

Learning Semantics-aware Distance Map with Semantics Layering Network for Amodal Instance Segmentation

Ziheng Zhang*
zhangzh@shanghaitech.edu.cn
ShanghaiTech University
Shanghai, China

Anpei Chen*
chenap@shanghaitech.edu.cn
ShanghaiTech University
Shanghai, China

Ling Xie
xieling@shanghaitech.edu.cn
ShanghaiTech University
Shanghai, China

Jingyi Yu
yujingyi@shanghaitech.edu.cn
ShanghaiTech University
Shanghai, China

Shenghua Gao[†]
gaoshh@shanghaitech.edu.cn
ShanghaiTech University
Shanghai, China

ABSTRACT

In this work, we demonstrate yet another approach to tackle the amodal segmentation problem. Specifically, we first introduce a new representation, namely a semantics-aware distance map (sem-dist map), to serve as our target for amodal segmentation instead of the commonly used masks and heatmaps. The sem-dist map is a kind of level-set representation, of which the different regions of an object are placed into different levels on the map according to their visibility. It is a natural extension of masks and heatmaps, where modal, amodal segmentation, as well as depth order information, are all well-described. Then we also introduce a novel convolutional neural network (CNN) architecture, which we refer to as semantic layering network, to estimate sem-dist maps layer by layer, from the global-level to the instance-level, for all objects in an image. Extensive experiments on the COCOA and D2SA datasets have demonstrated that our framework can predict amodal segmentation, occlusion, and depth order with state-of-the-art performance.

CCS CONCEPTS

• **Computing methodologies** → **Image segmentation; Scene understanding; Neural networks.**

KEYWORDS

amodal perception, image segmentation, convolutional neural networks

ACM Reference Format:

Ziheng Zhang, Anpei Chen, Ling Xie, Jingyi Yu, and Shenghua Gao. 2019. Learning Semantics-aware Distance Map with Semantics Layering Network

*indicates equal contributions.

[†]indicates the corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

MM '19, October 21–25, 2019, Nice, France

© 2019 Association for Computing Machinery.

ACM ISBN 978-1-4503-6889-6/19/10...\$15.00

<https://doi.org/10.1145/3343031.3350911>

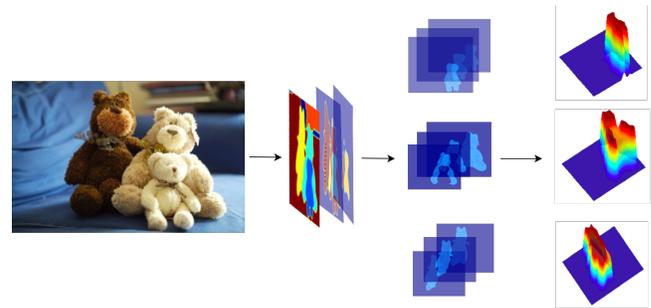


Figure 1: The pipeline of our framework. The pictures from left to right are: 1) an input image; 2) the global layering map, which describes the number of overlapped objects in each pixel of the input image; 3) the instance layering maps, which contain the pixel-wise depth layer indices for the corresponding object instance in the input image; 4) the sim-dist maps, which indicate the modal segmentation, the amodal segmentation and the relative depth order of all object instances in the input image. An amodal segmentation of an object is inferred layer by layer, from global layering maps to instance layering maps and the final sem-dist maps in SLN (Semantics Layering Network). (best viewed in color)

for Amodal Instance Segmentation. In *Proceedings of the 27th ACM International Conference on Multimedia (MM '19)*, October 21–25, 2019, Nice, France. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3343031.3350911>

1 INTRODUCTION

The recent years have witnessed great progress in visual understanding from image classification [16, 21, 33] and detection [23, 25, 30] to segmentation [2, 3, 26]. It appears that the performance of machine vision systems is stepping closer and closer to that of humans in terms of accuracy. Despite that, human vision has the strong ability to see beyond the visible, i.e., to perceive whole semantic concepts with only partial visibility. This ability, also known as amodal perception [41], had hardly been exploited to develop machine vision systems with similar capability until most recently, when researchers started making such attempts by modifying state-of-the-art segmentation models and training them with synthetic [10, 22, 41] and/or human-annotated datasets [10, 41].

Amodal perception implicitly requires vision systems to possess three critical abilities. The first is to recognize an object even if it is partially occluded by others. The second is to infer the most probable appearance of the invisible parts of an object given only the visible parts. And the third is to be aware of the relationship between overlapping objects. In other words, an amodal perception system should always keep a *consistent* whole picture of an object given arbitrary occlusion patterns; and it should understand the depth order between objects. Indeed, every single objective of building an amodal perception system is not new and has long been studied separately or jointly in the community. There is plenty of literature on visual understanding despite the presence of occlusion [11, 31, 35, 37], depth ordering [27, 36, 40], and object completion and inpainting [7, 20], that together show the feasibility and practicability of building machine vision systems with such capabilities. In this work, we try to solve all of these problems in a single amodal segmentation framework.

Amodal segmentation is not an easy task, as it needs a model to understand not only semantic concepts but also the relative depth order between them. Existing researches [10, 22, 41] have tended to directly use the formulation of traditional semantic segmentation to deal with amodal segmentation, where amodal masks, as well as visible/invisible masks, are regressed for amodal object proposals. Though depth order is retrievable by analyzing the amodal mask and the visible/invisible mask, it is not explicitly supervised in the training process; thus it is not clear whether the model has learned such information. In addition, even though predicting a segmentation mask for the invisible parts of an object is more difficult than for that for the visible parts, existing methods directly regress the full amodal mask without considering such differences in difficulties. In contrast, we believe that each part of an object with different levels of occlusion should be treated separately.

