

Prohibited Items Segmentation via Occlusion-aware Bilayer Modeling

Yunhan Ren^{1,2}, Ruihuang Li¹, Lingbo Liu², Changwen Chen^{1,*}

¹Department of Computing, The Hong Kong Polytechnic University, Hong Kong, China

²Research Institute of Multiple Agents and Embodied Intelligence, Pengcheng Laboratory, Shenzhen, China

yunhan.ren@connect.polyu.hk, csrhl@comp.polyu.edu.hk, liulingbo918@gmail.com, changwen.chen@polyu.edu.hk

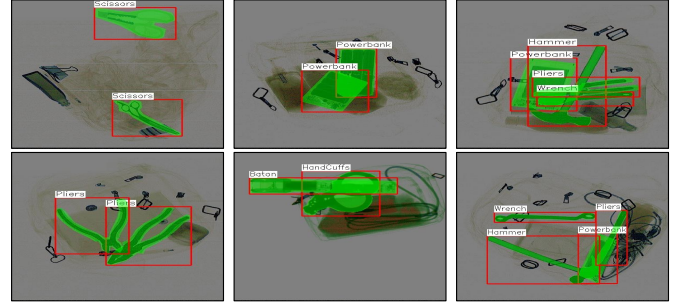
Abstract—Instance segmentation of prohibited items in security X-ray images is a critical yet challenging task. This is mainly caused by the significant appearance gap between prohibited items in X-ray images and natural objects, as well as the severe overlapping among objects in X-ray images. To address these issues, we propose an occlusion-aware instance segmentation pipeline designed to identify prohibited items in X-ray images. Specifically, to bridge the representation gap, we integrate the Segment Anything Model (SAM) into our pipeline, taking advantage of its rich priors and zero-shot generalization capabilities. To address the overlap between prohibited items, we design an occlusion-aware bilayer mask decoder module that explicitly models the occlusion relationships. To supervise occlusion estimation, we manually annotated occlusion areas of prohibited items in two large-scale X-ray image segmentation datasets, PIDray and PIXray. We then reorganized these additional annotations together with the original information as two occlusion-annotated datasets, PIDray-A and PIXray-A. Extensive experimental results on these occlusion-annotated datasets demonstrate the effectiveness of our proposed method. The datasets and codes are available at: <https://github.com/Ryh1218/Occ>.

Index Terms—X-ray inspection, prohibited item segmentation, Segment Anything Model, occlusion handling

I. INTRODUCTION

Security inspection is a critical process in various real-world contexts, such as airports and train stations [1], [2]. Typically, human inspectors are responsible for examining scanned X-ray images generated by security inspection machines to identify potentially prohibited items. With the advancement of deep learning technologies, the computer vision community has made efforts to achieve automatic object detection and segmentation by applying general instance segmentation models [3]–[5].

However, these methods, primarily designed for natural images, face two main challenges when applied to X-ray images: (1) Objects in X-ray images exhibit a significant appearance gap compared to those in natural images, which is primarily caused by different materials absorbing X-rays to varying degrees. (2) Objects in X-ray images often overlap with other items while showing substantial intra-class variations, as illustrated in Fig. 1. These complexities challenge traditional pipelines in understanding the semantic meanings and spatial relationships of prohibited items within X-ray images. Therefore, deep learning frameworks with strong generalization capabilities and specialized techniques for handling



(a) Intra-class variations (b) Overlapping items (c) Complex scenarios

Fig. 1. Unique characteristics of X-ray images where multiple prohibited items with intra-class variations overlap with each other, causing difficulties in distinguishing and segmenting them.

occlusion are essential for accurately detecting and segmenting prohibited items in X-ray images [6], [7].

To this end, we propose an occlusion-aware instance segmentation pipeline specifically designed for segmenting prohibited items in X-ray images. To precisely capture the general representation of X-ray images and fully adapt to the appearance patterns of prohibited items, we incorporate the Segment Anything Model (SAM) into our pipeline, leveraging its rich image priors and powerful generalization capabilities [8]. Recent advances in applying SAM across various research fields have demonstrated its effectiveness in object recognition and generalization. For example, recent work by [9] employs prompt learning methods to harness SAM’s capabilities for instance segmentation tasks in remote sensing applications. Additionally, studies have explored SAM’s application in medical image segmentation, showing promising results despite certain limitations when compared to specialized medical segmentation models [10]. Furthermore, integrating SAM with other advanced techniques, such as generative adversarial networks, has been proposed to enhance its performance on pose estimation [11]. These ongoing research efforts provide a solid foundation for adopting SAM to address specific challenges.

Specifically, inspired by [9], our pipeline employs a frozen SAM Image Encoder as its backbone to generate precise image representations. It then converts the Region of Interest (ROI) features generated by Mask R-CNN into prompts and encodes them using the SAM Prompt Encoder. The SAM Mask Decoder is utilized to regress the mask predictions based on the image representations and prompts. To minimize the

* Corresponding author.

