

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065

# Human De-occlusion: Invisible Perception and Recovery for Humans

Anonymous CVPR submission

Paper ID 1550

## Abstract

In this paper, we tackle the problem of human de-occlusion which reasons about the occluded segmentation masks and invisible appearance content of humans. In particular, a two-stage framework is proposed to estimate the invisible portions and paint content inside. For the stage of mask completion, a stacked network structure is proposed to refine inaccurate masks from the instance segmentation model and predict the integrated masks simultaneously. Additionally, the guidance from human parsing and typical pose masks are leveraged to bring prior information. For the stage of content recovery, a novel parsing guided attention module is applied to isolate body parts and capture context information across multiple scales. Besides, an Amodal Human Perception dataset (AHP) is proposed to settle the task of human de-occlusion. AHP has advantages of providing annotations from real-world scenes and the number of humans is comparatively larger than other amodal perception datasets. Based on this dataset, experiments demonstrate that our method performs over the state-of-the-arts in both tasks of mask completion and content recovery.

## 1. Introduction

Visual recognition tasks have witnessed significant advances driven by deep learning, such as classification [21, 14], detection [39, 10] and segmentation [29, 52]. Despite the achieved progress, *amodal perception*, i.e. to reason about the occluded parts of objects, is still challenging for vision models. In contrast, it is easy for humans to interpolate the occluded portions with human visual systems [60, 35]. To reduce the recognition ability gap between vision models and human, recent work has proposed methods to infer the occluded portions of objects, including estimating the invisible segmentation [26, 60, 38, 51, 55] and recovering the invisible content of objects [6, 51, 55].

In this paper, we aim at the problem of estimating the invisible masks and content for humans. We refer to such a task as *human de-occlusion*. Compared with general

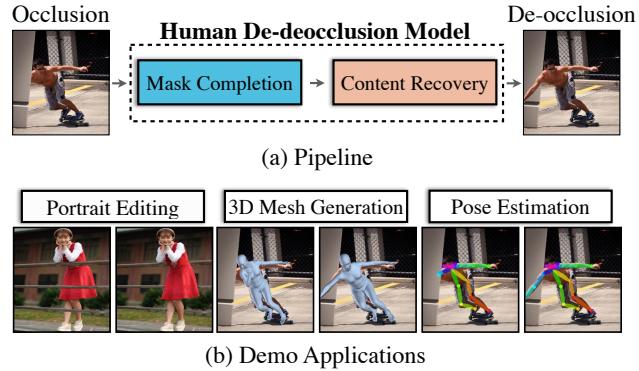


Figure 1. (a) The pipeline of our framework to tackle the task of human de-occlusion, which contains two stages of mask completion and content recovery. (b) Some demo applications which demonstrate better results can be obtained after human de-occlusion.

amodal perception, human de-occlusion is a more special and important task, since completing human body plays key roles in many vision tasks, such as portrait editing [1], 3D mesh generation [19], and pose estimation [11] as illustrated in Fig 1 (b) . Different from common object de-occlusion task (e.g., vehicle, building, furniture), human de-occlusion presents new challenges in three aspects. First, humans are non-rigid bodies and their poses vary dramatically. Second, the context relations between humans and backgrounds are inferior to common objects which are limited to specific scenes. Third, to segment and recover the invisible portions of humans, the algorithm should be able to recognize the occluded body parts to take symmetrical or interrelated patches to paint the content.

To precisely generate the invisible masks and content of humans, we propose a two-stage framework to accomplish human de-occlusion as shown in Fig 1 (a). The first stage segments the invisible portions of humans, and the second stage paints the content inside the invisible regions obtained in the previous stage. In the testing phase, the two stages are cascaded into an end-to-end manner.

Previous methods [55, 6, 51] adopt perfect modal mask as an extra input (e.g., ground truth annotation or well-

108 segmented mask), and expect to output the amodal mask.  
 109 However, obtaining ground truth or accurate segmented  
 110 mask is non-trivial in actual applications. To address the is-  
 111 sue, our first stage utilizes a stacked hourglass network [32]  
 112 to segment the invisible regions progressively. Specifically,  
 113 the network first applies a hourglass module to refine the  
 114 inaccurate input modal mask, then a secondary hourglass  
 115 module is applied to estimate the integrated amodal mask.  
 116 In addition, our network benefits from human parsing and  
 117 typical pose masks. Inspired by [22, 51], human parsing  
 118 pseudo labels are served as auxiliary supervision to bring  
 119 prior cues and typical poses are introduced as references.  
 120

121 In the second stage, our model paints visual content in-  
 122 side the invisible portions. Different from typical inpainting  
 123 methods [37, 28, 53, 31, 50], the context information inside  
 124 the human should be excessively explored. To this end, a  
 125 novel parsing guided attention (PGA) module is proposed  
 126 to capture and consolidate the context information from the  
 127 visible portions to paint the missing content across multiple  
 128 scales. Briefly, the module contains two attention streams.  
 129 The first path isolates different body parts to obtain the re-  
 130 lations of the missing part with others. The second spatially  
 131 calculates the relationship between the invisible and the  
 132 visible portions to capture contextual information. Then  
 133 the two streams are fused by concatenation and the mod-  
 134 ule outputs an enhanced deep feature. The module works at  
 135 different scales to possess our model of stronger recovery  
 136 capability and better performance.

137 As most existing amodal perception datasets [60, 38, 16]  
 138 focus on amodal segmentation <sup>1</sup> and few human images  
 139 are involved, an amodal perception dataset specified for hu-  
 140 man category is required to evaluate our method. On ac-  
 141 count of this, we introduce an Amodal Human Perception  
 142 dataset, namely AHP. There are three main advantages of  
 143 our dataset: **a)** the number of humans in AHP is compara-  
 144 tively larger than other amodal perception datasets. **b)** the  
 145 occlusion cases synthesized from AHP own amodal seg-  
 146 mentation and visual content ground-truths from real-world  
 147 scenes, hence subjective judgements and consistency of the  
 148 invisible portions from individual annotators are unneces-  
 149 sary; **c)** the occlusion cases with expected occlusion distri-  
 150 bution can be readily obtained. Existing amodal perception  
 151 datasets have fixed occlusion distributions and apparently it  
 152 will bring gaps in specified scenarios;

153 Our contributions are summarized as follows:

- 154 • We propose a two-stage framework for precisely gen-  
 155 erating the invisible parts of humans. To the best of our  
 156 knowledge, this is the first study which both considers  
 157 human mask completion and content recovery.
- 158 • A stacked network structure is proposed to do mask

159 <sup>1</sup>The term ‘modal segmentation’ and ‘amodal segmentation’ is used to  
 160 refer to the visible and integrated portions of an object respectively.  
 161

162 completion progressively, which specifically refines  
 163 inaccurate modal mask. Additionally, human parsing  
 164 and typical pose masks are introduced to bring prior  
 165 information.

- 166 • We propose a novel parsing guided attention (PGA)  
 167 module to capture body parts guidance and context in-  
 168 formation across multiple scales.
- 169 • A dataset, namely AHP, is proposed to settle the task  
 170 of human de-occlusion.

