

5.1 Setup

Datasets. We evaluated our system on three datasets: **COCOA** (Zhu et al., 2017), **KINS** (Qi et al., 2019) and the rendered **CSD**. **COCOA** is annotated from COCO2014 (Lin et al., 2014), a large scale natural image datasets, in which 5,000 images are selected to manually label with pairwise occlusion orders and amodal masks. **KINS** is derived from the outdoor traffic dataset KITTI (Geiger et al., 2012), in which 14,991 images were labeled with absolute layer orders and amodal masks. **CSD** is our rendered synthetic dataset, which contains 8,298 images, 95,030 instances for training and 1,012 images, 11,648 instances for testing. We conducted thorough experiments and ablation studies to assess

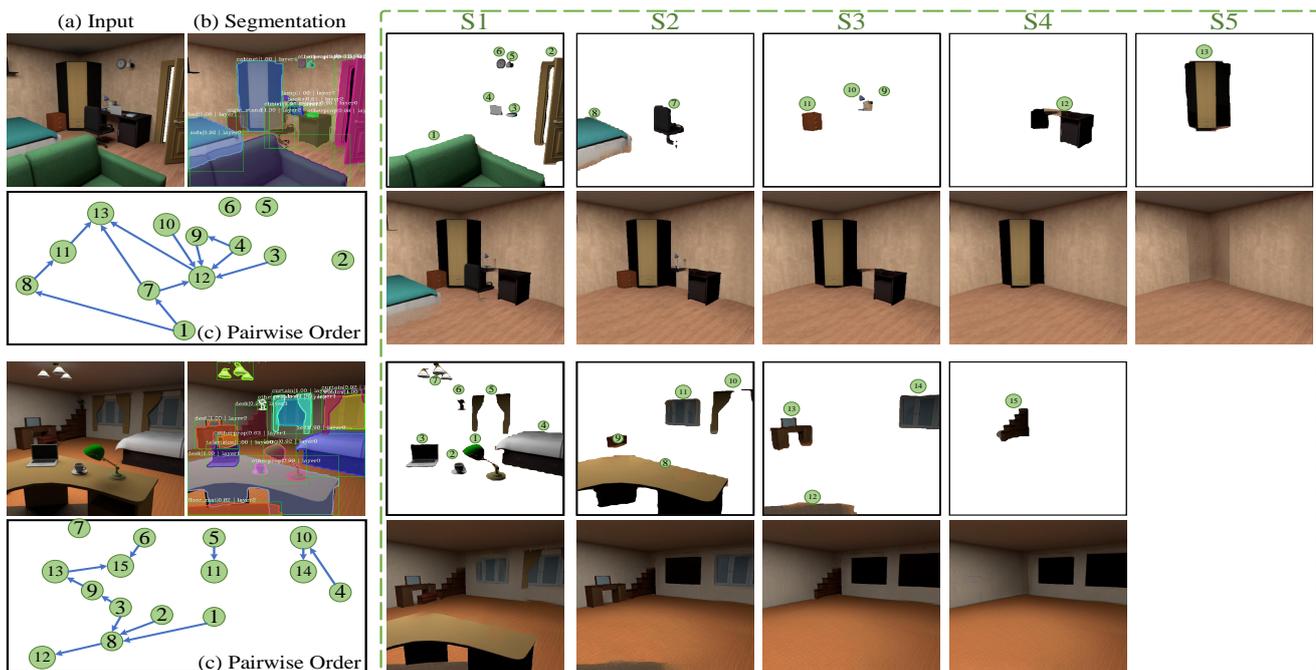


Fig. 8 Layer-by-Layer Completed Scene Decomposition results on rendered CSD testing set. (a) Input RGB images. (b) Final amodal instance segmentation. (c) Inferred directed graph for pairwise order. (d) Columns labeled S1-5 show the decomposed instances (top) and completed scene (bottom) based on the predicted non-occluded masks. Note that the originally occluded invisible parts are filled in with realistic appearance.

the quality of the completed results for invisible appearance estimation (since the in-the-wild datasets lack ground-truth for the occluded parts).

Metrics. For amodal instance segmentation, we report the standard COCO metrics (Lin et al., 2014), including AP (average over IoU thresholds), AP_{50} , AP_{75} , and AP_S , AP_M and AP_L (AP at different scales). Unless otherwise stated, the AP is for mask IoU. For appearance completion, we used RMSE, SSIM and PSNR to evaluate the quality of generated images. All images were normalized to the range $[0, 1]$.

Since the occlusion order is related to the quality of instance segmentation, we defined a novel metric for evaluating the occlusion order that uses the previous benchmark criterion for instance segmentation. Specifically, given a pairwise occlusion order $G = (\Omega, W)$ predicted by the model, we only evaluate the order for these valid instances that have IoU with ground-truth masks over a given threshold. For instance, if we set the threshold as 0.5, the predicted instance ω will be evaluated when we can identify a matched ground-truth mask with $\text{IoU} \geq 0.5$. Hence we can measure the **occlusion average precision (OAP)** as assessed with different thresholds.

5.2 Results on Synthetic CSD Dataset

We first present results that we obtained from our framework when experimenting on our synthetic CSD dataset.

5.2.1 Main Results

Completed scene decomposition. We show the qualitative results of CSDNet in Fig. 8. Given a single RGB image, the system has learned to decompose it into semantically complete instances (e.g. counter, table, window) and the background (wall, floor and ceiling), while completing RGB appearance for *invisible* regions. Columns labeled S1-5 show the completed results layer-by-layer. In each layer, fully visible instances are segmented out, and after scene completion some previously occluded regions become fully visible in the next layer, e.g. the table in the second example. The final amodal instance segmentation results shown in Fig. 8(b) consist of the fully visible amodal masks in each layer. Note that unlike MONet (Burgess et al., 2019), our model does not need predefined slots. The process will stop when it is unable to detect any more objects.

Amodal instance segmentation. We compare CSDNet to the state-of-the-art methods in amodal instance segmentation in Table 2. As the existing works Mask-RCNN (He et al., 2017) and HTC (Chen et al., 2019) are aimed at imodal perception for visible parts, we retrained their models for amodal perception task, by providing amodal ground-truths. We also retrained MLC (Qi et al., 2019) on our rendered dataset, which are the latest work for amodal perception. For PCNet (Zhan et al., 2020), we used the predicted visible mask as input, rather than the original visible ground-truth annotations. While the HTC (Chen et al., 2019) im-

Table 2 Amodal Instance Segmentation on CSD testing sets. Mask-RCNN (He et al., 2017) and HTC (Chen et al., 2019) are the state-of-the-arts of the COCO segmentation challenges. MLC (Qi et al., 2019) is the latest amodal instance segmentation work for outdoor scene. PCNet (Zhan et al., 2020) is the self-supervised amodal completion work. The **CSDNet-gt** holds same training environment as **CSDNet**, but is tested with completed ground-truths images I^* in each step. Best results used ground-truth annotations are marked with *, while best results only used RGB images are in bold.