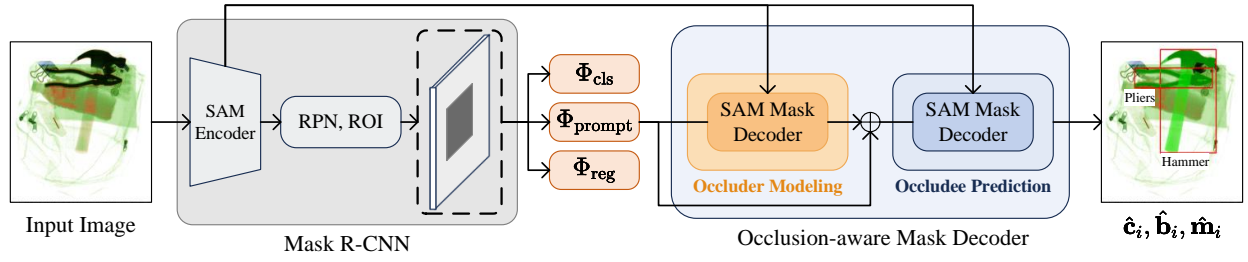


Fig. 2. Demonstration of the proposed pipeline. A frozen SAM Encoder extracts features from input images. Then, RPN and ROI Align predict the normalized ROI features $\hat{\mathbf{v}}_i$. Φ_{cls} , Φ_{reg} and Φ_{prompt} represents the semantic head, localization head, and prompt head. They obtain category predictions, IoU predictions, and sparse prompt embeddings from ROI features, respectively. The proposed occlusion-aware bilayer mask decoder then takes prompt embeddings and image embeddings to segment the target objects with the guidance of their occluding objects.

influence of overlap between prohibited items, inspired by [6], [7], we design an occlusion-aware bilayer mask decoder module that explicitly distinguishes overlapping areas among prohibited items. Specifically, we employ two sequential SAM Mask Decoders with separate training objectives: the first decoder estimates the occlusion area for each target object, while the second regresses the binary mask of the target item guided by the predicted occlusion area. The intermediate optimized image embeddings and prompt embeddings from the first decoder serve as guidance for explicitly modeling occlusion relationships. To effectively train the first decoder, corresponding ground-truth annotations of each occlusion area are required, while none of the existing prohibited instance segmentation X-ray datasets contain such information. Therefore, we introduce two large-scale occlusion-annotated datasets, PIDray-A and PIXray-A, derived from the large-scale X-ray image segmentation datasets PIDray and PIXray [12]–[14]. For each prohibited item, we manually annotate the overlapping areas of other prohibited objects and record the segmentation masks of these areas as additional annotation information. Consequently, the proposed datasets maintain the same number of images as their baseline versions but include additional occlusion annotations to supervise the occlusion-aware mask decoder. To demonstrate the effectiveness of our proposed pipeline, we conduct extensive experiments on the two occlusion-annotated datasets. The experimental results indicate the superiority of our pipeline compared to the state-of-the-art segmentation methods.

In summary, this paper will contribute in the following ways:

- We propose an occlusion-aware instance segmentation pipeline designed for automatic security inspection.
- We design an occlusion-aware bilayer mask decoder module with a guidance mechanism to effectively handle occlusion between prohibited items.
- We create two large-scale security datasets with extra occlusion annotations to meet the need for ground-truth annotations of occlusion areas.

II. RELATED WORK

A. Instance Segmentation

Instance segmentation is concerned with locating objects within images and generating a semantic mask for each indi-

vidual object [15]. Two-stage instance segmentation methods typically consist of a region proposal phase followed by a classification and refinement process. Among two-stage methods, the Mask R-CNN series [16] serve as robust baselines for subsequent advancements. This method first detects the bounding boxes of target objects and then performs semantic segmentation within each ROI area [6]. To enhance feature interaction, PANet [17] integrates a bottom-up pathway grounded in Feature Pyramid Network (FPN) [18]. Meanwhile, Cascade R-CNN [19] employs a sequential architecture of detection heads with different IoUs (Intersection over Union) to generate high-precision bounding boxes. Expanding the cascading strategy, Hybrid Task Cascade [20] introduces a multitask, multistage hybrid structure that enriches contextual information, leading to superior performance.

B. Occlusion Handling

Traditional occlusion handling methods typically involve sophisticated algorithms that exploit symmetry and consistency in stereo matching [21] or incorporate contextual information into random field models for object recognition and segmentation [22]. Advanced methods also include predicting occlusion overlap order by building a scoring histogram [23] or incorporating top-down category-specific reasoning and shape prediction into an energy minimization framework [24]. Additionally, [6], [7] explicitly model occluding and occluded objects separately, using graph convolutional networks to process ROI features of occluding objects and then guide the segmentation of target occluded objects.