To this end, we propose the formulation of the amodal segmentation problem as the learning of a *semantics-aware distance map* (sem-dist map) for each object, which describes the pixel-level modal, amodal and relative depth order, as demonstrated in Fig. 2{a,b,d}. The sem-dist map is a kind of level-set representation, which describes not only the confidence of occurrence but also the global visibility level of the corresponding object in a scene. Take two overlapping objects as an example (see Fig. 2{a,b,d}). The values in their sem-dist maps should at least fall into three intervals: pixels that do not belong to each object’s amodal segmentation in the backmost interval, pixels that belong to the visible parts of each object in the uppermost interval and pixels that belong to the invisible parts of occluded objects in the middle interval. In this case, imagine that we have a reference level plane that can move up and down along the intensity dimension of the sem-dist maps. If we move the plane from the highest level to the lowest, the whole object in front and the visible parts of the occluded objects, *i.e.*, their modal segmentation, will first emerge out of the plane. Then as we move the plane down to the lower bound of the middle interval, we will have the amodal segmentation of both objects.

This kind of formulation has several advantages: first, it unifies amodal segmentation and modal segmentation in a single framework, and one can easily get them by just thresholding the sem-dist map; second, since the relative visibility level is equivalent to the relative depth order, the output of our framework naturally contains

pixel-wise depth information for each object; third, each part of an object is naturally layered on the sem-dist map according to their visibility level, which explicitly reflects the difference in difficulties in predicting their segmentation. In order to estimate a sem-dist map from an image, we introduce a novel convolutional neural network (CNN) architecture, which we refer to as the semantic layering network (SLN). The SLN is a proposal-based two-stage framework consisting of four modules: an encoder network (ENC), a global layering module (GLM), a region proposal network (RPN), and an instance layering module (ILM). An input image first goes into the ENC for feature extraction. Then the global layering maps, which describe the pixel-level visibility level for all objects, are inferred by the GLM. Finally, the global layering maps and the extracted image features are fed into both the RPN and the ILM to predict the instance layering maps and the final sem-dist maps. The pipeline of our amodal segmentation framework is shown in Fig. 1, and the architecture of our proposed SLN is demonstrated in Fig. 4.

In sum, the main contributions of this work are as follows: first, we introduce the sem-dist map for amodal instance segmentation, which unifies all targets of amodal segmentation into a single compact representation; second, we present a novel CNN architecture to solve the amodal segmentation problem; and third, we conduct extensive experiments which demonstrate that our framework can predict amodal segmentation and depth order with state-of-the-art performance. We have released our code, pre-trained models and results¹.

2 RELATED WORK

As previously mentioned, several topics, including occlusion handling, depth ordering, and object completion, are related to amodal segmentation. Hence we first briefly review the literature on these topics, then introduce the existing methods on amodal segmentation.

2.1 Visual Understanding with Occlusion

There has been plenty of researches on object detection [23, 25, 30], semantic segmentation [2, 3, 26], and instance segmentation [5, 14, 28], yet most of them have merely focused on what we can see in an image, *i.e.*, predicting bounding boxes around or pixels belonging to the visible parts of objects. As such, occlusion is often treated as noise and is overcome implicitly through learning-based methods that train on huge amounts of data. There are also other researchers trying to explicitly eliminate the effect of occlusion. For instance, [37] used conditional random fields (CRFs) to describe the possible configurations of object parts, which imposes constraints on the appearance of partially occluded objects given their visible parts. [11] introduced binary cells for object detection, thereby explicitly indicating the visibility of object parts, thus making detection more robust to occlusion. [12] incorporated global scene priors and occluding object classes as cues to complete the labeling of partially occluded semantic regions. [18] proposed the modeling of occlusions by inferring the 3D interactions of objects for instance detection from an arbitrary viewpoint. And finally, [4] used class-specific likelihood maps and inferred visible and occluded regions to

¹<https://github.com/apchenstu/SLN-Amodal.git>

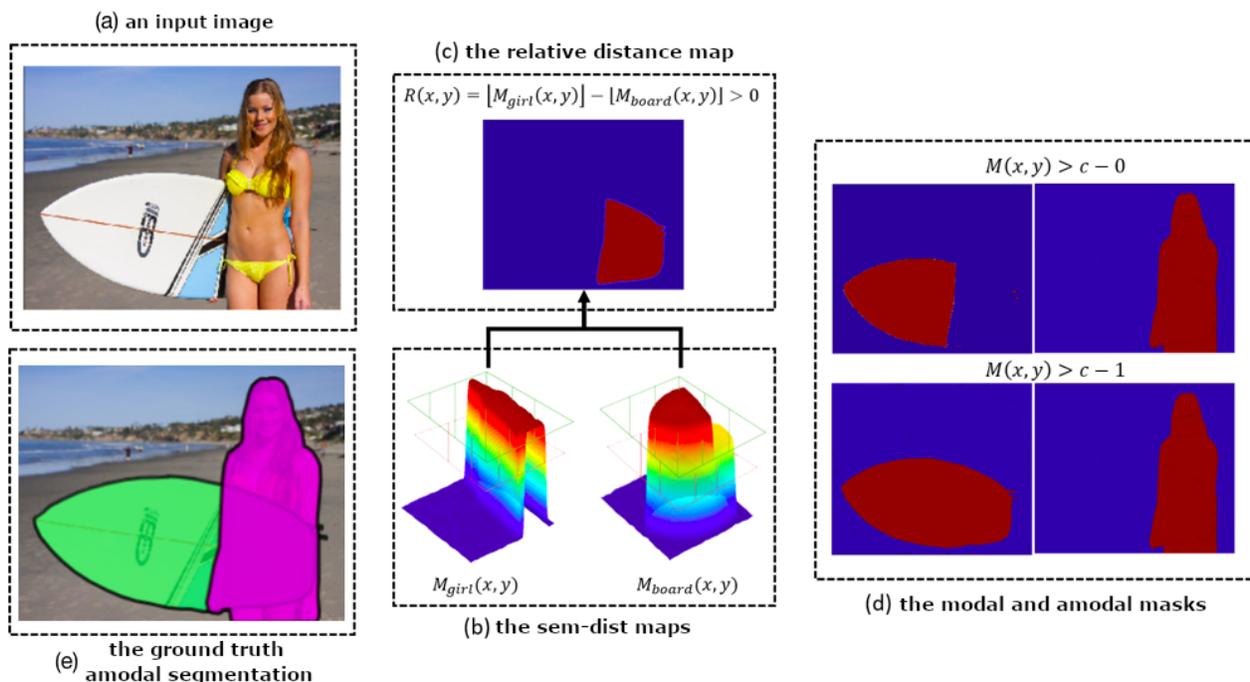


Figure 2: Illustrations of: a) a example image of a girl holding a surfboard; b) the sem-dist maps for the girl $M_{girl}(x,y)$ and the surfboard M_{board} ; c) the relative distance map $R(x,y)$ of the girl and the surfboard; d) the masks of the girl and the surfboard generated by moving the reference level plan from the highest visibility level to the lowest as shown in (b); e) the ground truth amodal segmentation of the girl and the surfboard. (best viewed in color)

obtain segmentation candidates via energy minimizing frameworks, which were then scored using a class-specific classifier.