	SegNet	box AP	mask AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L
Mask-RCNN (He et al., 2017)	Mask-RCNN	51.3	46.8	67.2	50.6	14.5	43.0	49.9
MLC (Qi et al., 2019)	Mask-RCNN	52.3	47.2	67.5	50.9	14.7	43.8	50.2
PCNet (Zhan et al., 2020)	Mask-RCNN	-	43.6	59.1	43.4	11.5	40.4	46.0
HTC (Chen et al., 2019)	HTC	52.9	47.3	65.9	51.6	12.2	41.3	51.0
MLC (Qi et al., 2019)	HTC	53.6	47.9	66.1	52.3	13.1	41.9	51.7
PCNet (Zhan et al., 2020)	HTC	-	45.7	60.6	49.2	10.2	39.3	48.4
CSDNet	Mask-RCNN	52.6	48.7	66.2	53.1	15.7	42.8	52.2
CSDNet	HTC	56.3	50.3	67.7	53.4	17.4	44.2	53.1
CSDNet-gt	Mask-RCNN	54.9	53.1	66.5	56.9	21.4*	49.9	57.0
CSDNet-gt	HTC	60.3*	56.0*	67.9*	59.3*	19.6	53.4*	59.5*

Table 3 Instance depth ordering on CSD testing sets. We report the pairwise depth ordering on occluded instance pairs OAP_{occ} . I^* = ground-truth completed image in each step s^* , V_{gt} = visible ground-truth mask, \hat{V}_{pred} = visible predicated mask, and \hat{F}_{pre} = full (amodal) predicated mask. layer order¹ only predicts the occlusion / non-occlusion labels in the original image (the first step in our model).

	Inputs		Ordering Algorithm	OAP	OAP ₅₀	OAP ₇₅	OAP ₈₅	OAP _S	OAP _M	OAP _L
	Amodal	Ordering								
SeGAN (Ehsani et al., 2018)	$I + V_{gt}$	$V_{gt} + \hat{F}_{pre}$	IoU Area	68.4	-	-	-	-	-	-
SeGAN (Ehsani et al., 2018)	$I + \hat{V}_{pred}$	$V_{gt} + \hat{F}_{pre}$	IoU Area	66.1	50.2	65.6	70.4	10.6	65.0	63.8
HTC + MLC (Qi et al., 2019)	I	$V_{gt} + \hat{F}_{pre}$	IoU Area	76.5	70.3	77.1	79.8	11.6	69.8	78.2
HTC + MLC (Qi et al., 2019)	I	$\hat{F}_{pre} + \text{layer}$	layer order ¹	51.9	44.3	50.8	54.6	11.7	60.8	46.2
HTC + PCNet (Zhan et al., 2020)	$I + \hat{V}_{pred}$	$V_{gt} + \hat{F}_{pre}$	IoU Area	70.8	56.9	71.3	76.0	11.3	67.1	68.6
CSDNet	I	\hat{F}_{pre}	Area	44.7	45.3	45.7	45.1	17.4	34.5	41.5
CSDNet	I	\hat{F}_{pre}	Y-axis	62.0	60.1	61.2	62.7	63.4	58.6	66.1
CSDNet	I	$V_{gt} + \hat{F}_{pre}$	IoU Area	80.7	77.2	81.0	82.9	61.1	73.7	80.5
CSDNet	I	$\hat{F}_{pre} + \text{layer}$	layer order	81.7	76.6	80.9	84.6	15.7	75.9	82.6
CSDNet-gt	I^*	$\hat{F}_{pre} + \text{layer}$	layer order	88.9*	85.2*	88.5*	90.1*	49.6	84.3*	89.9*

proves on Mask-RCNN’s (He et al., 2017) bounding box AP by about 1.6 points by refining the bounding box offsets in three cascade steps, the improvement for amodal mask segmentation is quite minor at 0.5 points. We believe this is an inherent limitation of methods that attempt amodal segmentation of occluded objects directly without first reasoning about occluding objects and masking their image features, as such the front objects’ features will distract the network. In contrast, our CSDNet is able to improve the amodal mask segmentation accuracy by a relative 6.3% with the same *backbone segmentation network* (HTC), by jointing segmentation and completion with layer-by-layer decomposition.

To further demonstrate that better completed images improve amodal segmentation, we consider a scenario with a completion oracle, by using ground-truth appearances to repair the occluded parts in each step. This is denoted as the **CSDNet-gt**, for which amodal instance segmentation accuracy increases from 47.3% to 56.0% (relative 18.4% improvement). We also note that, while the **CSDNet-gt** using Mask-RCNN achieves lower bounding box score than our **HTC CSDNet** (“54.9” vs “56.3”), the mask accuracy is

much higher (“53.1” vs “50.3”). This suggests that amodal segmentation benefits from better completed images.

Instance depth ordering. Following (Zhu et al., 2017), we report the pairwise instance depth ordering for correctly detected instances in Table 3. The original SeGAN and PCNet used ground-truth visible masks V_{gt} as input. For a fair comparison, we first retrained them on our synthetic data using the same segmentation network (HTC (Chen et al., 2019)) for all models. After predicting amodal masks, we assessed various instance depth ordering algorithms: two baselines proposed in AmodalMask (Zhu et al., 2017) of ordering by *area*¹ and by *y-axis* (amodal masks closest to image bottom in front), ordering by *incremental area* defined as the *IoU area* between visible and amodal masks², and our ordering by absolute *layer order* (Section 3.3).

As can be seen in Table 3, all instantiations of our model outperformed baselines as well as previous models. Unlike

¹ We used the heuristic in PCNet (Zhan et al., 2020) — larger masks are ordered in front for KINS, and behind for COCOA and CSD.

² See details in (Zhan et al., 2020), where the visible ground-truth masks V_{gt} are used for ordering.

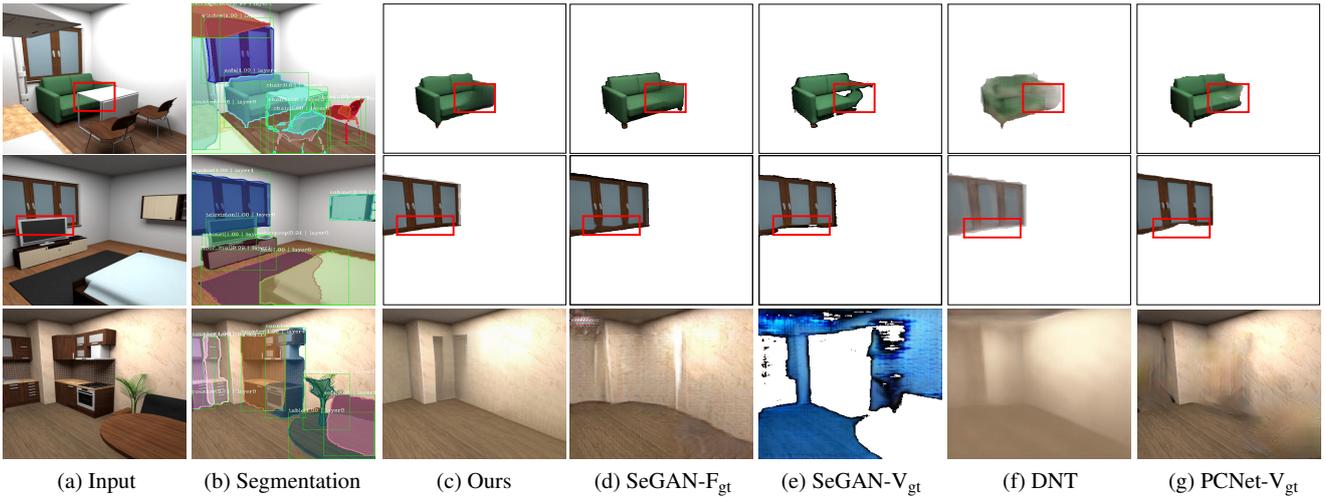


Fig. 9 Results for **Visiting the Invisible**. We show the input image, our amodal instance segmentation results, and the objects and background we try to visit. The red rectangles highlight the previously *invisible* regions of occluded objects.

Table 4 **Object Completion**. F_{gt} = full ground-truth mask, V_{gt} = visible ground-truth mask. For methods provided with F_{gt} , we only evaluate the completion networks.