## 2. Related Work

**Amodal Segmentation and De-occlusion.** The task of  
 175 modal segmentation is to assign categorical or instance la-  
 176 bel for each visible pixel in an image, including seman-  
 177 tic segmentation [29, 4, 58, 34] and instance segmenta-  
 178 tion [39, 13, 3]. Differently, amodal segmentation aims at  
 179 segmenting visible and recovering invisible parts, which is  
 180 equivalent to the integrated mask, for each instance. Re-  
 181 search on amodal segmentation emerges from [26] which  
 182 iteratively enlarges detected box and recomputes heatmap  
 183 of each instance based on modally annotated data. After-  
 184 wards, AmodalMask [60], MLC [38], ORCNN [8] and PC-  
 185 Nets [55] push forward the field of amodal segmentation by  
 186 releasing datasets or techniques.

187 De-occlusion [55, 6, 51] targets at predicting the invis-  
 188 ible content of instances. It is different from general inpaint-  
 189 ing [28, 53, 31, 50] which recovers the missing areas (man-  
 190 ual holes) to make the image look reasonable. At present,  
 191 de-occlusion usually depends on the invisible masks pre-  
 192 dicted by amodal segmentation. SeGAN [6] studies amodal  
 193 segmentation and de-occlusion of indoor objects with syn-  
 194 synthetic data. Yan *et al.* [51] propose an iterative framework  
 195 to complete cars. Xu *et al.* [49] aim at portrait completion  
 196 by a two-stage strategy without involving amodal segmen-  
 197 tation. Compared with them, we jointly tackle the task of  
 198 human amodal segmentation and de-occlusion, and appar-  
 199 ently humans have larger variations in shapes and texture  
 200 details than rigid vehicles and furniture.

**Amodal Perception Datasets.** Amodal perception is a  
 203 challenging task [8, 38, 60, 6] aiming at recognizing the oc-  
 204 cluded parts of instances. Zhu *et al.* [60] point out that hu-  
 205 mans are able to predict the occluded regions with high de-  
 206 grees of consistency though the task is ill-posed. There are  
 207 some pioneers building amodal segmentation datasets. CO-  
 208 COA [60] elaborately annotates modal and amodal masks  
 209 of totally 5,000 images originates from COCO [27]. To  
 210 alleviate the problem that the data scale of COCOA is too  
 211 small, Qi *et al.* [38] establish KINS with 14,991 images  
 212 for two main categories of ‘people’ and ‘vehicle’. In [6],  
 213 Ehsani *et al.* introduce a synthetic indoor dataset namely  
 214 DYCE. SAIL-VOS [16] is also synthetic but a video dataset

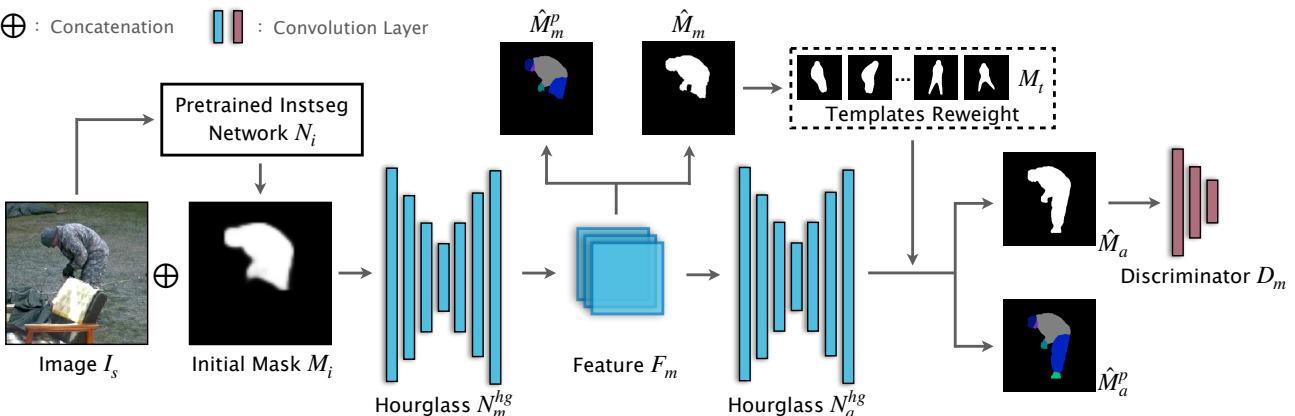


Figure 2. **Illustration of the refined mask completion network.** An pretrained instance segmentation network  $N_i$  is applied to obtain initial modal mask  $M_i$  from the input image  $I_s$  which contains an occluded human. Then one hourglass module  $N_m^{hg}$  outputs the refined modal mask  $\hat{M}_m^p$  and the corresponding parsing results  $\hat{M}_m^p$ . To obtain the complete mask, another hourglass module  $N_a^{hg}$  is stacked with collected template masks served as prior appearance cues and it finally outputs amodal mask  $\hat{M}_a$  and parsing results  $\hat{M}_a^p$ . At the end, a discriminator  $D_m$  is applied to improve the quality of generated amodal masks.

by leveraging game simulator. Our dataset AHP focuses on amodal perception of humans on account of its wide applications. Particularly, the differences with the above datasets mainly lie in two aspects: **a)** our dataset greatly extends the quantity ( $\sim 56k$ ) of humans in real-world scene and provides amodal mask annotation; **b)** our dataset emphasizes on completing both invisible mask and content with trustworthy supervision provided, to overcome the lame learning methods with only modal supervision in [55, 8].

**Human Related Works.** There are extensive literature on humans, such as detection [48, 59, 56], parsing [54, 57, 46, 41] and pose estimation [33, 2, 30, 45] etc. Recently MaskGAN [24] combines parsing with face image manipulation. Our work attempts to benefit from these studies for better amodal segmentation and de-occlusion.

### 3. Method

#### 3.1. Overview

We propose a two-stage framework to accomplish human de-occlusion. The first stage segments the invisible portions of humans and the second paints content inside, as shown in Fig 2 and Fig 3 respectively.

#### 3.2. Refined Mask Completion Network

The refined mask completion network targets at segmenting the invisible masks. As mentioned in Sec. 1, most methods first take perfect but non-trivial modal masks as input and the predicted amodal masks will be subtracted with them to get the expected invisible masks. Some of them also apply modal masks from instance segmentation methods, but the gap between amodal perception and instance segmentation annotations are not being noticed. To

this end, our mask completion network firstly refines the trustless initial modal masks. Specifically, an existing instance segmentation network  $N_i$  is applied to obtain the initial modal masks  $M_i$  from the input image  $I_s$  which contains an occluded human. Then the image and the initial mask is concatenated and fed into one hourglass module  $N_m^{hg}$  to obtain the refined binary modal mask  $\hat{M}_m^p$ . The supervision on corresponding parsing results  $\hat{M}_m^p$  is accompanied to strengthen the semantic understanding of body parts. The modal recognition process can be formulated as:

$$\hat{M}_m, \hat{M}_m^p = N_m^{hg}(I_s, N_i(I_s)). \quad (1)$$

Next, our network needs to complete the refined modal mask. One direct solution is to augment another amodal branch to accomplish this task. However, we empirically find it will significantly delegate the performance of the refined modal masks. We suspect that completing the integrated amodal mask requires homogeneous feature both in the occluded and visible portions of the human, which is opposite to the previous modal mask recognition paying attention on the visible regions only. Accordingly, another hourglass module  $N_a^{hg}$  is stacked behind the first one to estimate the amodal mask  $\hat{M}_a$  and parsing results  $\hat{M}_a^p$ .