2.2 Depth Ordering

Depth ordering has long been studied. Early research on object tracking [1, 8, 34, 38] incorporated shapes, boundaries, and occlusion patterns across frames as clues to predict depth order. In [19], depth order was inferred using features extracted from boundaries and junctions separately, on which a Markov Random Field (MRF) model and graph optimization were applied to get globally consistent ordering. [27] first identified occlusions between objects with T-junctions and highly convex contours, then used these occlusion cues to arrange objects according to depth order. [40] built a CNN-based architecture to jointly reason pixel-wise instance-level segmentation as well as depth order from multi-scale image patches and they combined predictions into the final labeling via the MRF. [36] used a fully convolutional network (FCN) to jointly predict pixel-level semantic labels, depths and the directions to object centers from a single street scene image.

2.3 Object Completion and Inpainting

Some of the research that was capable of handling occlusions [4, 11, 12, 37] also applied object completion to occluded objects. Those methods, however, have been shown to work only on specific or limited object categories and relied on available shape models or depth inputs. [20] took a step forward by utilizing a probabilistic

framework to learn category-specific object size distributions and then leveraged the model to estimate the veridical size of the occluded objects in new images. Also, a recent piece of research [7] made attempts to generate the appearance of the invisible part of an object through a generative adversarial network (GAN). For depth inpainting, research [17] used color images to guide the inpainting of the corresponding depth maps by aligning the edges. Xue *et al.* [39], on the other hand, used low-rank regularization to inpaint single depth images without using color images.

2.4 Amodal Segmentation

The concept of amodal segmentation has just emerged in the last few years [10, 22, 41], though similar problems had actually been addressed years before in many applications including detection [11, 12, 20], segmentation [4, 37], reconstruction [13, 32], and so on. Traditional approaches, however have usually relied on depth information or focused on specific object categories, while recent methods have solely used RGB images only and targeted objects of arbitrary categories. [22] presented the first method to tackle the amodal segmentation problem, where authors first generated data with amodal ground truths by randomly overlaying one object instance onto another and then using the generated data to train a network that could predict the segmentation heatmap. An amodal bounding box was then generated using their proposed Iterative Bounding Box Expansion strategy from the segmentation heatmap and the modal bounding box predicted by a general-purpose object detector. [41] built a large amodal segmentation dataset with

human-annotated masks of amodal, visible and invisible regions of each object together with relative depth orders. Multiple baseline models were designed and trained on the proposed dataset to predict amodal masks or the depth orders. [10] proposed a multi-task model that simultaneously predicted amodal masks, visible masks, and occlusion masks for each object instance. In addition, they also provided a new semantic amodal segmentation dataset D2S amodal and supplemented the class labels for the amodal datasets proposed in [41].

3 SEMANTICS-AWARE DISTANCE MAP

In this section, we will introduce the semantics-aware distance map (the sem-dist map) for amodal instance segmentation (the example sem-dist maps for a girl and a surfboard are shown in Fig. 2b). Unlike the commonly used heatmap in instance segmentation, where each pixel measures the confidence of the occurrence of the visible part of an object ranging from zero to one, the sem-dist map describes the pixel-level *visibility* of the *whole* object where each pixel value ranges from zero to positive infinity.

More concretely, we first define the visibility level L of a region $\Omega = \{(x, y)\}$ belonging to an object to be an integer l if and only if the region will become visible after we ‘remove’ at least l objects from the scene. Particularly, if a region does not belong to the object at any visibility level, we define its visibility level to be zero. Then, the sem-dist map M for an object instance can be defined as:

$$M(x, y) = C(x, y) - L(x, y) \quad (1)$$

where $L(x, y) \in \mathbb{N}^0$ and $C(x, y) \in [0, 1]$ stand for the visibility level and the confidence of occurrence of the object’s amodal segmentation at (x, y) on the image, respectively. The intuition behind this definition is that an amodal segmentation framework should have higher confidence in predicting the occurrence of an object part when it is visible in an image, and have lower confidence when it is occluded by more other objects, *i.e.*, when the visibility level of the object part is higher. The sem-dist map emphasizes such varying difficulties in predicting segmentation for an object by explicitly adding bias to the heatmap of an object part according to its visibility level.

The sem-dist map contains rich object information. First, it contains both modal and amodal segmentations of an object. To see that, if one only considers the region in the sem-dist map where the visibility level equals zero, then M is reduced to the modal heatmap of the object M_{modal} . In addition, the amodal heatmap M_{amodal} can also be easily calculated by taking the fractional part of each pixel in M , *i.e.*

$$M_{modal}(x, y) = \begin{cases} M(x, y), & M(x, y) \in [0, 1) \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

$$M_{amodal}(x, y) = M(x, y) - \lfloor M(x, y) \rfloor \quad (3)$$

where $\lfloor \cdot \rfloor$ stands for the floor function.

On the other hand, for the same region in the sem-dist maps of two objects where one occludes another, the one with a lower visibility level indicates that the corresponding part of the object is closer to the camera than another. Thus the pixel-level depth order of two objects can also be retrieved from the sem-dist map. Considering the sem-dist maps of two mutually intersecting objects

M_A and M_B , if we take the difference of the integer part of the two corresponding sem-dist maps, we will have another map which we denote as R_{AB} , *i.e.*

$$R_{AB}(x, y) = \begin{cases} \lfloor M_A(x, y) \rfloor - \lfloor M_B(x, y) \rfloor, & (x, y) \in \Omega_{AB} \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

where Ω_{AB} stands for the region within which object A and B are mutually intersecting, *i.e.*

$$\Omega_{AB} = \{(x, y) | (M_A(x, y) - \lfloor M_A(x, y) \rfloor) \cdot (M_B(x, y) - \lfloor M_B(x, y) \rfloor) > c^2\} \quad (5)$$

where c is the confidence threshold. Note that we use c^2 here because that $(M_A(x, y) - \lfloor M_A(x, y) \rfloor) \cdot (M_B(x, y) - \lfloor M_B(x, y) \rfloor)$ is homogeneous to the square of the confidence $C(x, y)$, according to Eq. 1.