	C1- F_{gt}				C2		
	RMSE	SSIM	PSNR		RMSE	SSIM	PSNR
SeGAN	0.1246	0.8140	21.42	C2a- V_{gt}	0.2390	0.6045	16.03
PCNet	0.1129	0.8267	23.16		0.2483	0.5931	15.54
DNT	0.1548	0.7642	20.32	C2b	0.2519	0.5721	15.10
PICNet	0.0927	0.8355	28.81		0.1401	0.7730	24.71
CSDNet	0.0614	0.9179	35.24		0.0914	0.8768	30.45

SeGAN (Ehsani et al., 2018) and PCNet (Zhan et al., 2020), our final model explicitly predicts the occlusion labels of instances, which improved the OAP substantially. While MLC (Qi et al., 2019) predicts the instance occlusion order in a network, it only contains one layer for binary occlusion / non-occlusion labeling. In contrast, our method provides a fully structural decomposition of a scene in multiple steps. Additionally, we observed that our model achieves better performance with a higher IoU threshold for selecting the segmented mask (closer match to the ground-truth masks). We further observed that the occlusion relationships of small objects are difficult to infer in our method. However, the *Y-axis* ordering method had similar performance under various metrics as it only depends on the locations of objects. Note that our depth ordering does *not* rely on any ground-truth that is used in (Ehsani et al., 2018; Zhan et al., 2020).

Object completion. We finally evaluated the quality of generated appearances. We compared our results to those from SeGAN (Ehsani et al., 2018), Dhamao *et al.* (Dhamao et al., 2019) (abbrev. as DNT), PCNet (Zhan et al., 2020) and PICNet (Zheng et al., 2019) (original point-attention) in Table 4. We evaluated different conditions of: C1) when the ground-truth full mask F_{gt} is provided to all methods, C2a) when

the ground-truth visible mask V_{gt} is the input to SeGAN and PCNet, and C2b) when an RGB image is the only input to other methods. C2a- V_{gt} is considered because SeGAN and PCNet assumes that a predefined mask is provided as input.

In C1- F_{gt} , CSDNet substantially outperformed the other methods. In C2, even when given only RGB images *without* ground-truth masks, our method worked better than SeGAN and PCNet with V_{gt} . One important reason for the strong performance of CSDNet is the *occlusion reasoning* component, which constrains the completed shapes of partly occluded objects based on the global scene context and other objects *during testing*.

Qualitative results are visualized in Fig. 9. We noted that SeGAN worked well only when ground-truth amodal masks F_{gt} were available to accurately label which parts were *invisible* that needed filling in, while DNT generated blurry results from simultaneously predicting RGB appearance and depth maps in one network, which is not an ideal approach (Zamir et al., 2018). The PCNet (Zhan et al., 2020) can not correctly repair the object shape as it trained without ground-truth object shape and appearance. Our CSDNet performed much better on background completion, as it only masked fully visible objects in each step instead of all objects at a go, so that *earlier completed information propagates to later steps*.

5.2.2 Ablation Studies

To demonstrate the two tasks can contribute to a better scene understanding system by jointly optimizing, instead of solving them isolated, we ran a number of ablations.

Does better completion help decomposition? We show quantitative results for a fixed decomposition network (layered HTC (Chen et al., 2019) with two completion methods in

Table 5 Ablations for joint optimization. In each table, we fixed one model for one subtask and trained different models for the other subtask. Better performance in one task can improve the performance in the other, which demonstrates the joint training of two tasks with layer-by-layer decomposition contributes to each other.

(a) **Effect of different completion methods on instance segmentation (HTC-based decomposition)**. “sep” = separate training of the 2 networks, “w/o” = without any completion, and “end” = joint training.

	train	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L
gt	-	56.0*	67.9*	59.3*	19.6*	53.4*	59.5*
w/o	-	36.8	52.6	38.3	10.8	31.6	38.2
PICNet-point	sep	40.8	63.0	43.5	12.6	37.2	43.7
PICNet-patch	sep	43.8	60.7	46.6	11.0	36.3	45.5
PICNet-point	end	47.7	63.2	50.6	14.9	41.7	51.3
PICNet-patch	end	50.3	67.7	53.4	17.4	44.2	53.1

Table 5(a). Without any completion (“w/o”), segmented results were naturally bad (“36.8” vs “50.3”) as it had to handle empty masked regions. More interestingly, even if advanced methods were used to generate visual completion, the isolated training of the decomposition and completion networks led to degraded performance. This suggests that even when generated imagery looks good visually, there is still a domain or semantic gap to the original visible pixels, and thus flaws and artifacts will affect the next segmentation step. The PICNet with patch attention provides better completed results than the original point attention PICNet (Zheng et al., 2019), resulting in a large improvement (“50.3” vs “47.7”) of amodal instance segmentation.

Does better decomposition help completion? To answer this, we report the results of using different scene segmentation networks with a same completion network (Patch-attention PICNet (Zheng et al., 2019)) in Table 5(b). We also first considered the ideal situation that ground-truth segmentation masks were provided in each decomposition step. As shown in Table 5(b), the completion quality significantly improved (RMSE: “0.0614”, SSIM: “0.9179” and PSNR: “35.24”) as occluded parts were correctly pointed out and the completion network precisely knows which parts need to be completed. HTC (Chen et al., 2019) provided better instance masks than Mask-RCNN (He et al., 2017), which resulted in more accurately completed scene imagery. The best results were with end-to-end jointly training.

(b) **Effect of different decomposition methods on scene completion (Patch-Attention PICNet)**. Better scene decomposition improved scene completion, likewise with joint training being most effective.

	train	RMSE	SSIM	PSNR
gt	-	0.0614*	0.9179*	35.24*
M-RCNN(He et al., 2017)	sep	0.1520	0.7781	22.34
HTC (Chen et al., 2019)	sep	0.1496	0.7637	26.75
M-RCNN(He et al., 2017)	end	0.1345	0.7824	27.31
HTC (Chen et al., 2019)	end	0.0914	0.8768	30.45

5.3 Results on Real Datasets

We now assess our model on real images. Since the ground-truth appearances are unavailable, we only provide the visual *scene manipulation* results in Section 5.4, instead of quantitative results for *invisible completion*.

Completed scene decomposition. In Fig. 10, we visualize the layer-by-layer completed scene decomposition results on real images. Our CSDNet is able to decompose a scene into completed instances with correct ordering. The originally occluded invisible parts of “suitcase”, for instance, is completed with full shape and realistic appearance. Note that, our system is a fully scene understanding method that only takes an image as input, without requiring the other manual annotations as (Ehsani et al., 2018; Zhan et al., 2020).

Amodal instance segmentation. Next, we compare with state-of-the-art methods on amodal instance segmentation. Among these, AmodalMask (Zhu et al., 2017) and ORCNN (Follmann et al., 2019) were trained for the COCOA dataset, MLC (Qi et al., 2019) works for the KINS dataset, and PC-Net (Zhan et al., 2020) is focused on amodal completion (mask completion) rather than amodal instance segmentation (requiring precise visible masks). For a fair comparison, when these methods do not provide results on a given dataset, we trained their models using publicly released code.