Since it can be asserted that the target amodal masks are integrated humans, some typical poses can be implanted into the network as prior appearance cues. We collect a batch of template masks by running k-means on the training set, denoted as  $M_t$ . Then the  $\ell_2$  distance  $D_{m,t}$  of  $\hat{M}_m$  with each template mask is calculated as the attention weight vector by  $W_{m,t} = 1/D_{m,t}$ . The vector  $W_{m,t}$  is multiplied back with the template masks to highlight suitable candidates and a convolution layer without activation is applied to combine these re-weighted templates as an additional

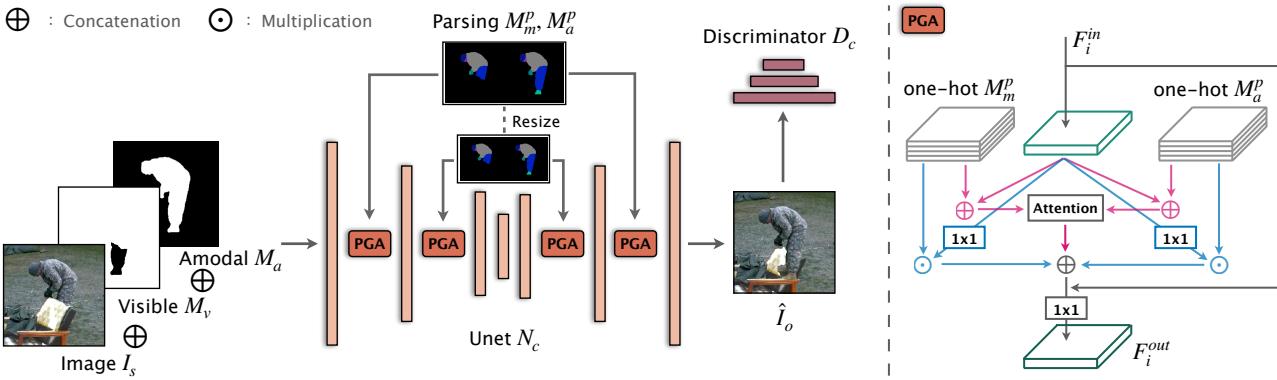


Figure 3. **Illustration of the parsing guided content recovery network.** **Left:** We adopt Unet [40] with partial convolution [28] as the basic architecture. The image  $I_s$  concatenated with visible mask  $M_v$  and amodal mask  $M_a$  is passed into the network  $N_c$  to paint inside the invisible portions. To leverage the guidance of body part cues  $M_m^p$  and  $M_a^p$  from the previous stage, our proposed Parsing Guided Attention (PGA) module enhances the deep features at multiple scales. Finally a discriminator  $D_c$  is applied to identify the quality of image  $I_o$ . **Right:** The details of our proposed PGA module. The module distributes one-hot parsing results to the input feature  $F_i^{in}$  with two attention streams. The first stream (cyan) decomposes the feature into different body parts and the second (magenta) tries to establish the pixel-level relationship between the visible context and the invisible regions.

feature before making predictions. The amodal completion process is given by:

$$\hat{M}_a, \hat{M}_a^p = N_a^{hg} (F_m \oplus \hat{M}_m, \text{Conv.}(M_t \odot W_{m,t})), \quad (2)$$

where  $F_m$  is the intermediate feature from the last hourglass module.

The cross-entropy loss  $\mathcal{L}_{CE}(\cdot)$  is applied to supervise the modal and amodal predictions and human parsing. The idea in adversarial learning can be borrowed to boost the performance by adopting a discriminator  $D_m$  and the perceptual loss [9]  $\mathcal{L}_{prec}(\cdot)$  which measures the distance on extracted features of the generated and the ground truth masks. The several loss functions are formulated as follows:

$$\begin{aligned} \mathcal{L}_{seg} &= \mathcal{L}_{CE}(\hat{M}_m, M_m) + \mathcal{L}_{CE}(\hat{M}_a, M_a) + \\ &\quad \mathcal{L}_{CE}(\hat{M}_m^p, M_m^p) + \mathcal{L}_{CE}(\hat{M}_a^p, M_a^p), \\ \mathcal{L}_{adv} &= \mathbb{E}_{\hat{M}_a} [\log(1 - D_m(\hat{M}_a))] + \mathbb{E}_{M_a} [\log D_m(M_a)], \\ \mathcal{L}_{gen} &= \mathcal{L}_{\ell 1}(\hat{M}_a, M_a) + \mathcal{L}_{prec}(\hat{M}_a, M_a), \end{aligned} \quad (3)$$

where  $\mathcal{L}_{\ell 1}(\cdot)$  denotes the reconstruction loss. Then the final loss can be summarized with proper coefficients:

$$\mathcal{L}_m = \lambda_1 \mathcal{L}_{seg} + \lambda_2 \mathcal{L}_{adv} + \lambda_3 \mathcal{L}_{gen}. \quad (4)$$

### 3.3. Parsing Guided Content Recovery Network

The mask completion network is able to localize the invisible regions and body parts by subtracting the refined modal results from the amodal ones, and the goal of this stage is to paint visual content inside. Different from general inpainting methods, the missing content is some parts of a human instead of other surroundings to make the image

look plausible. In spite of this distinction, some well studied methods can be referred as base structures and here we adopt Unet [40] with partial convolution [28] as the architecture, as shown on the left part of Fig. 3.

Our pipeline is straightforward. At first, the visible mask  $M_v$  and the amodal mask  $M_a$  are concatenated with the image  $I_s$  containing an occluded human as input. The visible mask tells the network  $N_c$  which pixels need to be painted and the amodal mask points out where the integrated human occupies in the image. Additionally, the body part cues  $M_m^p$  and  $M_a^p$  from the previous stage are aggregated into the deep features by our proposed Parsing Guided Attention (PGA) module at multiple scales. Finally the network outputs a human recovered image  $\hat{I}_o$  and a discriminator  $D_c$  is applied to identify the quality of image  $\hat{I}_o$ . This stage can be formulated as:

$$\hat{I}_o = N_c (I_s \oplus M_v \oplus M_a, M_m^p, M_a^p). \quad (5)$$

The PGA module is depicted on the right part of Fig 3. Take some scale  $i$  for example, the module takes in deep feature  $F_i^{in}$  and the parsing masks  $M_m^p$  and  $M_a^p$  which are converted into two one-hot logits and resized to the same spatial size with  $F_i^{in}$ . It contains two attention streams. The first stream (cyan) decomposes the feature into different body parts and compare them. In specific, the feature  $F_i^{in}$  is reduced to the same channel number with the parsing logits (*i.e.* 19), then it is element-wisely multiplied with the two parsing logits to distribute the feature in different body parts. Then the two distributed features are concatenated and a  $1 \times 1$  convolution layer is applied to yield feature of useful body parts. The second stream (magenta) tries to establish the pixel-level relationship between the visible context and the invisible regions. The difference with the self