We can easily get the pixel-wise depth order between A and B by examining the sign of each pixel in R_{AB} as shown in Fig. 2{b,c}: if $R_{AB}(x, y) > 0$ then A is closer to the camera than B and if $R_{AB}(x, y) < 0$ then B is closer to the camera than A at (x, y) . Taking the $R(x, y)$ in Fig. 2c as an example, we can conclude that the girl is in front of the surfboard where they are occluded because $R(x, y) > 0$ in that region. Note that the depth order between A and B is undefined where $R_{AB}(x, y) = 0$. In addition, one should also note that the depth order is defined on a regional basis in this work, rather than on an object basis as that in the amodal dataset [41]. The reason is that each part of an object can have a different depth order with respect to other objects. Considering a figure on horseback, if a picture is taken from the side of the horse, then the two legs of the rider should have different depth orders with respect to the horse (see Fig. 3). Thus the depth order of an object is not well-defined, and extra criteria are needed in order to get the depth order of the full object.

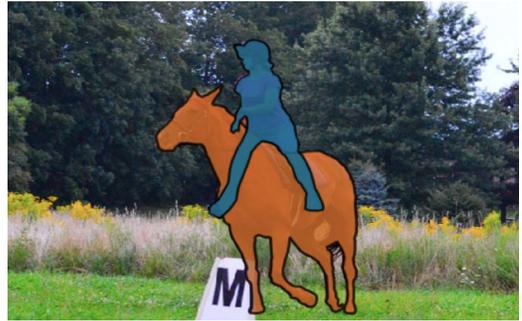


Figure 3: A figure on horseback. The pairwise depth order between the rider and the horse is unclear when two legs of the rider have different depth orders with respect to the horse. (best viewed in color)

From the discussions above, we can see that in addition to semantic information, the sem-dist map is also capable of describing the pixel-wise depth order between two objects. Given such properties of this kind of representation, we call it semantics-aware distance map (sem-dist map). Though we use the term distance to emphasize the fact that we can get the relative distance from the camera to mutually intersecting objects with a sem-dist map, one should note

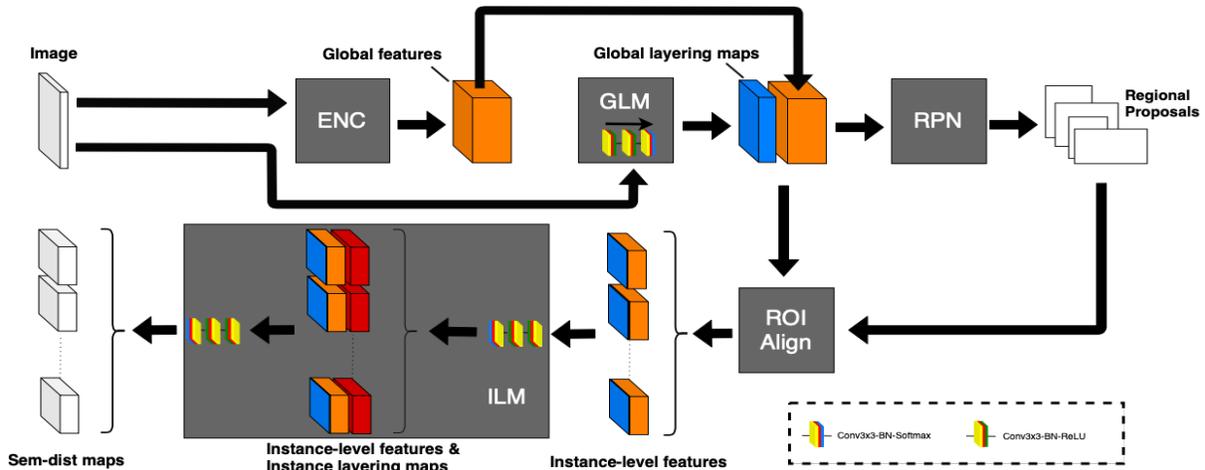


Figure 4: The architecture of semantics layering network. The input image first goes into the ENC for feature extraction. At the same time, GLM also takes the input image to predict the global layering maps. After that, both the feature maps from ENC and the global layer maps are fed into RPN to generate amodal bounding boxes proposals, which are used by ROI-Align layer [14] to extract both features and global layer maps for each instance. Finally, the outputs from ROI-Align layer are collected by the ILM to estimate the final instance-level sem-dist maps. (best viewed in color)

that the absolute distance cannot be retrieved from it. In addition, if an object has self-occluded parts on its amodal segmentation, relative distances will be determined by the occluding parts that are closer to the camera, according to our definition.

The sem-dist map is a compact representation, where the modal and amodal segmentation, occlusion and depth order of objects are well-described. Its completeness and compactness make it a desirable target for amodal perception, yet also make it difficult to train with existing methods. Therefore we introduce a carefully designed convolutional neural network (CNN) architecture to infer the sem-dist map from an image.

4 SEMANTICS LAYERING NETWORK

Here we introduce the semantics layering network (SLN), a two-stage CNN-based architecture which infers instance-level sem-dist maps from a single RGB image.

4.1 Overview

SLN is a member of proposal-based two-stage architectures like Mask R-CNN [14], and it is composed of four modules: an encoder network (ENC), a region proposal network (RPN), a global layering module (GLM) and an instance layering module (ILM), as shown in Fig. 4. An input image first goes into the ENC for feature extraction. At the same time, GLM also takes the input images to predict a global layering maps. After that, both the feature maps from ENC and the global layer maps are fed into RPN to generate amodal bounding box proposals, which are used by an ROI-Align layer [14] to extract both features and global layer maps for each instance. Finally, the outputs from the ROI-Align layer are collected by the ILM to estimate final instance-level sem-dist maps. Since the architecture of the RPN is directly borrowed from the Faster R-CNN [30], we only detail the remaining three modules in the following subsections.