C. Prohibited Item Segmentation Benchmarks

Recently, several datasets have been proposed to facilitate research on prohibited item detection. [25] and [26] contain prohibited items with complicated backgrounds and overlapping, but the numbers of images and prohibited items are insufficient. Similarly, [27] introduces a large-scale security inspection dataset named SIXray, but only a small proportion of images (0.84%) contain prohibited items. In response to the need for prohibited item segmentation in real-world scenarios, [12] and [13] present a large-scale benchmark, PIDray, consisting of 124,486 X-ray images, with 47,677 of them containing prohibited items across 12 categories. Additionally, the PIXray dataset [14] is proposed with 15 classes of 15,201

prohibited items, which are all annotated with instance-level masks. These datasets provide robust experimental references for our research.

III. METHODOLOGY

A. Model Architecture

Our proposed framework is based on Mask R-CNN [16], incorporating SAM Image Encoder [8] as the backbone, along with its prompt encoder to transform ROI features generated by Mask R-CNN into SAM-style prompt embedding. An occlusion-aware mask decoder module is proposed for decoding the prompted image features to semantic masks while distinguishing the occlusion relationships between prohibited items. Fig. 2 demonstrates our proposed pipeline. Formally, consider a training dataset $\mathcal{D}_{\text{train}} = (\mathbf{x}_1, \mathcal{Y}_1), (\mathbf{x}_2, \mathcal{Y}_2), \dots, (\mathbf{x}_N, \mathcal{Y}_N)$, where $\mathbf{x}_i \in \mathbb{R}^{H \times W \times 3}$ represents an image while $\mathcal{Y}_i = \{(\mathbf{b}_i, \mathbf{c}_i, \mathbf{m}_i, \mathbf{m}_i^e)\}$ represents corresponding ground-truth information set. Specifically, $\mathbf{b}_i \in \mathbb{R}^{n_i \times 4}$ stands for n_i bounding box annotations of prohibited items in image \mathbf{x}_i , $\mathbf{c}_i \in \mathbb{R}^{n_i \times C}$ represents the corresponding set of semantic labels, $\mathbf{m}_i \in \mathbb{R}^{n_i \times H \times W}$ stands for semantic mask of each prohibited item, and $\mathbf{m}_i^e \in \mathbb{R}^{n_i \times H \times W}$ represents the part of the semantic masks of other annotated object which overlapping the area of \mathbf{b}_i , if it has any. In our pipeline, for an input image \mathbf{x}_i , its image representation feature $\mathbf{f}_i^{\text{img}}$, predicted regional proposals $\hat{\mathbf{o}}_i$, and normalized ROI features $\hat{\mathbf{v}}_i$ are defined as:

$$\begin{aligned} \mathbf{f}_i^{\text{img}} &= \Phi_{\text{enc}}(\mathbf{x}_i) \\ \hat{\mathbf{o}}_i &= \Phi_{\text{rpn}}(\mathbf{f}_i^{\text{img}}) \\ \hat{\mathbf{v}}_i &= \Phi_{\text{roi}}(\mathbf{f}_i^{\text{img}} + \text{PE}, \hat{\mathbf{o}}_i) \end{aligned} \quad (1)$$

, where Φ_{enc} is a frozen SAM image encoder, PE is positional encoding to retain spatial information, Φ_{rpn} and Φ_{roi} are standard RPN and ROI align in Mask R-CNN pipeline.

Predicted ROI features $\hat{\mathbf{v}}_i$ are then processed by the semantic head Φ_{cls} , the localization head Φ_{reg} , and the prompt head Φ_{prompt} to produce the category prediction $\hat{\mathbf{c}}_i$, the bounding box prediction $\hat{\mathbf{b}}_i$ and the SAM-format prompt embedding $\mathbf{f}_i^{\text{sparse}}$.

To accurately estimate the target object mask predictions $\hat{\mathbf{m}}_i$ while avoiding being influenced by its occluding object mask $\hat{\mathbf{m}}_i^e$, we introduce an occlusion-aware bilayer mask decoder architecture, inspired by [6], which is designed to explicitly model the relationships between the occluding object (occluder) and the target object (occludee). The occlusion-aware mask decoder Φ_m explicitly models the relationships between $\hat{\mathbf{m}}_i$ and $\hat{\mathbf{m}}_i^e$ (if any exist):

$$\begin{aligned} \{\hat{\mathbf{c}}_i, \hat{\mathbf{b}}_i, \mathbf{f}_i^{\text{sparse}}\} &= \{\Phi_{\text{cls}}(\hat{\mathbf{v}}_i), \Phi_{\text{reg}}(\hat{\mathbf{v}}_i), \Phi_{\text{prompt}}(\hat{\mathbf{v}}_i)\} \\ \hat{\mathbf{m}}_i, \hat{\mathbf{m}}_i^e &= \Phi_m(\mathbf{f}_i^{\text{img}}, \mathbf{f}_i^{\text{sparse}}) \end{aligned} \quad (2)$$