attention [47] is that self attention wildly calculates pixel-wise relations regardless of the semantics of the pixels. In particular, we concatenate the two logits behind the input feature, and two  $1 \times 1$  convolution layers  $\phi(\cdot)$  and  $\psi(\cdot)$  are followed to extract key features for visible and amodal regions respectively, denoted as  $K_{vis} = \phi(F_i^{in} \oplus M_m^p)$  and  $K_{amo} = \psi(F_i^{in} \oplus M_a^p)$ . Then the relationship matrix  $R$  can be obtained by:

$$\begin{aligned}\tilde{R} &= (M_v \odot K_{vis})^T ((1 - M_v) \odot K_{amo}); \\ R &= \text{Softmax}(\tilde{R}, \dim = 0) \in \mathbb{R}^{HW \times HW},\end{aligned}\quad (6)$$

where  $H$  and  $W$  denote the height and width viewed as collapsible spatial dimension. The matrix means that given a point of the invisible regions, it gives the pair-wised relevance of the point with the points of all visible regions. Since the relationship matrix contains value for each point, a matrix multiplication of  $R$  and the input feature can be applied to extract related information from the visible. At last, the outputs of the two streams are concatenated with the input feature  $F_i^{in}$ , and a  $1 \times 1$  convolution layer is followed to reduce the channel number same with  $F_i^{in}$ .

The loss function to optimize our content recovery network is formulated as follows:

$$\begin{aligned}\mathcal{L}_c = \beta_1 (\mathbb{E}_{\hat{I}_o} [\log(1 - D_c(\hat{I}_o))] + \mathbb{E}_{I_o} [\log D_c(I_o)]) + \\ \beta_2 \mathcal{L}_{\ell 1}(\hat{I}_o, I_o) + \beta_3 \mathcal{L}_{prec}(\hat{I}_o, I_o) + \\ \beta_4 \mathcal{L}_{style}(\hat{I}_o, I_o),\end{aligned}\quad (7)$$

where  $\mathcal{L}_{style}(\cdot)$  denotes the style loss proposed in [28].

## 4. The Amodal Human Perception Dataset

In this section, we describe how we collect human images and obtain their annotations with minimal manual effort. Our method adopts a straightforward pipeline and takes advantage of existing works on human segmentation. As a result, the proposed Amodal Human Perception dataset, namely AHP, contains unoccluded human images with masks. Therefore the occlusion cases can be synthesized by pasting occluders from other datasets onto humans. Moreover, some informative statistics are summarized to analyze our dataset.

### 4.1. Data Collection and Filtering

Instead of collecting images of people from the Internet, we capitalize on several large instance segmentation and detection datasets to acquire human images. Then a segmentation model is applied to obtain masks for the humans with only box-level annotations. The reason to establish pixel-level annotations is that they can be utilized for occlusion synthesis. Finally, an image filtering scheme with manual effort is applied to construct our dataset. The detailed process is as follows.

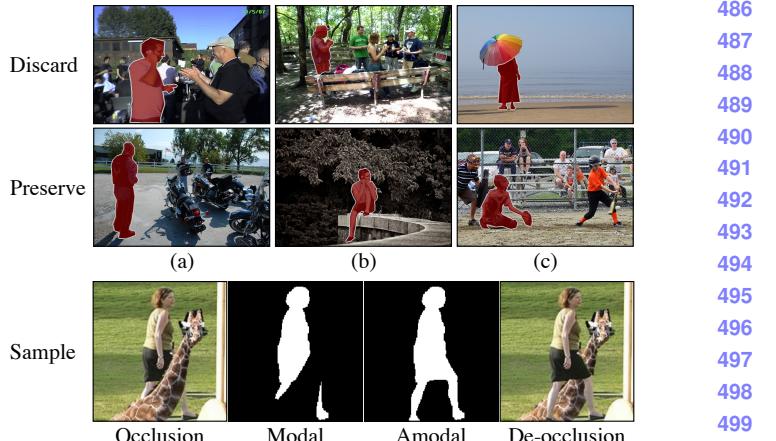


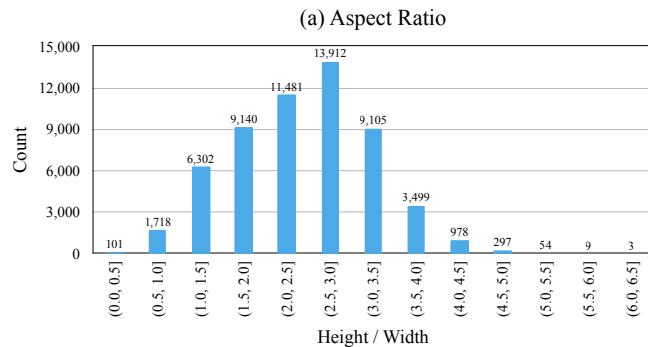
Figure 4. The first two rows show some ‘Discard’ and ‘Preserve’ exemplars and the last row demonstrates a synthesized occlusion image with ground-truths of the modal and amodal masks and the final de-occlusion image.

**(1) Image Acquisition.** We collect human images from several large-scale instance segmentation and detection datasets, including COCO [27], VOC [7] (with SBD [12]), LIP [11], Objects365 [42] and OpenImages [23]. Since each instance in these datasets has been annotated with category label, we simply keep the instances labeled with ‘human’ for further processing. In addition, the human instances that are too small (*e.g.* < 300 pixels) or have overlaps within the ‘gutter’ (*e.g.* 5 pixels) along four image boundaries are dropped because parts of them are very likely to be out of view.

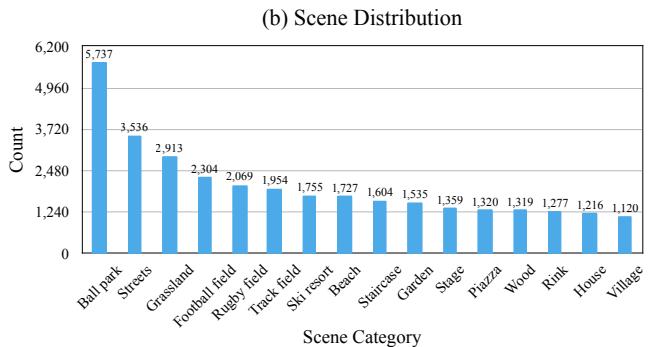
**(2) Human Segmentation.** The two largest datasets of Objects365 and OpenImages provide only box-level annotations, which means a great many human instances lack of pixel-level segmentation masks. Hence, a human segmentation model, *i.e.* Deeplab [5] trained on private massive human data, is applied to obtain high-quality segmentation results. In specific, the box for each instance is enlarged by 3 times in height and width to get a cropped image patch which will be fed into the model. Note we only keep segmentation results inside the original annotated box. Although the segmentation algorithm may fail on complicated cases, the issue will be addressed next.