4.2 Encoder Network

Just like most CNN-based architectures, we utilize a general-purpose encoder network to extract features from input images before task-specific modules. The ENC can be an arbitrary network which has both a large enough receptive field to capture global information and a powerful enough capability to extract low-level and high-level features. Since the requirements of the encoder for our task meet that for semantic segmentation, we choose a dilated ResNet-50 [16] as our ENC architecture. In addition, we also collect features of size $256, 128 \times 128, 64 \times 64$ and 32×32 from four stages of ResNet, rescale them to the same height and width, then concatenate them all together to provide the final feature maps used by the following task-specific modules.

4.3 Global Layering Module

In a proposal-based two-stage architecture, task-specific heads take cropped features after an ROI-Pooling [30] or an ROI-Align [14] layer to make instance-level inferences. Though such a scheme is suitable for object detection or modal segmentation, etc., where the cropped features basically contain all the required information, it is not suitable in our case. The reason is that the sem-dist map describes the *global* visibility level for an object, yet the instance-level feature maps mostly contain information on single objects. Therefore, in addition to using an encoder with a larger receptive field, we introduce a global layering module (GLM) to explicitly encode the global visibility level in a global layering map.

The global layering map M_G is a multi-layer heatmap, of which each layer describes the pixel-level confidence of the occurrence of an object at the corresponding visibility level determined by the index of the layer. In other words, in M_G , the amodal heatmap of all objects are ‘placed’ in different layers according to the visibility level. M_G will retain the information of the global visibility level for all of the objects after passing through the ROI-Align layer because

	all regions						things only						stuff only					
	AP	AR ¹⁰	AR ¹⁰⁰	AR ^N	AR ^P	AR ^H	AP	AR ¹⁰	AR ¹⁰⁰	AR ^N	AR ^P	AR ^H	AP	AR ¹⁰	AR ¹⁰⁰	AR ^N	AR ^P	AR ^H
AmodalMask [41]	5.74	13.5	29.23	34.4	31.0	21.3	6.12	16.5	33.1	36.2	37.0	23.6	0.78	5.4	18.1	22.3	16.1	18.0
ARCNN [10]	4.1	10.2	21.3	27.2	22.0	13.3	4.4	12.0	23.9	28.3	34.7	15.2	0.3	4.8	13.8	19.8	15.1	10.1
ARCNN with visible mask	6.6	15.3	32.4	42.5	34.8	17.1	7.8	19.5	37.6	45.5	40.8	19.9	0.5	3.3	17.1	22.9	19.9	12.5
SLN (ours)	8.4	16.6	36.5	44.8	40.1	22.5	9.6	20.5	40.5	47.2	43.6	24.9	0.8	5.3	25.0	28.8	31.3	18.6

Table 1: Amodal segmentation results on the COCOA validation set for our method and two state-of-the-art methods under no, partial, and heavy occlusion (AR^N, AR^P, AR^H) and for different object types (things and stuff).

such information is encoded in the channel dimension. In addition, M_G can also be viewed as a pixel-level amodal proposal, which could be used to enhance the ability of RPN to propose amodal bounding boxes. Therefore M_G is a desirable intermediate target for our task.

As for the architecture of GLM, we adopt a simple *conv-bn-relu-conv-bn-relu-conv-softmax* block, which stands for three convolution layers followed by Batch Normalization and ReLU activation, although the last convolution layer uses a softmax activation function. The kernel size and stride size is 3×3 and 1, respectively, for all convolution layers.

4.4 Instance Layering Module

The ILM aims to infer the sem-dist map for each object instance from the feature maps after the ROI-Align layer. Though the sem-dist map can be regressed directly, we decide to first predict the instance layering map as an intermediate representation, and then estimate the final sem-dist map from it. The definition of the instance layering map is similar to the global instance map, except that the former only predicts the layered amodal heatmap for a single object instance rather than for all objects. We used the same architecture as GLM to regress the intermediate instance layering map. After that, the sem-dist map is first derived from the instance layering map with Eq. 1, then we use a *conv-bn-relu-conv-bn-relu-conv* block which takes the derived sem-dist map and features from the ROI-Align layer to predict a refined sem-dist map. The kernel size and stride size is 3×3 and 1 respectively, for all convolution layers.

4.5 Loss Function

There are three intermediate targets, *i.e.*, the region proposals, the global layering map as well as the instance layering map, and one final target, *i.e.*, the sem-dist map, that need to be supervised. For the region proposals, we use the smooth L1 loss and denote this loss term as ℓ_R . For the global layering map, the instance layering map and the sem-dist map, we use the binary cross entropy loss and denote the three loss terms as ℓ_G , ℓ_I and ℓ_M respectively. All of the targets can be jointly optimized by minimizing the overall loss ℓ , which is the weighted sum of all of the loss terms, *i.e.*

$$\ell = \lambda_R \ell_R + \lambda_G \ell_G + \lambda_I \ell_I + \lambda_M \ell_M \quad (6)$$

where $\lambda_{\{R,G,I,M\}}$ stands for the weight of the corresponding loss term.

4.6 Implementation Details

We implement SLN with the PyTorch² framework and trained it on a single Geforce 1080 Ti GPU. The ENC is initialized with parameters pretrained for semantic instance segmentation task on the COCO2014 dataset [24], and other modules are initialized with values according to the strategy described in [15]. All the targets are optimized using Stochastic Gradient Descent (SGD), with $lr = 0.2$, $weight_decay = 5 \times 10^{-4}$, $momentum = 0.9$, except for all normalization layers, of which $weight_decay$ is set to zero. Although all the modules in our network could be jointly optimized with Eq. 6, in practice we optimize GLM, RPN and ILM one after another and fix the parameter of the formerly trained module before training the later one for faster convergence, then we fine-tune all of the modules by jointly minimizing Eq. 6 with $\lambda_{\{R,G,I,M\}} = 1$.

