

# Segmentation Mask Guided End-to-End Person Search

Dingyuan Zheng<sup>a</sup>, Jimin Xiao<sup>a,\*</sup>, Kaizhu Huang<sup>a</sup>, Yao Zhao<sup>b</sup>

<sup>a</sup>*Department of Electrical and Electronic Engineering, Xi'an Jiaotong-Liverpool University, Suzhou, China*

<sup>b</sup>*Institute of Information Science, Beijing Jiaotong University, Beijing, China*

---

## Abstract

Person search aims to search for a target person among multiple images recorded by multiple surveillance cameras, which faces various challenges from both pedestrian detection and person re-identification. Besides the large intra-class variations owing to various illumination conditions, occlusions and varying poses, background clutters in the detected pedestrian bounding boxes further deteriorate the extracted features for each person, making them less discriminative. To tackle these problems, we develop a novel approach which guides the network with segmentation masks so that discriminative features can be learned invariant to the background clutters. We demonstrate that joint optimization of pedestrian detection, person re-identification and pedestrian segmentation enables to produce more discriminative features for pedestrian, and consequently leads to better person search performance. Extensive experiments on benchmark dataset CUHK-SYSU, show that our proposed model achieves the state-of-the-art performance with 86.3% mAP and 86.5% top-1 accuracy respectively.

**Keywords:** person search, re-identification, pedestrian detection, segmentation masks, background clutters

---



---

\*Corresponding author

*Email address:* [jimin.xiao@xjtlu.edu.cn](mailto:jimin.xiao@xjtlu.edu.cn) (Jimin Xiao)

<sup>1</sup>D. Zheng, J. Xiao, K. Huang are with the Department of Electrical and Electronic Engineering, Xian Jiaotong-Liverpool University, Suzhou, China (e-mail: dingyuan.zheng, jimin.xiao, kaizhu.huang@xjtlu.edu.cn).

<sup>2</sup>Y. Zhao is with Institute of Information Science, Beijing Jiaotong University, Beijing, China (e-mail: yzhao@bjtu.edu.cn).

## 1. Introduction

Person re-identification has been widely applied in video surveillance systems with increasing demands in urban safety. It has attracted great attention in the computer vision community during the last decade. Person re-identification is generally solved as a retrieval problem [1, 2]. Given a probe image, person re-identification aims to find all the images in the gallery set with the same identity. However, person re-identification has not been fully addressed, since the images captured by cameras are usually with the characteristics of illumination variations, occlusions and low resolution owing to the shooting environment. These challenges potentially increase the intra-class variations and raise the recognition difficulty.

To this end, a great deal of research works on person re-identification devote to extract more discriminative features to represent human individuals, either by hand-crafted features [3, 4] or by CNN features [5, 6]. Most of the existing person re-identification methods engage on cropped pedestrian bounding boxes without considering background clutters. Specifically, human individuals are represented by the features extracted from the regions constrained with the detected pedestrian bounding boxes, and Euclidean distance is computed to evaluate the similarity level among those probe-gallery pairs. It may result in a situation that different persons with similar background are close in the learned feature space. For example, in Fig. 1, the person in the bounding box in the third figure is different from the probe image. However, it is ranked before the person in the bounding box of the fourth figure who has the same identification as the probe image; this is simply because its background is more similar with the probe image.

One straightforward yet effective solution to tackle the problem is to make the foreground part such as human body as the dominant region for feature extraction. In [7], it adopts pose estimation approach to locate the key body points, and then aggregates the local features extracted from the pre-defined body regions with the global features extracted from the whole image. Based



Figure 1: One example to show the negative effect of background clutters which deteriorate the person re-identification performance. Blue box in the first column is the probe image. Other columns are the searching results from rank-1 to rank-3. Green boxes indicate the correct searching result, while the red box indicates the incorrect searching result.

on the similar ideas, Tian *et al.* propose a person-region guided pooling network with the assist of human parsing maps to solve the background bias problem [8]. Recently, researchers also attempt to introduce the attention mechanism into the person re-identification task for pedestrian feature extraction [9, 10, 11].

Person search aims to search for the targeting person among multiple images recorded with different surveillance cameras, where the pedestrian bounding boxes are not available. Person search, different from person re-identification which assumes most of the pedestrian bounding boxes are manually cropped or perfectly detected by the state-of-the-art detectors, i.e. Faster R-CNN [12], handles the challenges from both pedestrian detection and re-identification. Specifically, considering the step of pedestrian detection, the misalignment and false alarm caused by detectors further decrease the recognition rate [13, 14]. Meanwhile person search also has the aforementioned problem resulting from the background clutters in the generated pedestrian bounding boxes.