B. Bilayer Mask Decoder

1) *Module structure*: Fig. 3 illustrates the structure of our proposed bilayer mask decoder, which consists of two sequentially connected SAM mask decoders. The first decoder predicts the mask of the area where the occluder overlaps

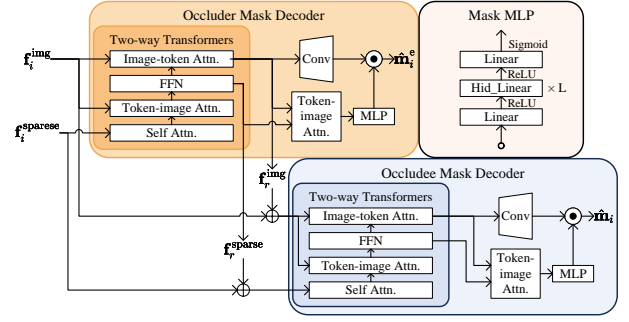


Fig. 3. A detailed implementation of the bilayer occlusion-aware mask decoder structure. Two mask decoders separately predict the masks of the occluder and occludee, providing information about occlusion relationships as one of the learning objectives for the pipeline.

the occludee, and its output guides the subsequent decoder in generating the occludee's prediction.

2) *The occluder decoder*: The occluder decoder Φ_r aims to explicitly model the occlusion between the occluder and occludee using the SAM mask decoder. Formally, it takes the sparse prompt embedding $\mathbf{f}_i^{\text{sparse}}$ and image embedding $\mathbf{f}_i^{\text{img}}$ as inputs to produce the overlapping area prediction of the occluder mask, denoted as $\hat{\mathbf{m}}_i^e$. Simultaneously, benefiting from the SAM mask decoder's involvement, $\mathbf{f}_i^{\text{sparse}}$ and $\mathbf{f}_i^{\text{img}}$ are optimized in terms of global information, with their optimized versions denoted as $\mathbf{f}_r^{\text{sparse}}$ and $\mathbf{f}_r^{\text{img}}$.

3) *The residual guidance mechanism*: Unlike BCNet, which uses only simple GCN-optimized ROI features for guidance, we incorporate information from both image and prompt embeddings. Specifically, we element-wise add these transformer-optimized embeddings ($\mathbf{f}_r^{\text{sparse}}$ and $\mathbf{f}_r^{\text{img}}$) to the original embeddings separately before feeding them into the occludee mask decoder for mask regression. This residual guidance mechanism ensures the occlusion area is explicitly emphasized through cross-attention, helping the model learn occlusion relationships and thereby improving overall performance.

4) *The occludee decoder*: Finally, the occludee decoder Φ_e employs another SAM mask decoder to generate the final occludee prediction, using the refined embeddings $\mathbf{f}_e^{\text{img}}$ and $\mathbf{f}_e^{\text{sparse}}$ as inputs. The above process can be expressed as follows:

$$\begin{aligned} \hat{\mathbf{m}}_i^e, \mathbf{f}_r^{\text{img}}, \mathbf{f}_r^{\text{sparse}} &= \Phi_r(\text{Cat}(\mathbf{t}_{\text{mask}}, \mathbf{t}_{\text{IoU}}, \mathbf{f}_i^{\text{sparse}}), \mathbf{f}_i^{\text{img}}) \\ \mathbf{f}_e^{\text{img}}, \mathbf{f}_e^{\text{sparse}} &= \mathbf{f}_i^{\text{img}} + \mathbf{f}_r^{\text{img}}, \mathbf{f}_i^{\text{sparse}} + \mathbf{f}_r^{\text{sparse}} \\ \hat{\mathbf{m}}_i &= \Phi_e(\text{Cat}(\mathbf{t}_{\text{mask}}, \mathbf{t}_{\text{IoU}}, \mathbf{f}_e^{\text{sparse}}), \mathbf{f}_e^{\text{img}}) \end{aligned} \quad (3)$$

, where \mathbf{t}_{mask} and \mathbf{t}_{IoU} are learnable tokens and are concatenated with prompt embeddings, $\text{Cat}(\cdot)$ represents concatenate operation. To ensure the accuracy of the occluder mask decoder during training, its intermediate occlusion area mask prediction $\hat{\mathbf{m}}_i^e$ is used to compare with ground-truth annotations as an extra learning objective.

C. End-to-end Parameter Learning

Based on the Mask R-CNN architecture [16], instance segmentation can be formed as an end-to-end parameter learning

TABLE I
STATISTICAL DETAILS OF THE PIDRAY-A AND PIXRAY-A DATASETS

Type	Count	PIDray-A				PIXray-A	
		Train	Test			Train	Test
			Easy	Hard	Hide		
Images	Total	76913	24758	9746	13069	3560	1486
	Anno	29454	9482	3733	5005	3560	1486
	Multi	7411	0	3733	3	2891	1193
	Occlu	4785	0	2434	0	1550	654
Annos	Total	39708	9482	8892	5008	10709	4508
	Extra	11080	0	5585	0	4114	1765

problem. This process is guided by a multi-task loss function \mathcal{L} , which simultaneously optimizes the RPN’s region proposal loss \mathcal{L}_{rpn} and the ROI’s object recognition loss. The ROI object recognition loss has three components: classification loss \mathcal{L}_{cls} , bounding box regression loss \mathcal{L}_{reg} , and mask prediction loss \mathcal{L}_{seg} . We adopt \mathcal{L}_{rpn} , \mathcal{L}_{cls} and \mathcal{L}_{reg} from [16].