5 EXPERIMENTS

We conduct several experiments to evaluate the performance of our framework on the COCO amodal dataset [41] and D2S amodal dataset [10], which we refer to as COCOA and D2SA, respectively. The COCOA dataset is the first amodal dataset consisting of 5000 images, of which 2500, 1250 and 1250 images are used for training, validation and testing, respectively. The annotations of the COCOA dataset include the amodal segmentation of each object, visible/invisible regions of each amodal segmentation as well as the relative depth orders of all objects in each image. All of the objects in the COCOA dataset are also classified into two categories: ‘things’ and ‘stuff’, where a ‘thing’ is an object with a canonical shape while ‘stuff’ has a consistent visual appearance but can be of an arbitrary extent. The D2SA dataset is a recently proposed class-specific amodal dataset, which contains 2000 and 3600 original and artificially augmented amodally annotated images from the D2S dataset [9] for training and validation. The D2SA dataset has class labels for each object and all kinds of annotations in the COCOA dataset except for depth order.

5.1 Evaluation Metrics

We evaluate the performance of our framework for amodal segmentation as well as depth order prediction, and the performance of amodal segmentation is reported for both the COCOA and D2SA datasets while the performance of depth order predictions is only reported for the COCOA dataset, due to a lacking of depth order ground truth data in the D2SA dataset.

Metrics for amodal segmentation: We use the mean average precision (AP) and mean average recall (AR) as our metrics as

²<https://pytorch.org/>

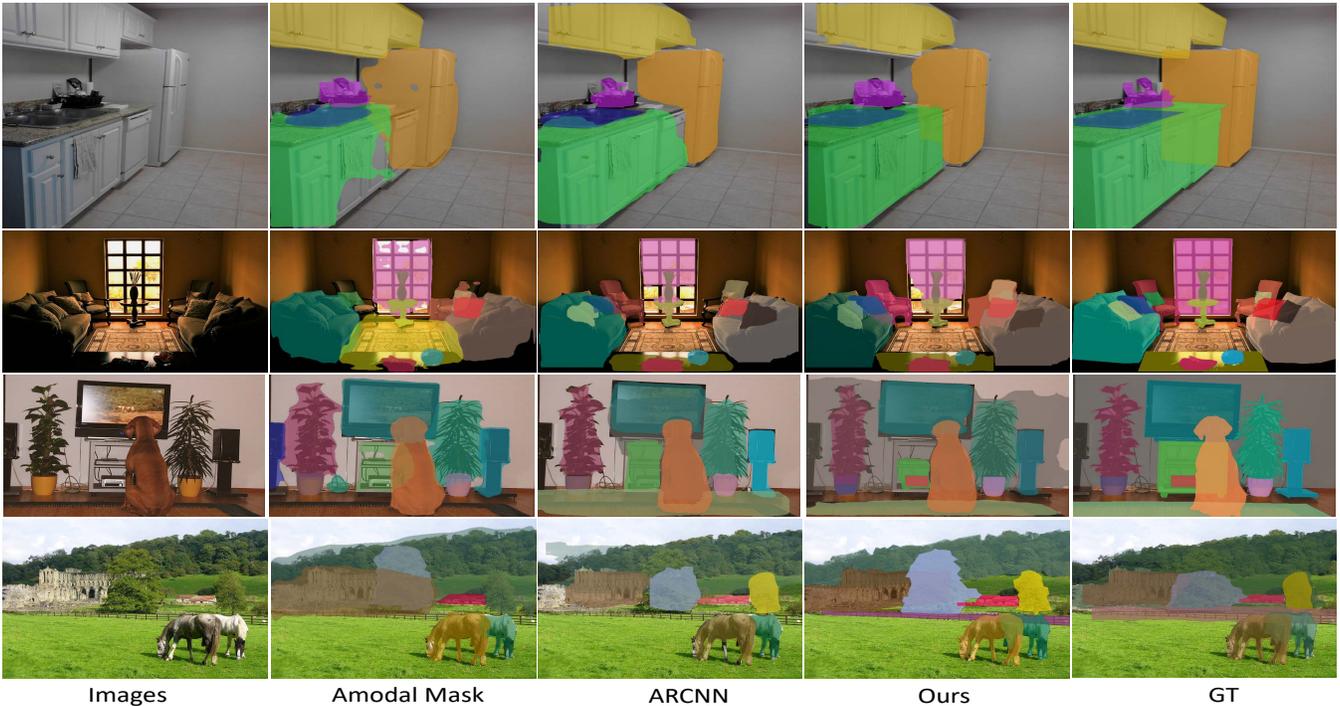


Figure 5: Some amodal segmentation results on the COCOA dataset. The columns from the left to the right are input images, results for AmodalMask [6], ARCNN [10], and our SLN, as well as the ground truth amodal segmentation, respectively (best viewed in color).

done in previous research. The AP is calculated by averaging the precision of mask prediction over ten equally spaced intersection-over-union (IOU) thresholds from 0.5 to 0.95, as is common practice. The AR is computed by averaging the segmentation recall over the same set of thresholds. The AR for both 10 and 100 segments per image is reported, which we denote as AR^{10} and AR^{100} , respectively. For the COCOA dataset, in addition to the overall AP and AR, we also report the AR for 100 segmentations of things and stuff per image separately as well as for three occlusion levels (none, partial or heavy occlusion), which we denote as AR^N , AR^R and AR^H , respectively. For the D2SA dataset, the AP, AR^{10} and AR^{100} for all objects and occluded objects are reported.

Metrics for pairwise depth order: For the COCOA dataset, we evaluate our depth order predictions. We report the accuracy in predicting which of two overlapping objects is in front of the other. There are 36k/23k overlapping objects in the train/val set.

5.2 Implementation Details

Amodal segmentation: We compare SLN with SOTA segmentation models, including AmodalMask in [41] and ARCNN in [10]. In all of our experiments, we use the released evaluation code³ used in [41] to compute the AP and the AR. For AmodalMask, we directly use the released segmentation results to conduct our evaluation. For ARCNN, because neither the model, the segmentation results nor the evaluation code have been released at the time of

writing this paper, we re-implement ARCNN according to the descriptions in [10] and evaluate the segmentation result using the same evaluation code in [41]. Although ARCNN is actually the standard Mask R-CNN [14] with a ResNet-101 [16] backbone trained on the COCOA/D2SA dataset, we can not reproduce similar AP and AR results as reported in [10]. Therefore, we report only the results of our re-implemented ARCNN and refer readers to their original paper [10] for their reported results. In addition, we also add another version of ARCNN where both amodal and visible masks rather than amodal masks only are regressed.