	COCO†	OpenImages	Objects365
Human Count	24,806	314,359	669,166
Preserved Count	5,940	17,677	32,982
Preserved Ratio	22.1%	5.62%	4.93%

Table 1. Statistics of the total number of human category and the preserved ratio from each original dataset. COCO† means COCO unites other datasets, *e.g.* VOC.



(a) Aspect Ratio



(b) Scene Distribution

Figure 5. Chart (a) shows the distribution of aspect ratio and chart (b) demonstrates the distribution of scene categories in our dataset.

**(3) Filtering Scheme.** Since we aim at collecting unoccluded human instances with accurate masks, there is a need to retain the satisfied samples. Three options are set up to route each sample: **Discard**: the human is occluded by other instances (*e.g.* desk, car or other humans) or parts of him/her are out of view; **Preserve**: the human is not occluded and the segmentation is fine; **Refine**: the human is not occluded but the segmentation result is not satisfied. Some ‘Discard’ and ‘Preserve’ exemplars are demonstrated in Fig 4. Note for the case of ‘Preserve-(c)’ that a baseball player wears a glove, we will keep it as the object affects the body structure not too much. But ‘Discard-(c)’ is an opposite case. There are 7 well-trained annotators filtering the collected human instances. Finally, the number of ‘Preserve’ samples reaches over 50,000 and we think it is already enough to do meaningful research. The samples with ‘Refine’ label also will be released for potential exploration.

Once our dataset is built, occlusion training samples can be synthesized. As shown in the last row of Fig 4, an giraffe is cut out from other datasets and pasted onto the woman to form an occlusion case. Since we have the mask annotation of the woman (amodal mask), the ground-truths of the modal mask and the invisible content can be easily inferred.

## 4.2. Data Statistics

The AHP dataset consists of 56,599 images following the aforementioned steps. Each sample costs about 3.26 seconds on average. Table 1 shows the total number of human category and the preserved ratio from each original dataset. It is noticeable that the preserved ratios of OpenImages and Objects365 are significantly smaller than COCO. Besides, the distributions of aspect ratio (height / width) and scene categories are shown in Fig 5. Since the preserved humans are usually standing, it is reasonable that there are nearly half of the samples whose aspect ratios are between 1.5 to 3.5. To estimate the distribution of scene categories, a scene recognition model [44] is applied. Here we only show categories that contain more than 1,000 samples.

## 5. Experiments

### 5.1. Implementation Details

**Network details and optimization.** The input sizes of mask completion and content recovery networks are both  $256 \times 256$ . We adopt CenterMask [25] as our pretrained instance segmentation model  $N_i$ . HRNet [43] is used to generate the pseudo labels of body parts. The hourglass modules of  $N_m^{hg}$  and  $N_a^{hg}$  and the Unet  $N_c$  are inherited from [55] for fair comparison, and they are initialized with random weights. Both  $D_m$  and  $D_c$  adopt PatchGAN [17] and share the same structure with four convolution layers. We set  $\lambda_1 = \lambda_2 = 1, \lambda_3 = 0.1$  and  $\beta_1 = 0.1, \beta_2 = \beta_3 = 1, \beta_4 = 40$  in all experiments. We use SGD with momentum (batch 32, lr  $1e-3$ , iterations 48k) and Adam [20] (batch 16, lr  $1e-4$ , iterations 230k) to optimize the completion and recovery networks respectively. We use PyTorch [36] framework with a NVIDIA Tesla V100 to conduct experiments. All of our code and dataset will be opened.

**Data splits.** One advantage of our AHP dataset is that the distribution of human occlusion ratio can be controlled manually. Based on our observation, we set the probability density of occlusion ratio to  $P_{0 \sim 0.1} = P_{0.1 \sim 0.2} = P_{0.3 \sim 0.4} = 1/3$  when training. We draw random instances from COCO [27] to synthesize occlusion cases. At present, our method adopts the simplest synthesis strategy and we leave the problem of how to generate better samples for future research. Our validation set contains 891 images, which is augmented from 297 integrated human samples each with three different occluders. The probability density of the occlusion ratio for the validation split is set to  $P_{0 \sim 0.1} = P_{0.1 \sim 0.2} = P_{0.3 \sim 0.4} = P_{0.4 \sim 0.5} = 1/4$  and this set is fixed after generated. To verify the effectiveness of our methods in real scenes, several photo editors collect a number of images and they move the foreground instances with similar depths onto the humans manually. Then these artificial occlusion cases are voted by another group of humans and only all passed samples are saved as our test set. The test set contains only 15 images due to its steep cost.

Method	Syn.		Real	
	$\ell_1 \downarrow$	IoU $\uparrow$	$\ell_1 \downarrow$	IoU $\uparrow$
Mask-RCNN [13]	0.2402	78.4/26.9	0.2961	72.4/26.3
Deeplab [5]	0.2087	70.7/20.9	0.3276	59.0/15.9
Pix2Pix [18]	0.2329	69.6/19.2	0.2993	64.9/17.6
SeGAN [6]	0.2545	76.7/23.6	0.2963	74.7/24.2
OVSR [51]	0.1830	80.2/28.1	0.2258	80.4/29.8
PCNets [55]	0.1959	83.1/29.1	0.2229	79.5/30.1
Ours	<b>0.1500</b>	<b>84.6/43.7</b>	<b>0.2138</b>	<b>83.6/41.4</b>

Table 2. The comparison results of mask completion task on our AHP dataset. ‘Syn.’ and ‘Real’ denote synthesized and real validation images. Our method improves over the other techniques both in  $\ell_1$  error of the amodal masks and IoUs of the amodal and invisible masks.

**Evaluation metrics.** To evaluate the quality of the completed mask, we adopt  $\ell_1$  distance and Intersection over Union (IoU) as our metrics. For the recovered image, we adopt  $\ell_1$  distance and Fréchet Inception Distance (FID) [15] score which measures the similarity between the ground-truth images and our generated results.

## 5.2. Results

**Quantitative comparison.** We compare our method to recent state-of-the-art techniques on our AHP dataset. For the mask completion task, some general methods widely applied in other fields like Mask-RCNN [13], Deeplab [5] and Pix2Pix [18] are adopted. And another group methods are specialized to solve the amodal perception task like SeGAN [6], OVSR [51] and PCNets [55]. Among them, OVSR and PCNets declare they have achieved or surpassed state-of-the-art performance. As shown in Table 2, our method has lower  $\ell_1$  error and better IoU results on amodal masks both on synthesized and real validation images. It is worth mentioning that we have a significant advantage on the invisible masks. For the task of content recovery, we also selected the two types of general and specialized methods to compare. As shown in Table 3, our method again performs over others on synthesized and real validation sets. Note that due to there are a few real samples, the metric FID has fluctuated widely as it measures the feature representation similarity between two sets of images.