In our framework, we modify the \mathcal{L}_{seg} with extra occluder mask loss. The original mask prediction loss, denoted as $\mathcal{L}_{\text{pred}}$, is typically computed using the Binary Cross Entropy loss, defined as:

$$\mathcal{L}_{\text{pred}}(\hat{\mathbf{m}}_i, \mathbf{m}_i) = -\frac{1}{w \cdot h} \sum_{j=1}^w \sum_{k=1}^h [\mathbf{m}_{i,j,k} \log(\hat{\mathbf{m}}_{i,j,k}) + (1 - \mathbf{m}_{i,j,k}) \log(1 - \hat{\mathbf{m}}_{i,j,k})] \quad (4)$$

, where $\hat{\mathbf{m}}_i$ and \mathbf{m}_i are the predicted mask and the ground-truth mask, w and h are the width and height of the mask, respectively, and the loss is taken over all pixels (j, k) in the mask.

In our proposed pipeline, an additional objective loss is implemented to let both the occludee mask prediction $\hat{\mathbf{m}}_i$ and occluder mask prediction $\hat{\mathbf{m}}_i^e$ contribute to the overall mask prediction loss. Formally, the segmentation loss \mathcal{L}_{seg} is defined as:

$$\mathcal{L}_{\text{seg}} = \mathcal{L}_{\text{pred}}(\hat{\mathbf{m}}_i, \mathbf{m}_i) + \lambda \cdot \mathcal{L}_{\text{pred}}(\hat{\mathbf{m}}_i^e, \mathbf{m}_i^e) \quad (5)$$

, where λ is a hyperparameter that balances the contributions of the occluder and occludee losses.

The total loss is the weighted sum of the above losses:

$$\mathcal{L} = \frac{1}{M} \sum_i \mathcal{L}_{\text{rpn}}^i + \frac{1}{N} \sum_j (\mathcal{L}_{\text{cls}}^j + \mathbb{1}^j (\mathcal{L}_{\text{reg}}^j + \mathcal{L}_{\text{seg}}^j)) \quad (6)$$

, where $\mathbb{1}$ represents the indicator function that is used to validate positive matches.

IV. EXPERIMENTS

A. Datasets

We primarily evaluate our instance segmentation model on the proposed occlusion-annotated PIDray-A and PIXray-A datasets, which are modified based on the large-scale X-ray image datasets PIDray and PIXray. The PIDray dataset [12], [13] consists of 124,486 X-ray images, including 47,677 images that contain prohibited items in 12 categories. The

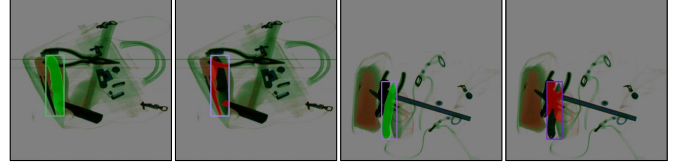


Fig. 4. Occlusion annotations in the proposed PIDray-A dataset. Green masks denote the target object (occludee) mask annotations, while red masks represent the background object (occluder) annotations for each occludee.

test set of PIDray is additionally divided into three subsets: the easy subset contains images with only one exposed prohibited item, the hard subset includes images with multiple prohibited items, and the hidden subset contains images with one prohibited object that is deliberately concealed to test the model’s ability under extreme circumstances. As for the PIXray dataset [14], it contains 5,046 X-ray images, all of which have one or more prohibited items. In total, the PIXray dataset includes 15,201 prohibited items with bounding boxes and semantic mask annotations across 15 object categories. In our experiment, we use an 8:2 train-test split for PIXray.

To optimize the occluder mask decoder, ground-truth occluder annotations are necessary for calculating \mathcal{L}_{seg} . Consequently, we add additional occlusion annotations for contraband items in PIDray and PIXray, reorganizing them into two occlusion-annotated datasets. Specifically, we first identify items whose masks have at least 5% of their area covered by other objects’ bounding boxes and mark them as occludees. Then, we verify whether these occludees’ masks intersect with their occluders’ masks, recording the intersecting mask area as occlusion annotations. For cases without intersections, we discard them as non-occluding instances. Fig. 4 demonstrates the example visualization of occlusion annotations on representative X-ray images, indicating the extra emphasis of our pipeline on occluding areas. Statistic details of the proposed PIDray-A and PIXray-A datasets are presented in Tab. I. Specifically, the ‘Anno,’ ‘Multi,’ and ‘Occlu’ rows indicate whether the images contain prohibited items, multiple prohibited items, or occluded items, respectively. The ‘Extra’ row indicates the number of prohibited items with additional occlusion annotations. It can be concluded that more than half of images with more than one prohibited item suffer from overlapping of prohibited items, indicating that an occlusion-aware segmentation pipeline is highly demanded.