Depth ordering: We compare SLN with multiple baseline methods proposed in [41], including SharpMask [29] and ExpandMask [41], for depth ordering. In addition, the pairwise depth order predicted with the ground truth amodal mask and the sem-dist map are also evaluated. Pairwise depth orders are extracted from amodal masks or sem-dist maps in different ways. For the amodal mask, the best matching mask is first selected for each object (with largest IOU). Then the selected masks for each pair of overlapping objects are fed into the OrderNet [41], which is a pre-trained Resnet-50 model slightly modified for varying number of input channels, to predict the pairwise depth order. The OrderNet is trained and tested separately for each set of masks. For the sem-dist map, we first extract the amodal masks and find the best matching sem-dist maps for each object. Then we use both $R(x, y)$ in Eq. 4 and OrderNet to predict the pairwise depth orders from their sem-dist map pairs. Recall that the depth order derived from $R(x, y)$ is defined on a regional basis rather than an object basis as in the COCOA dataset.

³<https://github.com/Wakeupbuddy/amodalAPI>



Figure 6: Some amodal segmentation results for SLN on D2SA dataset.

After we get a region-level relative depth order with $R(x, y)$, an object-level relative depth order is chosen to be consistent with the relative depth order of the largest region.

5.3 Results

COCOA: The amodal segmentation results for AmodalMask [41], ARCNN [10], and our proposed SLN are listed in Table 1. In this experiment, our framework generally shows to achieve better recall than all of the existing methods. Interestingly, we found that the AP for all methods that rely on bounding box-based amodal proposals is lower than the one for methods that rely on mask-based amodal proposals. Besides, the former set of methods seems to work better for zero occluded and partially occluded things, while the later works better for heavily occluded stuff. We suspect that pixel-level proposals could be more easy to be generalized to the amodal case, especially for stuff, which have consistent appearance but can be of arbitrary extent. The fact that our framework achieves superior performance to ARCNN [10] demonstrates the effectiveness of the sem-dist map to guide the network to perceive amodal concepts. We also show some typical qualitative results for AmodalMask, ARCNN and our method on COCOA dataset in Fig. 5. From the first row, one can observe that our SLN segments the refrigerator and the cupboard better and the washbowl is also well recognized. For the second row, the SLN not only accurately segments the sofa, but also identified the bolster on it. And for the forth row, our SLN is the only method that correctly completes the occluded part of the horse.

As for depth ordering, the comparison of accuracy between SLN with multiple baselines are listed in Table 2. Because sem-dist maps contain rich depth information, the OrderNet predicts depth order

	Sharp Mask [29]	Expand Mask [41]	Amodal Mask [41]	SLN Map	Ground TruthM	Ground TruthS
OrderNet [41]	.786	.785	.791	.854	.817	.872
$R(x, y)$	-	-	-	.764	-	1.00

Table 2: Accuracy of pairwise depth order predicted by multiple baselines proposed in [41] and our method. Note that the inputs of the OrderNet are masks for all mask-based methods and sem-dist maps for our method. The GroundTruthM stands for the ground truth masks and the GroundTruthS stands for the ground truth sem-dist maps.

	all			partial occlusion			heavy occlusion		
	AP	AR ¹⁰	AR ¹⁰⁰	AP	AR ¹⁰	AR ¹⁰⁰	AP	AR ¹⁰	AR ¹⁰⁰
ARCNN[10]	23.4	47.3	73.1	6.2	38.0	72.6	0.6	17.8	45.0
SLN (ours)	25.3	41.3	78.6	6.5	42.4	77.5	0.6	22.7	57.2

Table 3: The quantitative results for ARCNN and our method on D2SA validation set for all/occluded objects.

better with the predicted sem-dist maps than masks. Besides, $R(x, y)$ is also able to give us reasonable depth order predictions.

D2SA: The results on D2SA dataset are listed in Table 3. We can observe that our framework achieves better overall results than ARCNN for both AP and AR. Examples of amodal mask predictions for ARCNN and our method are shown in Fig. 6. As one can see, our method is able to retrieve the occluded parts of objects and to predict reasonable amodal segmentation. Interestingly, SLN works better for heavy occlusions on D2SA dataset than on COCOA dataset. This is because there are only things in the D2SA dataset and the experiments on both datasets indicate that SLN (and other SOTA methods that rely on bounding box proposals) works pretty well for "things" rather than "stuff".

6 CONCLUSION

In this work, we proposed to tackle the amodal segmentation problem by learning a semantics-aware distance map (sem-dist map) for each object in an image. Compared with the commonly used mask representation, the semantics-aware distance map describes pixel-level visibility level of an object, from which the modal, amodal segmentation and relative depth order of the object can be derived elegantly. In order to estimate the sem-dist map, we introduced the semantics layering network (SLN), in which sem-dist maps for all objects are inferred layer by layer, from global-level to instance-level, from an image. Extensive experiments on COCOA and D2SA datasets have demonstrated that our framework can predict amodal segmentation and pairwise depth order with state-of-the-art performance.

In our experiments, we observed that the performance bottleneck of SLN is the low quality of object proposals from the bounding box-based proposal modules (RPN) It seems that RPN is not good at making large area amodal proposals. To tackle this problem, one may try to use pixel-level proposal modules or proposal-free frameworks to further boost the performance of SLN for amodal segmentation.