Table 6 shows that our results (34.8 mAP and 32.2 mAP) are 0.4 points and 0.6 points higher than the recent MLC

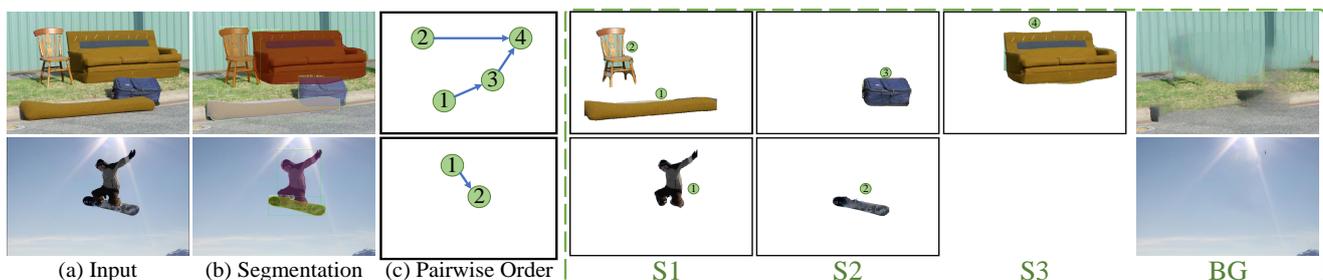


Fig. 10 Layer-by-layer completed scene decomposition on natural images. (a) Inputs. (b) Final amodal instance segmentation. (c) Inferred directed graph for pairwise occlusion order. (d) Columns labeled S1-3 show the decomposed instances with completed appearance in each step.

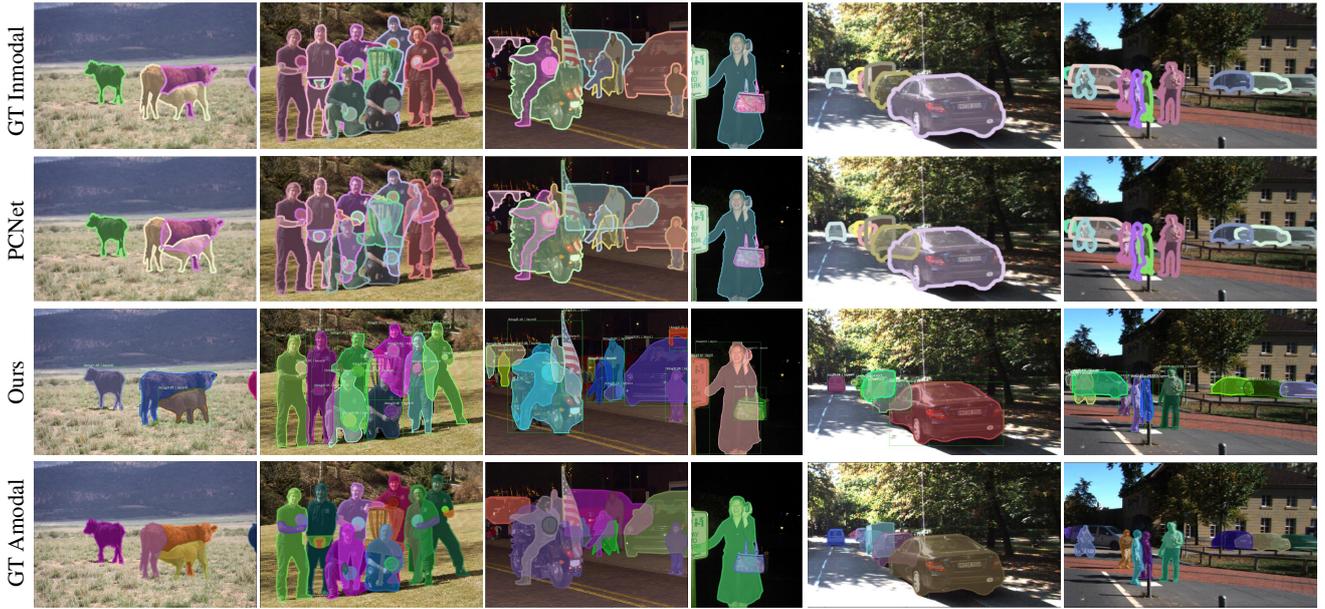


Fig. 11 Amodal instance segmentation results on natural images. Our CSDNet learns to predict the intact mask for the occluded objects (*e.g.* animals and human). Note that, unlike PCNet (Zhan et al., 2020), our model does *not* depend on the visible mask (first row) as input. Hence it can handle some objects without ground-truth annotation, such as two ‘humans’ in the third column and the ‘smartphone’ in the fourth column.

Table 6 Amodal Instance Segmentation on COCOA and KINS sets. The gray color shows results reported in existing works and the others are our reported results by using the released codes and our CSDNet.

	Inputs	SegNet	COCOA (%mAP)	KINS (%mAP)
Amodel (Zhu et al., 2017)	I	Sharp (Pinheiro et al., 2016)	7.7	-
Mask-RCNN (He et al., 2017)	I	Mask-RCNN (He et al., 2017)	31.8	29.3
ORCNN (Follmann et al., 2019)	I	Mask-RCNN (He et al., 2017)	33.2	29.0
MLC (Qi et al., 2019)	I	Mask-RCNN (He et al., 2017)	34.0	31.1
MLC (Qi et al., 2019)	I	HTC (Chen et al., 2019)	34.4	31.6
PCNet (Zhan et al., 2020)	$I + \hat{V}_{pred}$	Mask-RCNN (He et al., 2017)	30.3	28.6
PCNet (Zhan et al., 2020)	$I + \hat{V}_{pred}$	HTC (Chen et al., 2019)	32.6	30.1
CSDNet	I	Mask-RCNN (He et al., 2017)	34.1	31.5
CSDNet	I	HTC (Chen et al., 2019)	34.8	32.2

Table 7 Instance depth ordering on COCOA and KINS sets. The blue rows show the results that uses ground-truth annotations as inputs.

	Ordering Inputs	Ordering Algorithm	COCOA (OAP)	KINS (OAP)
OrderNet (Zhu et al., 2017)	$I + \hat{F}_{gt}$	Network	88.3	94.1
PCNet (Zhan et al., 2020)	$V_{gt} + \hat{F}_{pre}$	IoU Area	84.6	86.0
MLC (Qi et al., 2019)	$V_{gt} + \hat{F}_{pre}$	IoU Area	80.3	82.3
CSDNet	$V_{gt} + \hat{F}_{pre}$	IoU Area	84.7	86.4
MLC (Qi et al., 2019)	$\hat{V}_{pred} + \hat{F}_{pre}$	IoU Area	74.2	80.2
MLC (Qi et al., 2019)	$\hat{F}_{pre} + \text{layer}$	layer order ¹	66.5	71.8
PCNet (Zhan et al., 2020)	$\hat{V}_{pred} + \hat{F}_{pre}$	IoU Area	72.4	79.8
CSDNet	$\hat{V}_{pred} + \hat{F}_{pre}$	IoU Area	75.4	81.6
CSDNet	$\hat{F}_{pre} + \text{layer}$	layer order	80.9	82.2

using the same segmentation structure (HTC) in COCOA and KINS, respectively. PCNet (Zhan et al., 2020) considers amodal perception in two steps and assumes that visible masks are available. We note that their mAP scores were very high when the visible ground-truth masks were pro-

vided. This is because all initial masks were matched to the annotations (without detection and segmentation errors for instances, as shown in Fig. 11). However, when we used a segmentation network to obtain visible masks \hat{V}_{pred} for PCNet, the amodal instance segmentation results became lower