B. Implementation details

In line with the commonly used COCO metrics, we employ mean average precision (mAP) to evaluate our model’s performance. Specifically, we focus on mAP at three IoU thresholds: $T = 0.5$, $T = 0.75$, and $\{T\} = (0.50 : 0.05 : 0.95)$, for both bounding boxes and masks, denoted as AP_b and AP_m . All input images are resized to a resolution of 512×512 to generate visual patches for the SAM Image Encoder. To balance computational cost and model performance, we adopt ViT-Large as the backbone of SAM. The hyperparameter λ for segmentation loss calculation is set to 0.25. For optimization, we use the AdamW optimizer with an initial learning rate of

TABLE II
EXPERIMENTAL RESULTS OF OBJECT DETECTION AND INSTANCE SEGMENTATION ON THE PIDRAY-A AND PIXRAY-A DATASETS

Model	PIDray-A						PIXray-A					
	AP_b	AP_b^{50}	AP_b^{75}	AP_m	AP_m^{50}	AP_m^{75}	AP_b	AP_b^{50}	AP_b^{75}	AP_m	AP_m^{50}	AP_m^{75}
Mask R-CNN	58.7	76.2	67.2	50.2	73.4	58.2	65.2	91.1	75.3	56.1	86.7	61.7
Cascade MR-CNN	65.0	78.1	72.1	53.6	76.0	62.4	71.8	91.2	81.0	57.1	86.2	63.1
PIDrayNet	66.6	81.7	74.3	53.4	78.6	61.4	72.4	91.6	81.3	57.7	87.1	64.1
RSPrompter	66.0	83.1	74.3	56.5	81.4	64.4	72.0	94.3	83.1	58.9	90.1	66.4
Our Model	67.5	84.3	75.8	57.6	82.6	65.8	72.0	94.4	83.7	58.9	90.6	66.2

$1e^{-6}$. After a warm-up phase of 50 iterations, during which the learning rate linearly increases to $1e^{-4}$, we apply a Cosine Annealing scheduler [28] to gradually decay it back to $1e^{-6}$ over the training period, ensuring stable training. To augment the training data, we employ horizontal flipping and random jittering, implemented via the MMDetection framework. All experiments are conducted on a single NVIDIA A800 Tensor Core GPU with 80GB VRAM, using a batch size of 64. We train the model for 50 epochs on the PIDray-A dataset and 60 epochs on the PIXray-A dataset.

C. Experiment Results

We quantitatively compare our proposed methods with other state-of-the-art methods. Specifically, we compare our models with Mask R-CNN [16], Cascade Mask R-CNN [19], SDANet [13], and the PIDray baseline [12] (namely PIDrayNet). Moreover, we run RSPrompter [9] with the same backbone to indicate the impact of the proposed occlusion-aware bilayer mask decoder. The results of these comparisons are detailed in Tab. II, where our proposed model achieves the optimal results on 10 out of 12 metrics while obtaining suboptimal results on the other 2 metrics. Specifically, on the PIDray-A dataset, our proposed model achieves 0.9% improvements compared with suboptimal results on both bounding box mAP and semantic mask mAP. Regarding the PIXray-A dataset, our model obtains either optimal or suboptimal results. These results demonstrate the effectiveness of our proposed model on automatic instance segmentation of prohibited items.

We also demonstrate the results on subsets of the PIDray-A dataset, as shown in Tab. III. Our proposed model performs the best results on 4 out of 6 subsets. However, our model shows no clear advantage in detecting and segmenting hidden objects. We assume that this is primarily because the PIDrayNet uses specialized X-ray image optimizing tricks for hidden object filtering, while our model relies on a simple RPN for detection.

In addition, we provide a qualitative comparison between PIDrayNet, RSPrompter, and our proposed pipeline in Fig. 5. When distinguishing two overlapping power banks, our pipeline is the only model that can identify the occlusion and successfully segment the two objects separately. Additionally, ours is the only model that correctly predicts the mask of the covered hammer handle. These results indicate the effectiveness of our pipeline in distinguishing occlusions among prohibited items.