REFERENCES

- [1] L. Bergen and F. Meyer. 2000. A novel approach to depth ordering in monocular image sequences. In *Proceedings IEEE Conference on Computer Vision and Pattern Recognition. CVPR 2000 (Cat. No. PR00662)*, Vol. 2. 536–541 vol.2. <https://doi.org/10.1109/CVPR.2000.854907>
- [2] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L. Yuille. 2014. Semantic image segmentation with deep convolutional nets and fully connected crfs. *arXiv preprint arXiv:1412.7062* (2014).
- [3] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L. Yuille. 2018. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence* 40, 4 (2018), 834–848.
- [4] Yi-Ting Chen, Xiaokai Liu, and Ming-Hsuan Yang. 2015. Multi-Instance Object Segmentation With Occlusion Handling. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [5] Jifeng Dai, Kaiming He, and Jian Sun. 2016. Instance-aware semantic segmentation via multi-task network cascades. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 3150–3158.
- [6] Zhuo Deng and Longin Jan Latecki. 2017. Amodal detection of 3d objects: Inferring 3d bounding boxes from 2d ones in rgb-depth images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 5762–5770.
- [7] Kiana Ehsani, Roozbeh Mottaghi, and Ali Farhadi. 2018. Segan: Segmenting and generating the invisible. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 6144–6153.
- [8] D. Feldman and D. Weinshall. 2008. Motion Segmentation and Depth Ordering Using an Occlusion Detector. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 30, 7 (July 2008), 1171–1185. <https://doi.org/10.1109/TPAMI.2007.70766>
- [9] Patrick Follmann, Tobias Bottger, Philipp Hartinger, Rebecca König, and Markus Ulrich. 2018. MVTeC D2S: Densely Segmented Supermarket Dataset. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 569–585.
- [10] Patrick Follmann, Rebecca Kö Nig, Philipp Hä Rtinger, Michael Klostermann, and Tobias Bö Ttger. 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, 1328–1336.
- [11] Tianshi Gao, Benjamin Packer, and Daphne Koller. 2011. A segmentation-aware object detection model with occlusion handling. In *CVPR 2011*. IEEE, 1361–1368.
- [12] Ruiqi Guo and Derek Hoiem. 2012. Beyond the line of sight: labeling the underlying surfaces. In *European Conference on Computer Vision*. Springer, 761–774.
- [13] Saurabh Gupta, Pablo Arbelaez, and Jitendra Malik. 2013. Perceptual organization and recognition of indoor scenes from RGB-D images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 564–571.
- [14] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. 2017. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*. 2961–2969.
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification. In *Computer Vision (ICCV), 2015 IEEE International Conference on*. IEEE, 1026–1034.
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [17] Daniel Herrera C., Juho Kannala, Laazubor Ladicka, and Janne Heikkila. 2013. Depth Map Inpainting under a Second-Order Smoothness Prior. Vol. 7944. 555–566. https://doi.org/10.1007/978-3-642-38886-6_52
- [18] Edward Hsiao and Martial Hebert. 2014. Occlusion reasoning for object detection under arbitrary viewpoint. *IEEE transactions on pattern analysis and machine intelligence* 36, 9 (2014), 1803–1815.
- [19] Z. Jia, A. Gallagher, Y. Chang, and T. Chen. 2012. A learning-based framework for depth ordering. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*. 294–301. <https://doi.org/10.1109/CVPR.2012.6247688>
- [20] Abhishek Kar, Shubham Tulsiani, Joao Carreira, and Jitendra Malik. 2015. Amodal completion and size constancy in natural scenes. In *Proceedings of the IEEE International Conference on Computer Vision*. 127–135.
- [21] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*. 1097–1105.
- [22] Ke Li and Jitendra Malik. 2016. Amodal instance segmentation. In *European Conference on Computer Vision*. Springer, 677–693.
- [23] Tsung-Yi Lin, Piotr Dollar, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. 2017. Feature pyramid networks for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2117–2125.
- [24] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollar, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *European conference on computer vision*. Springer, 740–755.
- [25] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. 2016. Ssd: Single shot multibox detector. In *European conference on computer vision*. Springer, 21–37.
- [26] Jonathan Long, Evan Shelhamer, and Trevor Darrell. 2015. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3431–3440.
- [27] G. Palou and P. Salembier. 2013. Monocular Depth Ordering Using T-Junctions and Convexity Occlusion Cues. *IEEE Transactions on Image Processing* 22, 5 (May 2013), 1926–1939. <https://doi.org/10.1109/TIP.2013.2240002>
- [28] Pedro O Pinheiro, Ronan Collobert, and Piotr Dollar. 2015. Learning to segment object candidates. In *Advances in Neural Information Processing Systems*. 1990–1998.
- [29] Pedro O Pinheiro, Tsung-Yi Lin, Ronan Collobert, and Piotr Dollar. 2016. Learning to refine object segments. In *European Conference on Computer Vision*. Springer, 75–91.
- [30] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*. 91–99.
- [31] Andrew Senior, Arun Hampapur, Ying-Li Tian, Lisa Brown, Sharath Pankanti, and Ruud Bolle. 2006. Appearance models for occlusion handling. *Image and Vision Computing* 24, 11 (2006), 1233–1243.
- [32] Nathan Silberman, Lior Shapira, Ran Gal, and Pushmeet Kohli. 2014. A contour completion model for augmenting surface reconstructions. In *European Conference on Computer Vision*. Springer, 488–503.
- [33] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).
- [34] P. Smith, T. Drummond, and R. Cipolla. 2004. Layered motion segmentation and depth ordering by tracking edges. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 26, 4 (April 2004), 479–494. <https://doi.org/10.1109/TPAMI.2004.1265863>
- [35] Joseph Tighe, Marc Niethammer, and Svetlana Lazebnik. 2014. Scene parsing with object instances and occlusion ordering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 3748–3755.
- [36] Jonas Uhrig, Marius Cordts, Uwe Franke, and Thomas Brox. 2016. Pixel-level encoding and depth layering for instance-level semantic labeling. In *German Conference on Pattern Recognition*. Springer, 14–25.
- [37] John Winn and Jamie Shotton. 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)*, Vol. 1. IEEE, 37–44.
- [38] J. Xiao and M. Shah. 2005. Motion layer extraction in the presence of occlusion using graph cuts. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 27, 10 (Oct 2005), 1644–1659. <https://doi.org/10.1109/TPAMI.2005.202>
- [39] Hongyang Xue, Shengming Zhang, and Deng Cai. 2017. Depth image inpainting: Improving low rank matrix completion with low gradient regularization. *IEEE Transactions on Image Processing* 26, 9 (2017), 4311–4320.
- [40] Ziyu Zhang, Alexander G. Schwing, Sanja Fidler, and Raquel Urtasun. 2015. Monocular Object Instance Segmentation and Depth Ordering With CNNs. In *The IEEE International Conference on Computer Vision (ICCV)*.
- [41] Yan Zhu, Yuandong Tian, Dimitris Metaxas, and Piotr Dollar. 2017. Semantic amodal segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1464–1472.