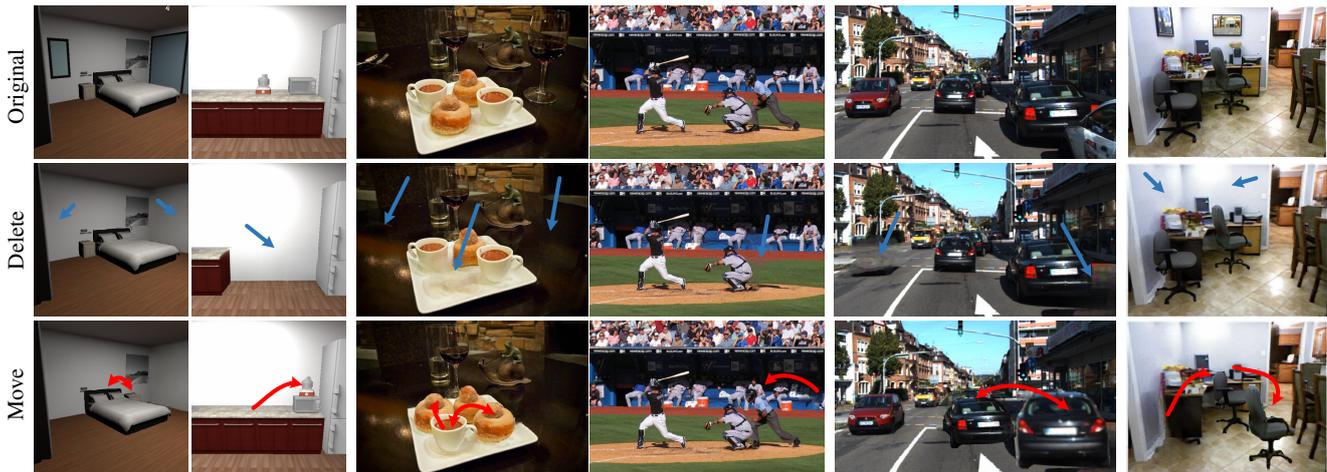


Fig. 12 Free editing based on the results of our system on images from various datasets. Note that our method is able to automatically detect, segment and complete the objects in the scene, *without the need for manual interactive masking*, with interactive operations limited to only “delete” and “drag-and-drop”. The blue arrows show object removal, while red arrows show object moving operations. We can observe that the originally *invisible* regions are fully visible after editing.

than other methods, suggesting that it is much harder to segment a visible mask and then complete it.

In Fig. 11, we compare our CSDNet and PCNet (Zhan et al., 2020). PCNet only completes the given visible annotated objects which had visible masks. In contrast, our CSDNet produces more contiguous amodal instance segmentation maps even for some unlabeled objects, for instance, the two “humans” in the third column.

Instance depth ordering. Finally, we report the instance depth ordering results in Table 7. In order to compare with existing work, we considered two settings: ground-truths provided (blue rows in Table 7), and only RGB images given. The OrderNet obtained the best results as the ground-truth full masks F_{gt} were given. We note that PCNet and our model achieved comparable performance when the visible ground-truth masks were used. Note that, we only used V_{gt} for depth ordering, while PCNet utilized the visible mask as input for both mask predication and depth ordering. Furthermore, when no ground-truth annotation was provided as input, our model performed better than MLCand PCNet.

5.4 Applications

We illustrate some image editing / re-composition applications of this novel task, after we learned to decompose a scene into isolated completed objects together with their spatial occlusion relationships. In Fig. 12, we visualize some recomposed scenes on various datasets, including our CSD, real COCOA (Zhu et al., 2017), KITTI (Qi et al., 2019) and NYU-v2 (Nathan Silberman and Fergus, 2012).

In these cases, we directly modified the positions and occlusion ordering of individual objects. For instance, in

the first bedroom example, we *deleted* the “window”, and *moved* the “bed” and the “counter”, which resulted in also *modifying* their occlusion order. Note that all original *invisible* regions were filled in with reasonable appearance. We also tested our model on real NYU-v2 (Nathan Silberman and Fergus, 2012) images which do *not* belong to any of the training sets used. As shown in the last column of Fig. 12, our model was able to detect and segment the object and complete the scene. The “picture”, for instance, is deleted and filled in with background appearance.

6 Conclusions

Building on previous inmodal and amodal instance perception work, we explored a higher-level structure scene understanding task that aims to decompose a scene into semantic instances, with completed RGB appearance and spatial occlusion relationships. We presented a layer-by-layer CSDNet, an iterative method to address this novel task. The main motivation behind our method is that fully visible objects, at each step, can be relatively easily detected and selected out without concern about occlusion. To do this, we simplified this complex task to two subtasks: instance segmentation and scene completion. We analyzed CSDNet and compared it with recent works on various datasets. Experimental results show that our model can handle an arbitrary number of objects and is able to generate the appearance of occluded parts. Our model outperformed current state-of-the-art methods that address this problem in one pass. The thorough ablation studies on synthetic data demonstrate that the two subtasks can contribute to each other through the layer-by-layer processing.

Acknowledgements This research was supported by the BeingTogether Centre, a collaboration between Nanyang Technological University (NTU) Singapore and University of North Carolina (UNC) at Chapel Hill. The BeingTogether Centre was supported by the National Research Foundation, Prime Minister's Office, Singapore under its International Research Centres in Singapore Funding Initiative. This research was also conducted in collaboration with Singapore Telecommunications Limited and partially supported by the Singapore Government through the Industry Alignment Fund – Industry Collaboration Projects Grant and the Monash FIT Start-up Grant.

References

- Autodesk Maya (2019) Autodesk Maya. <https://www.autodesk.com/products/maya/overview>
- Badrinarayanan V, Kendall A, Cipolla R (2017) Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine intelligence* 39(12):2481–2495
- Burgess CP, Matthey L, Watters N, Kabra R, Higgins I, Botvinick M, Lerchner A (2019) MONet: Unsupervised scene decomposition and representation. *arXiv preprint arXiv:190111390*
- Chen K, Pang J, Wang J, Xiong Y, Li X, Sun S, Feng W, Liu Z, Shi J, Ouyang W, et al. (2019) Hybrid task cascade for instance segmentation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp 4974–4983
- Chen LC, Papandreou G, Kokkinos I, Murphy K, Yuille AL (2017) Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 40(4):834–848
- Cordts M, Omran M, Ramos S, Rehfeld T, Enzweiler M, Benenson R, Franke U, Roth S, Schiele B (2016) The cityscapes dataset for semantic urban scene understanding. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp 3213–3223
- Dai J, He K, Sun J (2016) Instance-aware semantic segmentation via multi-task network cascades. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp 3150–3158
- Dhamo H, Navab N, Tombari F (2019) Object-driven multi-layer scene decomposition from a single image. In: *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*
- Dinh L, Krueger D, Bengio Y (2014) Nice: Non-linear independent components estimation. *arXiv preprint arXiv:14108516*
- Dinh L, Sohl-Dickstein J, Bengio S (2017) Density estimation using real nvp. In: *International Conference on Learning Representations*
- Ehsani K, Mottaghi R, Farhadi A (2018) SeGAN: Segmenting and generating the invisible. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp 6144–6153
- Everingham M, Van Gool L, Williams CK, Winn J, Zisserman A (2010) The pascal visual object classes (voc) challenge. *International journal of computer vision* 88(2):303–338
- Follmann P, Nig RK, Rtinger PH, Klostermann M, Ttger TB (2019) Learning to see the invisible: End-to-end trainable amodal instance segmentation. In: *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, IEEE, pp 1328–1336
- Geiger A, Lenz P, Urtasun R (2012) Are we ready for autonomous driving? the kitti vision benchmark suite. In: *2012 IEEE Conference on Computer Vision and Pattern Recognition*, IEEE, pp 3354–3361
- Girshick R (2015) Fast R-CNN. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp 1440–1448
- Girshick R, Donahue J, Darrell T, Malik J (2014) Rich feature hierarchies for accurate object detection and semantic segmentation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp 580–587
- Goodfellow I, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, Courville A, Bengio Y (2014) Generative adversarial nets. In: *Advances in Neural Information Processing Systems*, pp 2672–2680
- Gould S, Fulton R, Koller D (2009) Decomposing a scene into geometric and semantically consistent regions. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp 1–8
- Gulrajani I, Ahmed F, Arjovsky M, Dumoulin V, Courville AC (2017) Improved training of wasserstein gans. In: *Advances in neural information processing systems*, pp 5767–5777
- Guo R, Hoiem D (2012) Beyond the line of sight: labeling the underlying surfaces. In: *European Conference on Computer Vision*, Springer, pp 761–774
- He K, Zhang X, Ren S, Sun J (2015) Spatial pyramid pooling in deep convolutional networks for visual recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 37(9):1904–1916
- He K, Gkioxari G, Dollár P, Girshick R (2017) Mask R-CNN. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp 2961–2969
- Hu YT, Chen HS, Hui K, Huang JB, Schwing AG (2019) Sail-vos: Semantic amodal instance level video object segmentation—a synthetic dataset and baselines. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp 3105–3115
- Iizuka S, Simo-Serra E, Ishikawa H (2017) Globally and locally consistent image completion. *ACM Transactions on Graphics (TOG)* 36(4):107