In a recent work, Chen *et al.* adopt segmentation mask to solve the background clutter problem in the person search task [15]. Specifically, a two-stream model is established to extract the pedestrian features with one stream to emphasize the foreground information for the regions covered by the segmentation mask, and second stream to retain the global information for the original image. However, in [15] the foreground regions are heuristically fixed annotation, in other words, to what extent the background should be removed is decided

by the pedestrian segmentation mask. Besides, this work separates the steps of pedestrian masking, pedestrian detection and person re-identification, which ignores the fact that jointly optimizing these steps can further bring in performance gain.

Inspired by the previous works [16, 17, 18] that solve the person search task in an end-to-end manner, we propose an novel end-to-end person search framework that uses the segmentation mask to mitigate the negative effect of background clutters. Different from the previous work [15] that designates the foreground regions by the segmentation masks explicitly, we utilize the segmentation mask to guide the feature extraction network to learn the enriched foreground features through a parallel mask branch. To do this, segmentation masks are precisely labeled in our new created dataset. Our proposed person search approach jointly optimizes pedestrian detection, person re-identification and pedestrian segmentation, which obtains more discriminative features for pedestrians benefiting from end-to-end learning.

We summarize our contributions are as follows.

- We propose a segmentation masks guided person search framework so as to mitigate the negative effect of the background clutters in the detected pedestrian bounding boxes. Our proposed person search framework is trained end-to-end which considers the inherent relations among pedestrian detection, person re-identification and pedestrian segmentation, and hence more discriminative features for pedestrians can be learned, which effectively enhance the person search performance.
- We create a new dataset which contains precise pedestrian segmentation mask annotations for 1,833 images from the existing CUHK-SYSU dataset. The dataset will be released for the future segmentation mask based person search research, which can be downloaded from the link: <https://github.com/Dingyuan-Zheng/maskPS>. Meanwhile, it is found that our approach only requires partial annotations for the segmentation masks rather than that for the whole dataset.

- Extensive experiments on benchmark dataset CUHK-SYSU show that our proposed segmentation masks guided end-to-end person search framework outperforms a wide range of state-of-the-art person search methods, obtaining 86.3% mAP and 86.5% top-1 accuracy, respectively.

## 2. Related Work

In this section, we first review the existing works for the two sub-tasks in person search, pedestrian detection and person re-identification respectively. We then review the recent achievements on person search.

### 2.1. Pedestrian Detection

Pedestrian detection has witnessed significant improvement in the past few decades. The first landmark work achieved by Dalal *et al.*[19] adopts the architecture of HOG+SVM, and then DPM [20] is developed to better address the occlusion issue. After that, ICF [21] and its variants [22, 23] outperform the previous hand-crafted feature based pedestrian detection methods. More recently, great progress has been made on the realm of general object detection benefiting from the convolutional neural networks [24, 25, 26, 27, 28, 12]. Further, [29] discussed the feasibility of Faster R-CNN on pedestrian detection task. In this paper, we also adopt Faster R-CNN as our pedestrian detector.

### 2.2. Person Re-identification

With the great success of convolutional neural networks, researchers have proposed numerous deep learning based person re-identification solutions [30, 31, 32, 33, 34]. The re-identification system is typically composed of two categories, feature extraction and similarity metrics learning. Some researchers attempt to improve the person re-identification performance by taking the advantage of enhanced feature representation. For instance, in [35, 36], the original image is horizontally split into patches, and part matching is then applied among these generated local patches. In [7], local features of the body sub-regions defined by the pose estimation results are merged with the whole body

features to improve the robustness of the final feature representation. Other researchers propose better person re-identification solutions by using well-designed similarity metrics learning. Generally, one category adopts verification loss, for example, contrastive loss [32], triplet loss [37] or quadruplet loss [38], while another category utilizes identification loss [30, 39] or both [40]. In this paper, our person search framework is built upon the identification model.

### 2.3. Person Search

As an extension of the conventional person re-identification, person search retrieves the target person from the raw scene images, where pedestrian bounding boxes are not available [41]. In the pioneer work [17], Xiao *et al.* show that pedestrian detection and person identification could be solved in an end-to-end framework. Following this work, Xiao *et al.* [18] enhance the discriminability of the pedestrian features by introducing center loss. Liu *et al.* [42] recursively shrink the attentive regions till the target person is retrieved. In [43], global context of query-gallery pairs are emphasized by establishing a query-guided region proposal network and similarity sub-network in a siamese structure. Yan *et al.* further improve the person search performance by exploiting co-travelers as global context. A recent person search approach [15] uses segmentation mask to filter the foreground person from the original input and aggregates the features of both foreground and whole image which are extracted from a two-stream model. Besides, the authors also state that better person search performance can be achieved by solving pedestrian detection and identification separately with off-line pedestrian masks. Different from [15], we optimize jointly these three tasks in an end-to-end framework. In particular, we use the segmentation mask to guide the network to learn the discriminative regions automatically rather than explicitly specifying these regions.