TABLE III
EXPERIMENTAL RESULTS ON THE SUBSETS OF PIDRAY-A DATASET

Method	AP_b			AP_m		
	Easy	Hard	Hide	Easy	Hard	Hide
Mask R-CNN	66.2	58.6	43.8	59.2	50.1	35.5
Cascade MR-CNN	71.9	63.2	46.8	60.7	52.0	36.2
SDANet	72.5	63.7	48.0	61.1	51.7	37.0
PIDrayNet	74.5	64.8	53.0	61.4	51.9	39.7
RSPrompter	74.1	67.2	45.3	65.7	58.8	33.2
Our Model	75.7	68.3	47.1	67.1	60.0	34.3

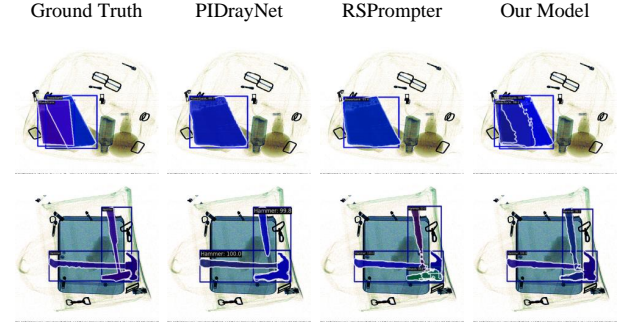


Fig. 5. Qualitative results on the PIDray-A dataset to demonstrate the model's ability in segmenting occluded prohibited items.

D. Ablation Study

1) *Impact on occluded prohibited items*: To demonstrate the effect of the proposed occlusion handling module, we further divide the training set of the PIDray-A into a smaller training set, the validation subset, and the occlusion subset. The validation subset contains images and annotations without occlusion, while the occlusion subset contains only occluding images and annotations. For simplicity, we discard X-ray images with no prohibited items. The approximate proportions of the split training set, validation subset, and occlusion subset are (0.78 : 0.12 : 0.1).

We train our proposed pipeline, RSPrompter, and PIDrayNet on the split training set and test their performances on the validation subset and the occlusion subset to show their effectiveness in dealing with common and occluded prohibited items separately. As shown in Tab. IV, our proposed pipeline outperforms all other models on the validation subset. Compared to PIDrayNet and RSPrompter, our model achieves 1.5% and 0.1% enhancements in bounding box precision, as well as 3.2% and 0.1% improvements in semantic mask precision. For the occlusion subset, our proposed pipeline achieves optimal

TABLE IV
ABLATION STUDY ON OCCLUDING AND NON-OCCLUSION SUBSETS

Set	Model	AP _b	AP _b ⁵⁰	AP _b ⁷⁵	AP _m	AP _m ⁵⁰	AP _m ⁷⁵
Val	PIDNet	77.2	88.1	84.8	67.3	87.5	80.0
	RSP	78.6	91.9	88.1	70.6	91.1	82.8
	Ours	78.7	92.0	88.2	70.7	91.4	83.1
Occ	PIDNet	49.8	64.0	55.2	37.8	60.7	40.3
	RSP	46.7	67.7	54.4	39.4	63.0	42.5
	Ours	47.5	67.9	55.4	40.3	64.1	44.1

TABLE V
ABLATION STUDY OF THE OCCLUDER MASK DECODER AND RESIDUAL FEATURE GUIDANCE

Occluder	Guide	AP _b	AP _b ⁵⁰	AP _b ⁷⁵	AP _m	AP _m ⁵⁰	AP _m ⁷⁵
✗	✗	66.1	83.2	74.2	56.6	81.3	64.7
✓	✗	66.6	83.5	74.9	56.9	81.8	64.8
✓	✓	67.5	84.3	75.8	57.6	82.6	65.8

performances on 5 out of 6 metrics, highlighting the significant potential of our proposed bilayer mask decoder structure and occluder guidance operation.

2) *Impact of the occlusion-aware bilayer mask decoder module:* We further examine the impact of the occlusion-aware bilayer mask decoder module in our proposed pipeline. We first evaluate the performances with and without the residual feature addition strategy and also explore the performance when the occluder mask decoder is removed entirely, resulting in a simple SAM Mask Decoder. Tab. V demonstrates the results, where the ‘Occluder’ column represents whether the occluder mask decoder is added, and the ‘Guide’ column indicates whether the residual feature guidance is applied. The results indicate that the occluder mask decoder and residual feature guidance operation together improve the overall performance of our proposed pipeline by more than 1% on every metric, indicating the effectiveness of our proposed module.

V. CONCLUSION

In this paper, we first highlight two critical differences between prohibited items in X-ray images and natural objects: the appearance gap and the severe overlapping. To address these challenges, we present an occlusion-aware instance segmentation pipeline based on Mask R-CNN. Specifically, to generalize the priors of natural objects to prohibited X-ray items, we leverage the extensive pre-trained knowledge and strong zero-shot generalization capability of the Segment Anything Model. For handling occlusion, we propose an occlusion-aware bilayer mask decoder that explicitly models the occlusion relationships. Specifically, it uses one decoder to estimate the occlusion areas directly and adopts another standalone decoder to segment the target object guided by the estimated occlusion areas. To supervise occlusion area estimation, we introduce two large-scale X-ray datasets, PIDray-A and PIXray-A, which are annotated with additional ground-truth occlusion information. Experimental results on these datasets demonstrate the effectiveness of our pipeline.