**Qualitative comparison.** Fig 6 shows a sample of predicted amodal masks and recovered images on our AHP dataset. For the mask completion task, Pix2Pix [18] and SeGAN [6] have difficulties at completing the occluded humans and their results are not satisfying at all. Our method has more symmetrical and better predicted masks compared to the two methods of OVSR [51] and PCNets [55]. For the task of content recovery, Deepfillv2 [53] seems blurring and there are apparent artifacts in the results of Pix2Pix [18] and SeGAN [6]. Again our method has better recovery perfor-

Method	Syn.		Real	
	$\ell_1 \downarrow$	FID $\downarrow$	$\ell_1 \downarrow$	FID $\downarrow$
Pix2Pix [18]	0.1126	19.66	0.0977	44.57
Deepfillv2 [53]	0.1127	21.61	0.0978	34.39
SeGAN [6]	0.1122	23.01	0.0980	37.68
OVSR [51]	0.0940	27.15	0.0859	51.26
PCNets [55]	0.0936	18.50	0.0853	36.02
Ours	<b>0.0519</b>	<b>13.85</b>	<b>0.0562</b>	<b>21.52</b>

Table 3. The comparison results of content recovery task on our AHP dataset. ‘Syn.’ and ‘Real’ denote synthesized and real validation images.

mance compared to OVSR [51] and PCNets [55].

## 5.3. Ablation Study

To understand our framework further, we conduct extensive experiments on the synthesized and fixed validation images (Sec 5.1) to prove the effectiveness of our method.

**Refined mask completion network.** Table 4 shows the results. The baseline takes the input image  $I_s$  and the initial mask  $M_i$ , then outputs amodal mask  $\hat{M}_a$  by networks  $N_m^{hg}$  and  $N_a^{hg}$  (line 1). The discriminator  $D_m$  improves the quality of the amodal mask  $M_a$  by 0.3% (line 2). To obtain precise invisible mask, the introduction of the modal segmentation  $M_m$  significantly refines the modal mask result by 4.2% but degrades the performance of the amodal by 1.2%. Fortunately the final invisible mask has 4.3% improvements (line 3). It shows that the modal segmentation has negative effects on the amodal completion, and we speculate obtaining the amodal mask requires homogeneous feature both in the occluded and visible portions of the humans. To solve the problem, the accompanying parsing branches ( $M_m^p, M_a^p$ ) are leveraged to bring in extra semantic guidance and template pose masks are utilized. The parsing improves 2.3% and 0.8% for the modal and amodal tasks respectively (line 4) and the network benefits from the templates by 0.7% and 0.2% (line 5). It is noticeable that the IoU metric of the invisible mask boosts 5.4% and 3.9%

Discriminator	Modal	Parsing	Templates	IoU $\uparrow$
1				77.5/84.0/30.3
2	✓			77.5/84.3/30.0
3	✓	✓		81.7/83.1/34.6
4	✓	✓	✓	84.0/83.9/40.0
5	✓	✓	✓	82.4/83.3/38.5
6	✓	✓	✓	<b>84.7/84.6/43.7</b>

Table 4. Ablation study of the refined mask completion network. The three columns of each IoU result represent the modal, amodal and invisible masks respectively.

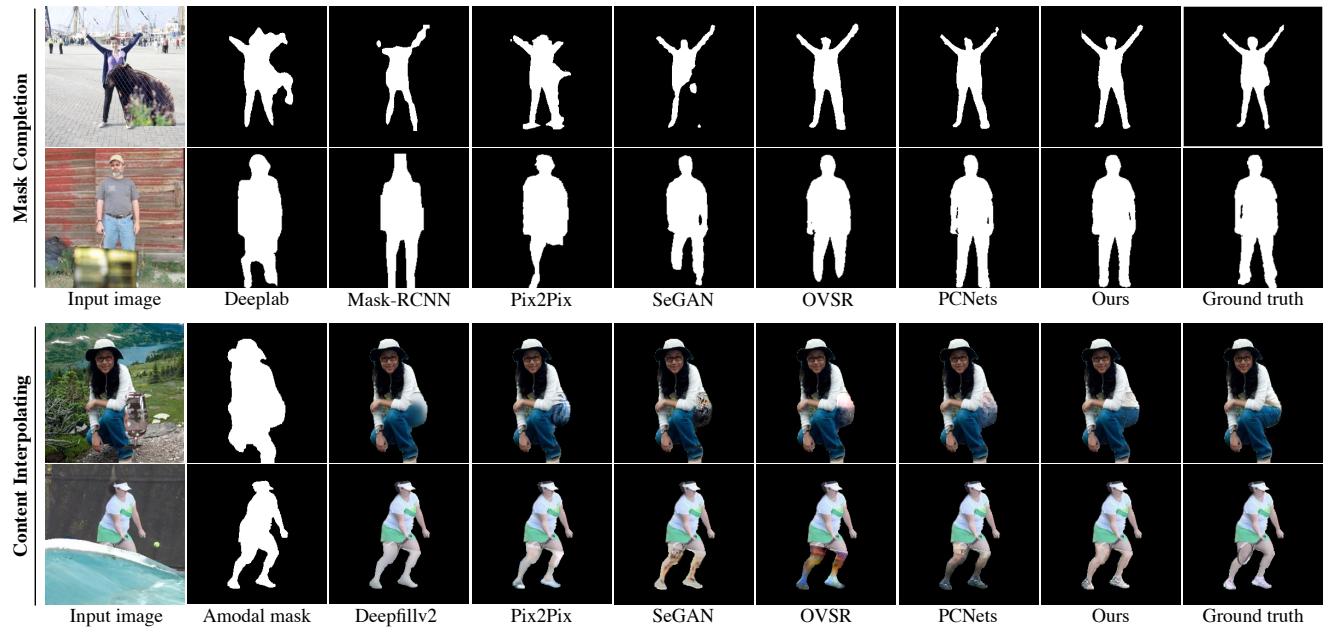


Figure 6. Qualitative comparison results of the mask completion and content recovery tasks on our AHP dataset.

<i>w</i>	0.0	0.1	0.3	0.5	0.7	0.9	1.0
FID $\downarrow$	18.04	18.19	<b>17.85</b>	18.74	18.61	19.11	19.71

Table 5. Ablation study of the proportion of the background.

of the two proposed modules. Finally, after we unite these parts together, the IoUs of the modal and invisible masks are remarkably improved by 7.2% and 13.4% respectively (line 6 vs. line 1).

**Parsing guided content recovery network.** In Table 5, we analyze the effect of the proportion of the background first. Intuitively the network should interpolate the invisible portions by referring to the visible parts of the human only. However, the background may support context information to know where the content can be imitated. Therefore, the new input image can be written as:  $I_s' = I_s * M_a + I_s * (1 - M_a) * w$ , where  $w$  is the proportion of the background. And the baseline is the Unet with partial convolution as [55] with  $w = 1$ . The experiments show that it gains 1.9 points when  $w = 0.3$ . Further, our proposed PGA module is disassembled to two attention streams and analyzed as shown in Table 6. The two streams are evaluated individually and the comparison with simply cascading the two streams shows our structure has better performance. Specifically, the first attention stream (Attention.B) separating different body parts boosts the performance by 2.94 points and the second stream (Attention.T) of establishing the relationship between the visible context and the invisible regions gains 2.3 points. Assembling the two streams in a cascade manner yields 5.04 points improvement. Lastly,

Bg (0.3)	Attention.B	Attention.T	Structure	FID $\downarrow$
1				19.66
2	✓			17.85
3		✓		16.76
4			✓	17.37
5	✓	✓	✓	Cascade
6	✓	✓	✓	Fusion

Table 6. Ablation study of the parsing guided content recovery network. ‘Bg’ denotes the background, and the columns ‘Attention.B’ and ‘Attention.T’ correspond to the two attention streams. There are two structures to assemble them: ‘Cascade’ and ‘Fusion’.

our proposed structure depicted on the right part of Fig 3 further boosts extra 0.8 points.