## 3. Proposed Method

In this section, we propose a novel partially labeled segmentation masks guided person search framework, as shown in Fig. 3. We first introduce our

new dataset which contains partially labeled segmentation masks, and then elaborate our end-to-end person search framework.

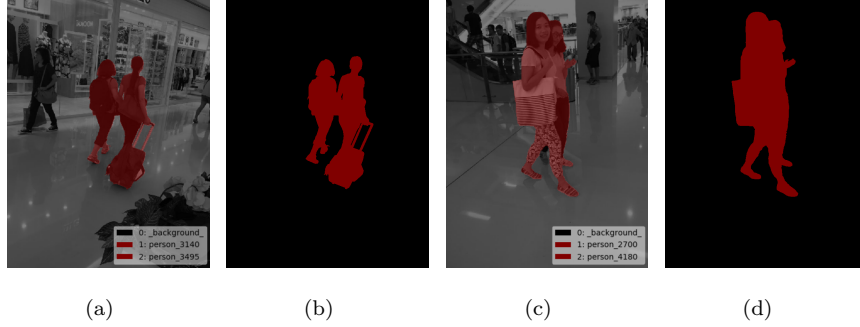


Figure 2: Two examples in our created dataset. We only provide the segmentation masks for the labeled persons (the persons labeled with [1-5532] in CUHK-SYSU dataset). The shadow regions in first and third columns indicate the labeled persons. The second and fourth columns are their segmentation masks.

### 3.1. A New Dataset with Partially Labeled Segmentation Masks

To the best of our knowledge, all current segmentation masks based person search/re-ID approaches are based on off-line pedestrian masks generated from Fully Convolutional Networks (FCN) [44] or Fully Convolutional Instance-aware Semantic Segmentation (FCIS) [45] without considering the benefit from joint optimization of pedestrian segmentation and person re-identification tasks. To extract more discriminative features and mitigate the negative effect of background clutters as well as to build an end-to-end framework that jointly optimize pedestrian detection, person identification and pedestrian segmentation, we created a new person search dataset to provide the precise annotations of pedestrian segmentation masks. We labeled the pedestrian segmentation masks for a portion of images in CUHK-SYSU dataset [16].

CUHK-SYSU dataset [16] is a large-scale dataset for person search, and the data is collected from diverse scenes. Specifically, it contains 18,184 images, 6,057 query persons in 12,490 images are captured from the street, while the rest 2,375 query persons in 5,694 images are collected from the movies and

dramas. The dataset is split into the training set and test set, and guarantees no overlap occurs on images and query persons between the training split and test split. Training split contains 11,206 images with 5,532 query persons, while test split contains 6,978 images with 2,900 query persons. Each query person appears in at least two images. The dataset also contains two subsets to evaluate the person search performance under low resolution and occlusion. The person identities of training split are in the range of  $[-1, 5532]$ , where -1 indicates the unlabeled persons and 0 indicates the background (non-person).

In our dataset, to guarantee the uniformity of data distribution, we divide the training set into  $N$  portions ( $N=2,241$  in our case), with 5 images in each portion (except the last portion, which contains 6 images), and we randomly select one image from each portion, and filter out the images with only unlabeled persons (the persons labeled with -1). Finally 1,833 images in the training set are selected for the segmentation mask labeling.

To the best of our knowledge, accessories, i.e, handbags, luggage cases and baby carriage, might act as suggestive context in person re-identification. In a consequence, we treat these objects as foreground during the mask annotating process. It should be noticed that we provide the mask annotations for only the labeled persons (the persons labeled with  $[1\sim 5532]$ ) in a raw scene image. The samples of the image with segmentation masks are shown in Fig. 2.

We utilize Labelme [46] as the annotation tool. All our segmentation masks follow the unified annotation rules. When a person is occluded by non-person objects, we only keep the visible part of the occluded person, and the accessories are kept as well. We also give the statistics for our created dataset, as shown in Table 1. The selected 1,833 images from the CHUK-SYSU training set contain 9,084 pedestrians in total, with 2,815 labeled persons and 6,269 unlabeled persons. In particular, we only annotated the segmentation masks for the labeled persons. The rest 9,373 images in the training set contains 12,270 labeled persons and 33,918 unlabeled persons.