ACKNOWLEDGMENTS

This research has been supported by the National Natural Science Foundation of China under Grant No. 62306258, the Major Key Project of PCL (Grant No. PCL2024A04), and the Hong Kong Polytechnic University start-up fund ZVVK.

REFERENCES

- [1] Samet Akcay and Toby Breckon, “Towards automatic threat detection: A survey of advances of deep learning within x-ray security imaging,” *Pattern Recognition*, vol. 122, pp. 108245, 2022.
- [2] Domingo Mery, Daniel Saavedra, et al., “X-ray baggage inspection with computer vision: A survey,” *IEEE Access*, vol. 8, pp. 145620–145633, 2020.
- [3] Ruyi Ji, Longyin Wen, et al., “Attention convolutional binary neural tree for fine-grained visual categorization,” in *CVPR*, 2020.
- [4] Zhi Tian, Chunhua Shen, et al., “FCOS: Fully convolutional one-stage object detection,” in *CVPR*, 2019.
- [5] Congcong Li, Dawei Du, et al., “Spatial attention pyramid network for unsupervised domain adaptation,” in *ECCV*, 2020.
- [6] Lei Ke, Yu-Wing Tai, et al., “Deep occlusion-aware instance segmentation with overlapping BiLayers,” in *CVPR*, 2021.
- [7] Lei Ke, Yu-Wing Tai, et al., “Occlusion-aware instance segmentation via BiLayer network architectures,” *IEEE TPAMI*, vol. 45, no. 8, pp. 10197–10211, 2023.
- [8] Alexander Kirillov, Eric Mintun, et al., “Segment anything,” in *ICCV*, 2023.
- [9] Keyan Chen, Chenyang Liu, et al., “Rsprompter: Learning to prompt for remote sensing instance segmentation based on visual foundation model,” *IEEE TGRS*, 2024.
- [10] Jun Ma, Yuting He, et al., “Segment anything in medical images,” *Nature Communications*, vol. 15, pp. 654, 2024.
- [11] Jiehong Lin, Lihua Liu, et al., “Sam-6d: Segment anything model meets zero-shot 6d object pose estimation,” in *CVPR*, 2024.
- [12] Libo Zhang, Luta Jiang, et al., “PIDray: A large-scale x-ray benchmark for real-world prohibited item detection,” *IJCV*, vol. 131, no. 12, pp. 3170–3192, 2023.
- [13] Boying Wang, Libo Zhang, et al., “Towards real-world prohibited item detection: A large-scale x-ray benchmark,” in *ICCV*, 2021.
- [14] Bowen Ma, Tong Jia, et al., “Automated segmentation of prohibited items in x-ray baggage images using dense de-overlap attention snake,” *IEEE TMM*, vol. 25, pp. 4374–4386, 2022.
- [15] Shervin Minaee, Yuri Boykov, et al., “Image segmentation using deep learning: A survey,” *IEEE TPAMI*, vol. 44, no. 7, pp. 3523–3542, 2022.
- [16] Kaiming He, Georgia Gkioxari, et al., “Mask r-cnn,” in *ICCV*, 2017.
- [17] Shu Liu, Lu Qi, et al., “Path aggregation network for instance segmentation,” in *CVPR*, 2018.
- [18] Tsung-Yi Lin, Piotr Dollár, et al., “Feature pyramid networks for object detection,” in *CVPR*, 2017.
- [19] Zhaowei Cai and Nuno Vasconcelos, “Cascade r-cnn: Delving into high quality object detection,” in *CVPR*, 2018.
- [20] Kai Chen, Jiangmiao Pang, et al., “Hybrid task cascade for instance segmentation,” in *CVPR*, 2019.
- [21] Jian Sun, Yin Li, et al., “Symmetric stereo matching for occlusion handling,” in *CVPR*, 2005.
- [22] John Winn and Jamie Shotton, “The layout consistent random field for recognizing and segmenting partially occluded objects,” in *CVPR*, 2006.
- [23] Joseph Tighe, Marc Niethammer, et al., “Scene parsing with object instances and occlusion ordering,” in *CVPR*, 2014.
- [24] Yi-Ting Chen, Xiaokai Liu, et al., “Multi-instance object segmentation with occlusion handling,” in *CVPR*, 2015.
- [25] Samet Akcay and Toby P. Breckon, “An evaluation of region based object detection strategies within x-ray baggage security imagery,” in *ICIP*, 2017.
- [26] Yanlu Wei, Renshui Tao, et al., “Occluded prohibited items detection: An x-ray security inspection benchmark and de-occlusion attention module,” in *ACM MM*, 2020.
- [27] Caijing Miao, Lingxi Xie, et al., “Sixray: A large-scale security inspection x-ray benchmark for prohibited item discovery in overlapping images,” in *CVPR*, 2019.
- [28] Ilya Loshchilov and Frank Hutter, “Sgdr: Stochastic gradient descent with warm restarts,” in *ICLR*, 2022.