## 6. Conclusion

In this paper, we tackle the problem of *human deocclusion* which is a more special and important task compared to de-occluding general objects. We propose a two-stage framework to settle it. By refining the initial mask from the pretrained model and completing the modal mask, our network is able to precisely predict the invisible regions. Then the content recovery network equipped with our proposed PGA module recovers the invisible details. In addition, our collected AHP dataset has advantages compared to the current amodal perception datasets. Extended studies on how to generate more realistic samples and how to exchange deep features of the two tasks will be explored in our future work.

864

## References

- [1] Marcelo Bertalmio, Guillermo Sapiro, Vincent Caselles, and Coloma Ballester. Image inpainting. In *Proceedings of the 27th annual conference on Computer graphics and interactive techniques*, pages 417–424, 2000. 1
- [2] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *IEEE conference on computer vision and pattern recognition*, pages 7291–7299, 2017. 3
- [3] Hao Chen, Kunyang Sun, Zhi Tian, Chunhua Shen, Yongming Huang, and Youliang Yan. BlendMask: Top-down meets bottom-up for instance segmentation. In *IEEE international conference on computer vision*, 2020. 2
- [4] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. 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, 2017. 2
- [5] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *European conference on computer vision (ECCV)*, pages 801–818, 2018. 5, 7
- [6] Kiana Ehsani, Roozbeh Mottaghi, and Ali Farhadi. Segan: Segmenting and generating the invisible. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 1, 2, 7
- [7] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010. 5
- [8] Patrick Follmann, Rebecca Kö Nig, Philipp Hä Rtinger, Michael Klostermann, and Tobias Bö Ttger. Learning to see the invisible: End-to-end trainable amodal instance segmentation. In *IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1328–1336. IEEE, 2019. 2, 3
- [9] Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. 4
- [10] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 580–587, 2014. 1
- [11] Ke Gong, Xiaodan Liang, Dongyu Zhang, Xiaohui Shen, and Liang Lin. Look into person: Self-supervised structure-sensitive learning and a new benchmark for human parsing. In *IEEE Conference on Computer Vision and Pattern Recognition*, July 2017. 1, 5
- [12] Bharath Hariharan, Pablo Arbelaez, Lubomir Bourdev, Subhransu Maji, and Jitendra Malik. Semantic contours from inverse detectors. In *International Conference on Computer Vision*, 2011. 5
- [13] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *IEEE international conference on computer vision*, pages 2961–2969, 2017. 2, 7
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. 1
- [15] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in neural information processing systems*, pages 6626–6637, 2017. 7
- [16] Y.-T. Hu, H.-S. Chen, K. Hui, J.-B. Huang, and A. G. Schwing. Sail-vos: Semantic amodal instance level video object segmentation – a synthetic dataset and baselines. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 2
- [17] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017. 6
- [18] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. *CVPR*, 2017. 7
- [19] Angjoo Kanazawa, Michael J. Black, David W. Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *Computer Vision and Pattern Recognition (CVPR)*, 2018. 1
- [20] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 6
- [21] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Neural Information Processing Systems*, pages 1097–1105, 2012. 1
- [22] Weicheng Kuo, Anelia Angelova, Jitendra Malik, and Tsung-Yi Lin. Shapemask: Learning to segment novel objects by refining shape priors. *IEEE International Conference on Computer Vision*, 2019. 2
- [23] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Tom Duerig, et al. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *arXiv preprint arXiv:1811.00982*, 2018. 5
- [24] Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive facial image manipulation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 5549–5558, 2020. 3
- [25] Youngwan Lee and Jongyoul Park. Centermask: Real-time anchor-free instance segmentation. 2020. 6
- [26] Ke Li and Jitendra Malik. Amodal instance segmentation. In *European Conference on Computer Vision*, pages 677–693. Springer, 2016. 1, 2
- [27] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 2, 5, 6
- [28] Guilin Liu, Fitsum A Reda, Kevin J Shih, Ting-Chun Wang, Andrew Tao, and Bryan Catanzaro. Image inpainting for