- Johnson J, Alahi A, Fei-Fei L (2016) Perceptual losses for real-time style transfer and super-resolution. In: Proceedings of the European Conference on Computer Vision, pp 694–711
- Kar A, Tulsiani S, Carreira J, Malik J (2015) Amodal completion and size constancy in natural scenes. In: Proceedings of the IEEE International Conference on Computer Vision, pp 127–135
- Karras T, Laine S, Aila T (2019) A style-based generator architecture for generative adversarial networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 4401–4410
- Kingma DP, Dhariwal P (2018) Glow: Generative flow with invertible 1x1 convolutions. In: Advances in neural information processing systems, pp 10215–10224
- Kingma DP, Welling M (2014) Auto-encoding variational bayes. In: editor (ed) Proceedings of the International Conference on Learning Representations (ICLR)
- Li K, Malik J (2016) Amodal instance segmentation. In: Proceedings of the European Conference on Computer Vision, Springer, pp 677–693
- Li Y, Qi H, Dai J, Ji X, Wei Y (2017) Fully convolutional instance-aware semantic segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 2359–2367
- Lin TY, Maire M, Belongie S, Hays J, Perona P, Ramanan D, Dollár P, Zitnick CL (2014) Microsoft coco: Common objects in context. In: European conference on computer vision, Springer, pp 740–755
- Lin TY, Dollár P, Girshick R, He K, Hariharan B, Belongie S (2017) Feature pyramid networks for object detection. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 2117–2125
- Ling H, Acuna D, Kreis K, Kim SW, Fidler S (2020) Variational amodal object completion. *Advances in Neural Information Processing Systems* 33
- Liu C, Kohli P, Furukawa Y (2016) Layered scene decomposition via the Occlusion-CRF. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 165–173
- Long J, Shelhamer E, Darrell T (2015) Fully convolutional networks for semantic segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 3431–3440
- Mirza M, Osindero S (2014) Conditional generative adversarial nets. *arXiv preprint arXiv:14111784*
- Nathan Silberman PK, Derek Hoiem, Fergus R (2012) Indoor segmentation and support inference from RGBD images. In: Proceedings of the European Conference on Computer Vision
- Nazeri K, Ng E, Joseph T, Qureshi F, Ebrahimi M (2019) EdgeConnect: Structure guided image inpainting using edge prediction. In: Proceedings of the IEEE International Conference on Computer Vision Workshops, pp 0–0
- Pathak D, Krahenbuhl P, Donahue J, Darrell T, Efros AA (2016) Context encoders: feature learning by inpainting. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 2536–2544
- Pinheiro PO, Collobert R, Dollár P (2015) Learning to segment object candidates. In: Advances in Neural Information Processing Systems, pp 1990–1998
- Pinheiro PO, Lin TY, Collobert R, Dollár P (2016) Learning to refine object segments. In: Proceedings of the European Conference on Computer Vision, Springer, pp 75–91
- Qi L, Jiang L, Liu S, Shen X, Jia J (2019) Amodal instance segmentation with kins dataset. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 3014–3023
- Ren S, He K, Girshick R, Sun J (2015) Faster R-CNN: Towards real-time object detection with region proposal networks. In: Advances in Neural Information Processing Systems, pp 91–99
- Sengupta S, Jayaram V, Curless B, Seitz SM, Kemelmacher-Shlizerman I (2020) Background matting: The world is your green screen. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 2291–2300
- Shade J, Gortler S, He Lw, Szeliski R (1998) Layered depth images. In: Proceedings of the 25th annual conference on Computer graphics and interactive techniques, pp 231–242
- Silberman N, Hoiem D, Kohli P, Fergus R (2012) Indoor segmentation and support inference from rgb-d images. In: European conference on computer vision, Springer, pp 746–760
- Simonyan K, Zisserman A (2014) Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:14091556*
- Song S, Yu F, Zeng A, Chang AX, Savva M, Funkhouser T (2017) Semantic scene completion from a single depth image. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 1746–1754
- Sun D, Sudderth EB, Black MJ (2010) Layered image motion with explicit occlusions, temporal consistency, and depth ordering. In: Advances in Neural Information Processing Systems, pp 2226–2234
- Tighe J, Niethammer M, Lazebnik S (2014) Scene parsing with object instances and occlusion ordering. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 3748–3755
- Vahdat A, Kautz J (2020) NVAE: A deep hierarchical variational autoencoder. In: *Neural Information Processing Systems (NeurIPS)*
- Van Den Oord A, Vinyals O, et al. (2017) Neural discrete representation learning. In: *Advances in Neural Information Processing Systems*, pp 6306–6315

- Winn J, Shotton J (2006) The layout consistent random field for recognizing and segmenting partially occluded objects. In: 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06), IEEE, vol 1, pp 37–44
- Yan X, Wang F, Liu W, Yu Y, He S, Pan J (2019) Visualizing the invisible: Occluded vehicle segmentation and recovery. In: Proceedings of the IEEE International Conference on Computer Vision, pp 7618–7627
- Yang C, Lu X, Lin Z, Shechtman E, Wang O, Li H (2017) High-resolution image inpainting using multi-scale neural patch synthesis. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), vol 1, p 3
- Yang Y, Hallman S, Ramanan D, Fowlkes C (2010) Layered object detection for multi-class segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 3113–3120
- Yang Y, Hallman S, Ramanan D, Fowlkes CC (2011) Layered object models for image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 34(9):1731–1743
- Yu J, Lin Z, Yang J, Shen X, Lu X, Huang TS (2018) Generative image inpainting with contextual attention. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 5505–5514
- Zamir AR, Sax A, Shen W, Guibas LJ, Malik J, Savarese S (2018) Taskonomy: Disentangling task transfer learning. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 3712–3722
- Zhan X, Pan X, Dai B, Liu Z, Lin D, Loy CC (2020) Self-supervised scene de-occlusion. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 3784–3792
- Zhang Z, Schwing AG, Fidler S, Urtasun R (2015) Monocular object instance segmentation and depth ordering with cnns. In: Proceedings of the IEEE International Conference on Computer Vision, pp 2614–2622
- Zheng C, Cham TJ, Cai J (2019) Pluralistic image completion. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 1438–1447
- Zhu Y, Tian Y, Metaxas D, Dollár P (2017) Semantic amodal segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 1464–1472

A Experimental Details

Training. We trained our model on the synthetic data into three phases: 1) the layered scene decomposition network (Fig. 2(b)) is trained with loss L_{decomp} for 24 epochs, where at each layer, re-composited layered ground-truths are used as input. 2) Separately, the completion network (Fig. 2(c)) is trained with loss L_{comp} for 50 epochs, wherein the ground-truth layer orders and segmented masks are used to designate the *invisible* regions for completion. 3) Both decomposition and completion networks were trained jointly for 12 epochs, *without relying on ground-truths as input* at any layer (Fig. 2(a)). Doing so allows the scene decomposition network to *learn to cope with flaws* (e.g. texture artifacts) in the scene completion network, and vice versa. For each scene, the iteration ends when no more objects are detected, or a maximum 10 iterations is reached.

The training on real data only involved phases 1 and 3, as no ground-truth appearances are available for the invisible parts. The layered decomposition network is trained only for one layer (original image) in phase 1 due to *no* re-composed ground-truth images. Since phase 3 does not rely on ground-truths as input, we trained it layer-by-layer on real images by providing the “pseudo ground truth” appearances to calculate the reconstruction loss. To reduce the effect of progressively introduced artifacts in image completion, we used bounding boxes detected in the first layer as proposals for remaining decomposition steps.

Inference. During testing, fully visible instances were selected out and assigned an absolute layer order corresponding to the step index s_k . In each layer, the decomposition network selects the highest scoring 100 detected boxes for mask segmentation and non-occlusion predication. As we observed that higher object classification scores provided more accurately segmented boundaries, we only selected non-occluded objects with high object classification scores and non-occlusion scores (thresholds of 0.5 for synthetic images and 0.3 for real images) among these 100 candidates. We further observed that in some cases, we detected multiple objects with high object classification confidences, yet none were classified as fully visible due to low non-occlusion scores, especially in complex scenes with steps larger than 5. We will then choose the instance with the highest non-occlusion score so that *at least one object is selected at each layer*. When no objects are detected, the iteration stops.

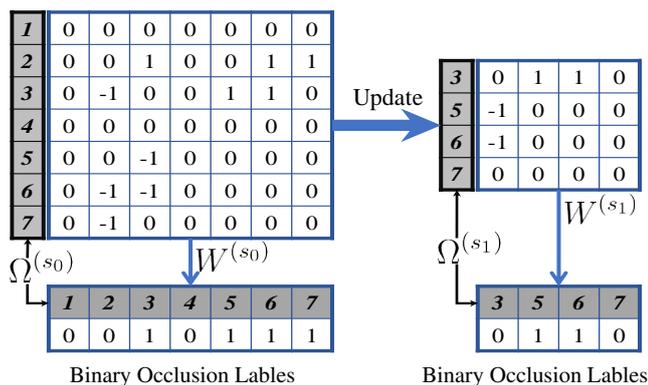


Fig. A.1 An illustration of obtaining the ground-truth binary occlusion labels from pairwise order graph $G = (\Omega, W)$ in each step s_k . If the indegree of a vertex is 0, it will be labeled as 0, a fully visible instance. Otherwise, the instance will be labeled as 1, being occluded. When some objects are detected and selected out in the previous step, the object indexes and the corresponding occlusions will be eliminated.

Instance depth ordering update. As illustrated in Fig. A.1, we calculate the indegree $deg^-(\omega)$ (counts of -1) of each instance in the matrix. If $deg^-(\omega) = 0$, meaning no objects are in front of it, its binary occlusion label will be 0. Otherwise, the object is occluded, labeled as 1. At each step, the fully visible objects will be eliminated from the directed graph G , and the ground-truth binary occlusion labels will be updated in each step. So if the table (instance #2) was selected in the previous step, the vertex index Ω will be updated after the corresponding object ω_2 is deleted from the occlusion matrix.

B Rendering Dataset

B.1 Data Rendering

Our **completed scene decomposition (CSD) dataset** was created using Maya (Autodesk Maya, 2019), based on the SUNCG CAD models (Song et al., 2017). The original SUNCG dataset contains 45,622 different houses with realistically modeled rooms. As realistically rendering needs a lot of time (average 1 hour for each house), we only selected 2,456 houses in current work. The set of camera views was based on the original OpenGL-rendering method in SUNCG, but further filtered such that a camera view was only be picked when at least 5 objects appeared in that view. We then realistically rendered RGB images for the selected views. Eight examples are shown in Fig. A.2 for various room types and lighting environments. Notice that our rendered images are much more realistic than the OpenGL rendered versions from the original SUNCG and likewise in (Dhamo et al., 2019).

To *visit the invisible*, the supervised method needs ground truth for the original occluded regions of each object. One possible way is to remove the fully visible objects in one layer and re-render the updated scene for the next layer, repeating this for all layers. However, during the training, the fully visible objects are not always correctly detected by the models. Thus, for more robust learning, we need to consider all different combinations of objects and the background. Given N objects, we would need to render 2^N images for each view. As we can see from the data statistics presented in Fig. B.4, an average of 11 objects are visible in each view. Due to slow rendering, we do not have the capacity to render all such scenes (average $2^{11} = 2048$ images per view). Instead, we separately rendered each isolated object with the full RGB appearance, as well as the empty room.

During training, the image of a scene is created by using a combination of the rendered images of these individual objects and the background to create a composed image, based on the remaining objects left after applying the scene decomposition network at each step. Since the room environment is empty for each individual objects during the rendering, the re-composited scenes have lower realism than the original scenes, due to missing shadows and lack of indirect illumination from other objects. In this project, we do not consider the challenges of working with shadows and indirect illumination, leaving those for future research.

B.2 Data Annotation

In Fig. B.3, we show one example of a rendered image with rich annotations, consisting of a semantic map, a depth map, visible annotations and full (amodal) annotations. For the semantic maps, we transferred the SUNCG class categories to NYUD-V2 40 categories so that this rendered dataset can be tested on real world images. The depth map is stored in 16-bit format, with the largest indoor depth value at 20m. The class category and layer order (both absolute layer order and pairwise occlusion order) are included for visible annotations and full annotations. The visible annotations also contain the visible bounding-box offset and visible binary mask for each instance. Additionally, we also



Fig. A.2 Realistic rendered images in the CSD dataset with various environment and lighting.