- 972       irregular holes using partial convolutions. In *Proceedings*  
973       of the European Conference on Computer Vision (ECCV),  
974       pages 85–100, 2018. 2, 4, 5
- 975       [29] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully  
976       convolutional networks for semantic segmentation. In *IEEE*  
977       conference on computer vision and pattern recognition,  
978       pages 3431–3440, 2015. 1, 2
- 979       [30] Yue Luo, Jimmy Ren, Zhouxia Wang, Wenxiu Sun, Jinshan  
980       Pan, Jianbo Liu, Jiahao Pang, and Liang Lin. Lstm pose  
981       machines. In *IEEE conference on computer vision and pattern  
982       recognition*, pages 5207–5215, 2018. 3
- 983       [31] Kamyar Nazeri, Eric Ng, Tony Joseph, Faisal Qureshi, and  
984       Mehran Ebrahimi. Edgeconnect: Structure guided image  
985       inpainting using edge prediction. In *IEEE International Con-  
986       ference on Computer Vision (ICCV) Workshops*, Oct 2019.  
987       2
- 988       [32] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hour-  
989       glass networks for human pose estimation. In *European con-  
990       ference on computer vision*, pages 483–499. Springer, 2016.  
991       2
- 992       [33] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hour-  
993       glass networks for human pose estimation. In *European con-  
994       ference on computer vision*, pages 483–499. Springer, 2016.  
995       3
- 996       [34] Hyeonwoo Noh, Seunghoon Hong, and Bohyung Han.  
997       Learning deconvolution network for semantic segmentation.  
998       In *IEEE international conference on computer vision*, pages  
999       1520–1528, 2015. 2
- 1000       [35] Stephen E Palmer. *Vision science: Photons to phenomenol-  
1001       ogy*. MIT press, 1999. 1
- 1002       [36] Adam Paszke, Sam Gross, Soumith Chintala, Gregory  
1003       Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Al-  
1004       ban Desmaison, Luca Antiga, and Adam Lerer. Automatic  
1005       differentiation in pytorch. 2017. 6
- 1006       [37] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor  
1007       Darrell, and Alexei A Efros. Context encoders: Feature  
1008       learning by inpainting. In *Proceedings of the IEEE con-  
1009       ference on computer vision and pattern recognition*, pages  
1010       2536–2544, 2016. 2
- 1011       [38] Lu Qi, Li Jiang, Shu Liu, Xiaoyong Shen, and Jiaya Jia.  
1012       Amodal instance segmentation with kins dataset. In *IEEE  
1013       Conference on Computer Vision and Pattern Recognition*,  
1014       pages 3014–3023, 2019. 1, 2
- 1015       [39] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun.  
1016       Faster r-cnn: Towards real-time object detection with region  
1017       proposal networks. In *Advances in neural information pro-  
1018       cessing systems*, pages 91–99, 2015. 1, 2
- 1019       [40] Olaf Ronneberger, Philipp Fischer, and Thomas Brox.  
1020       U-net: Convolutional networks for biomedical image segmen-  
1021       tation. In *Medical image computing and computer-assisted  
1022       intervention*, pages 234–241. Springer, 2015. 4
- 1023       [41] Tao Ruan, Ting Liu, Zilong Huang, Yunchao Wei, Shikui  
1024       Wei, and Yao Zhao. Devil in the details: Towards accurate  
1025       single and multiple human parsing. In *AAAI Conference on  
Artificial Intelligence*, volume 33, pages 4814–4821, 2019. 3
- 1026       [42] Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang  
Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365:  
1027       A large-scale, high-quality dataset for object detection. In  
1028       *IEEE international conference on computer vision*, pages  
1029       8430–8439, 2019. 5
- 1030       [43] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep  
1031       high-resolution representation learning for human pose esti-  
1032       mation. In *CVPR*, 2019. 6
- 1033       [44] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and  
1034       Alex Alemi. Inception-v4, inception-resnet and the im-  
1035       pact of residual connections on learning. *arXiv preprint  
arXiv:1602.07261*, 2016. 6
- 1036       [45] Manchen Wang, Joseph Tighe, and Davide Modolo. Com-  
1037       bining detection and tracking for human pose estimation in  
1038       videos. In *IEEE Conference on Computer Vision and Pattern  
1039       Recognition (CVPR)*, June 2020. 3
- 1040       [46] Wenguan Wang, Zhijie Zhang, Siyuan Qi, Jianbing Shen,  
1041       Yanwei Pang, and Ling Shao. Learning compositional neural  
1042       information fusion for human parsing. In *IEEE International  
1043       Conference on Computer Vision*, pages 5703–5713, 2019. 3
- 1044       [47] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaim-  
1045       ing He. Non-local neural networks. In *IEEE conference on  
1046       computer vision and pattern recognition*, pages 7794–7803,  
1047       2018. 5
- 1048       [48] Xinlong Wang, Tete Xiao, Yuning Jiang, Shuai Shao, Jian  
1049       Sun, and Chunhua Shen. Repulsion loss: Detecting pedes-  
1050       trians in a crowd. In *IEEE Conference on Computer Vision  
and Pattern Recognition*, pages 7774–7783, 2018. 3
- 1051       [49] Xian Wu, Rui-Long Li, Fang-Lue Zhang, Jian-Cheng Liu,  
1052       Jue Wang, Ariel Shamir, and Shi-Min Hu. Deep portrait im-  
1053       age completion and extrapolation. *IEEE Transactions on Im-  
age Processing*, 29:2344–2355, 2019. 2
- 1054       [50] Chaohao Xie, Shaohui Liu, Chao Li, Ming-Ming Cheng,  
1055       Wangmeng Zuo, Xiao Liu, Shilei Wen, and Errui Ding. Im-  
1056       age inpainting with learnable bidirectional attention maps.  
1057       In *The IEEE International Conference on Computer Vision  
(ICCV)*, October 2019. 2
- 1058       [51] Xiaosheng Yan, Feigege Wang, Wenxi Liu, Yuanlong Yu,  
1059       Shengfeng He, and Jia Pan. Visualizing the invisible: Oc-  
1060       cluded vehicle segmentation and recovery. In *IEEE Inter-  
1061       national Conference on Computer Vision*, pages 7618–7627,  
1062       2019. 1, 2, 7
- 1063       [52] Fisher Yu and Vladlen Koltun. Multi-scale context  
1064       aggregation by dilated convolutions. *arXiv preprint  
arXiv:1511.07122*, 2015. 1
- 1065       [53] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and  
1066       Thomas S Huang. Free-form image inpainting with gated  
1067       convolution. In *IEEE International Conference on Computer  
1068       Vision*, pages 4471–4480, 2019. 2, 7
- 1069       [54] Yuhui Yuan, Xilin Chen, and Jingdong Wang. Object-  
1070       contextual representations for semantic segmentation. *arXiv  
preprint arXiv:1909.11065*, 2019. 3
- 1071       [55] Xiaohang Zhan, Xingang Pan, Bo Dai, Ziwei Liu, Dahua  
1072       Lin, and Chen Change Loy. Self-supervised scene de-  
1073       occlusion. In *IEEE conference on computer vision and pat-  
1074       tern recognition (CVPR)*, June 2020. 1, 2, 3, 6, 7, 8
- 1075       [56] Shifeng Zhang, Longyin Wen, Xiao Bian, Zhen Lei, and  
1076       Stan Z Li. Occlusion-aware r-cnn: detecting pedestrians in a  
1077       crowd. In *European Conference on Computer Vision*, pages  
1078       637–653, 2018. 3

- 1080 [57] Ziwei Zhang, Chi Su, Liang Zheng, and Xiaodong Xie. 1134  
1081 Correlating edge, pose with parsing. In *IEEE Conference 1135*  
1082 on Computer Vision and Pattern Recognition, pages 8900– 1136  
1083 8909, 2020. 3 1137
- 1084 [58] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang 1138  
1085 Wang, and Jiaya Jia. Pyramid scene parsing network. In 1139  
1086 *IEEE conference on computer vision and pattern recogni- 1140*  
1087 tion, pages 2881–2890, 2017. 2 1141
- 1088 [59] Chunluan Zhou and Junsong Yuan. Bi-box regression for 1142  
1089 pedestrian detection and occlusion estimation. In *European 1143*  
1090 Conference on Computer Vision, pages 135–151, 2018. 3 1144
- 1091 [60] Yan Zhu, Yuandong Tian, Dimitris Metaxas, and Piotr 1145  
1092 Dollár. Semantic amodal segmentation. In *IEEE Conference 1146*  
1093 on Computer Vision and Pattern Recognition, pages 1464– 1147  
1094 1472, 2017. 1, 2 1148
- 1095 1149
- 1096 1150
- 1097 1151
- 1098 1152
- 1099 1153
- 1100 1154
- 1101 1155
- 1102 1156
- 1103 1157
- 1104 1158
- 1105 1159
- 1106 1160
- 1107 1161
- 1108 1162
- 1109 1163
- 1110 1164
- 1111 1165
- 1112 1166
- 1113 1167
- 1114 1168
- 1115 1169
- 1116 1170
- 1117 1171
- 1118 1172
- 1119 1173
- 1120 1174
- 1121 1175
- 1122 1176
- 1123 1177
- 1124 1178
- 1125 1179
- 1126 1180
- 1127 1181
- 1128 1182
- 1129 1183
- 1130 1184
- 1131 1185
- 1132 1186
- 1133 1187