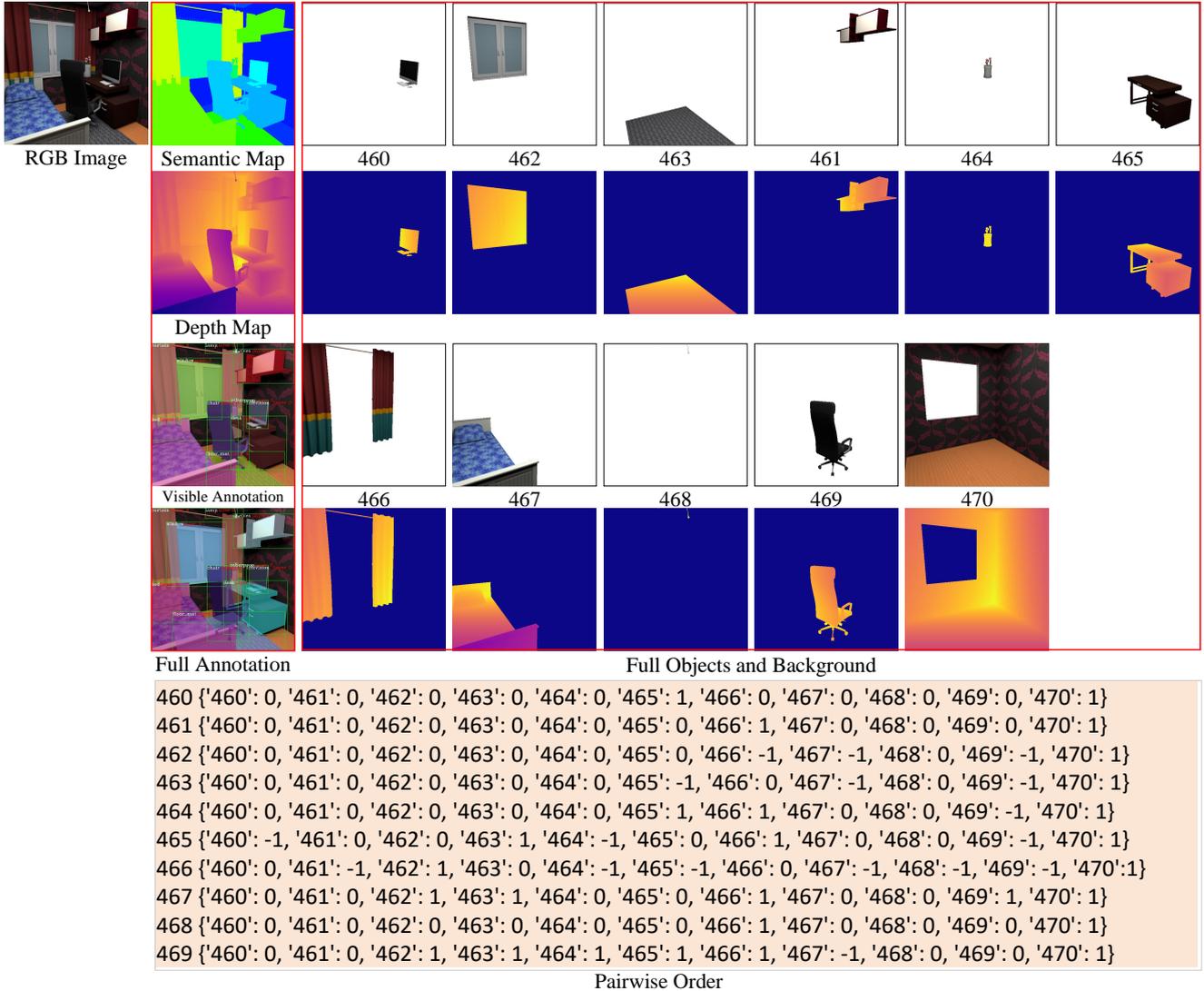


Fig. B.3 Illustration of Data Annotation. For each rendered image, we have a corresponding semantic map, a depth map, and dense annotation, including class category, bounding box, instance mask, absolute layer order and pairwise order. In addition, for each object, we have a full RGBA image and depth map.

have the full (amodal) bounding-box offset and completed mask for each individual object.

Pairwise Occlusion Order. The pairwise order for each object is a vector storing the occlusion relationship between itself and all other objects. We use three numbers $\{-1, 0, 1\}$ to encode the occlusion relationship between two objects — -1: occluded, 0: no relationship, 1: front (*i.e.* occluding). As can be seen in Fig. B.3, the computer (object number: #460) does not overlap the shelves (object number: #461), so the pairwise order is “0”, indicating these two objects have no occlusion relationship. The computer is however on top of the desk (object number: #465), hence the pairwise order for $W_{460,465}$ is “1”, and con-

versely the pairwise order for $W_{465,460}$ is “-1”, representing that the desk is occluded by the computer.

B.3 Data Statistics

In total, there are 11,434 views encompassing 129,336 labeled object instances in our rendered dataset. On average, there are 11 individual objects per view. Among these, 63.58% objects are partially occluded by other objects and the average occlusion ratio (average IoU between two objects) is 26.27%.

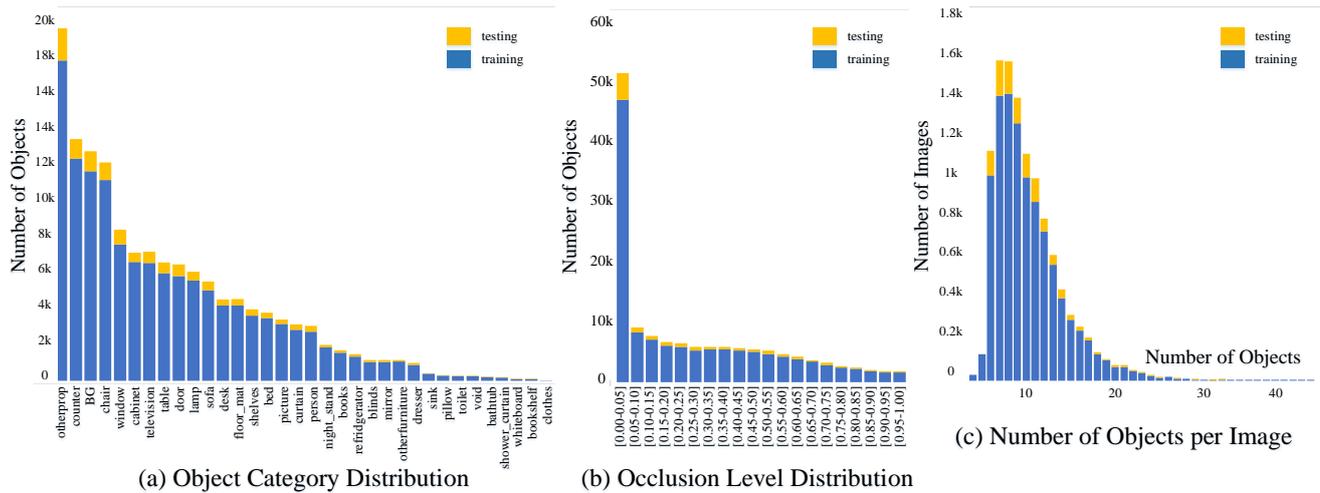


Fig. B.4 Data Statistics. Left: the object category distribution. Middle: the occlusion level distribution. Right: distribution of number of objects per image. On average there are 11 objects in each room.

Object Category Statistics. Fig. B.4(a) shows the overall object category distribution in our CSD dataset. Overall, the distribution is highly similar to the object distribution of NYUD-V2 dataset Nathan Silberman and Fergus (2012), containing a diverse set of common furniture and objects in indoor rooms. “Other props” and “Other furniture” are atypical objects that do not belong in a common category. In particular, “Other props” are small objects that can be easily removed, while “Other furniture” are large objects with more permanent locations. Additionally, we merge floors, ceilings, and walls as “BG” in this paper. If the user wants to obtain the separated semantic maps for these structures, these are also available.

Occlusion Statistics. The occlusion level is defined as the fraction of overlapping regions between two objects (Intersection over Union, or IOU). We divide the occlusion into 20 levels from highly visible (denoted as [0.00-0.05] fraction of occlusion) to highly invisible (denoted as [0.95-1.00] fraction of occlusion), with 0.05 increment in the fraction of occlusion for each level. Fig. B.4(b) shows the occlusion level in our dataset. In general, the distribution of occlusion levels is similar to the distribution in (Zhu et al., 2017), where a vast number of the instances are slightly occluded, while only a small number of instances are heavily occluded.

Object Count Distribution. Fig. B.4(c) shows the distribution of the number of objects present per view. On average, there are more than 11 objects in each view. This supports the learning of rich scene contextual information for a completed scene decomposition task, instead of processing each object in isolation.

B.4 Data Encoding

After we get the views and corresponding dense annotations, we encode the data annotation to COCO format³. The annotations are stored using JSON, and the CSD API will be made available for visualizing and utilizing the rendered dataset. The JSON file contains a series of fields, including “categories”, “images” and “annotations”. One short example is included in the supplementary file named `csd_short.json`.

³ <http://cocodataset.org>