Dataset	Number of images	Number of pedestrians	
Images with masks	1,833	LP	<b>2,815</b>
		UP	6,269
Images without masks	9,373	LP	12,270
		UP	33,918

Table 1: Statistics of our created dataset. LP: Labeled persons, UP: Unlabeled persons. The labeled persons (2,815) in the selected 1,833 images are annotated with pedestrian segmentation masks.

### 3.2. Our Proposed Person Search Framework

Person search aims to retrieve the target person across raw scene images without pedestrian bounding boxes. Our proposed approach jointly optimizes three sub-tasks including pedestrian detection, person identification and pedestrian segmentation in an end-to-end person search framework. Apart from the pedestrian detection module to produce online pedestrian bounding boxes and person identification module to categorize person identities, we further establish a parallel pedestrian segmentation branch to predict pedestrian masks. Benefiting from the end-to-end optimization of three tasks, more discriminative pedestrian features can be extracted. The overall schematic of the proposed segmentation masks guided end-to-end person search framework is shown in Fig. 3. The network is elaborated as follows.

Arbitrary size images are resized such that the shorter side has 600 pixels. An image is then fed into the first part of the residual backbone network [47]. Specifically, we divide the residual network into two parts, i.e, for ResNet-50, the first part contains the layers from Conv1 to Res4, and the rest Res5 forms the second part.

To address pedestrian detection, we adopt the region proposal network [12] (RPN) to produce online pedestrian proposals. RPN is trained with cross en-

tropy loss to distinguish pedestrians and background, we express it as  $\mathcal{L}_{cls}$ :

$$\mathcal{L}_{cls} = - \sum_{i=1}^N y_i \log(s_i), \quad (1)$$

where  $N$  is the number of the generated proposals,  $s_i$  is the prediction score and  $y_i$  is the related ground truth which indicates person or non-person. We use the Smoothed-L1 loss [28],  $\mathcal{L}_{reg}$ , to regress the precise location for each pedestrian, as defined as follows:

$$\mathcal{L}_{reg} = \begin{cases} 0.5\mathcal{D}^2 & |\mathcal{D}| < 1 \\ |\mathcal{D}| - 0.5 & otherwise, \end{cases} \quad (2)$$

where  $\mathcal{D}$  denotes the coordinate differences between the predicted box and its related ground truth location, and these two losses together are denoted as  $\mathcal{L}_{RPN}$ . The generated candidate boxes are either associated with background or a foreground part (the ground truth bounding boxes). Since we only provide the mask annotation for the labeled persons (persons labeled with [1~5532]) in a raw scene image, those generated candidate boxes associated with foreground parts are consequently divided into two types. The first is the candidate boxes associated with labeled persons, and it is denoted as proposals with mask, while the second is the candidate boxes associated with unlabeled persons, which we denote as proposals without mask, as shown in Fig. 3.

All the proposals generated from RPN and the feature maps generated from the first part of the residual network are input into the ROIAlign layer [48] to produce the fixed size feature map for each ROI.

Targeting for person identification, once the fixed size feature maps are obtained, these feature maps are further convolved into the second part of the residual network and the output,  $\mathcal{F}_p \in \mathcal{R}^{c \times m \times m}$ , are summarized into 2,048 dimensional feature vectors  $f_p \in \mathcal{R}^c$  through an average pooling layer. Here  $c$  is the channel width and  $m$  denotes the size of the feature maps. To further reduce the false alarm caused by RPN and refine the predicted locations of the candidate pedestrians,  $f_p$  are then fed into two fully connected layers

respectively and again supervised by  $\mathcal{L}_{reg}$  and  $\mathcal{L}_{cls}$  losses. Following [17], we denoted these two losses together as  $\mathcal{L}_{RCNN}$ . Besides,  $f_p$  is projected into a 256 dimensional feature vector  $f_{id} \in \mathcal{R}^d$  through the third fully connected layer followed by L2-normalization, which is used as the final feature representation for each retrieved pedestrian. In the training phase, we adopt OIM loss [17] to supervise the person identification module, where  $p_{id}$  indicates the probability of the identification features,  $f_{id}$ , belonging to  $id$ -th class,

$$p_{id} = \frac{\exp(v_{id}^T f_{id} / \tau)}{\sum_{j=1}^L \exp(v_j^T f_{id} / \tau) + \sum_{k=1}^Q \exp(u_k^T f_{id} / \tau)}. \quad (3)$$

Here  $\tau$  is a parameter to control the softness of the probability function. The features of the labeled identities are stored in a lookup table with dimension  $L$ , with  $v_{id}$  denoting the current feature for class  $id$  among 5,532 categories, and it is continuously updated during the training phase as follows:

$$v_{id} = \beta v_{id} + (1 - \beta) f_{id}, \quad (4)$$

where  $\beta$  is a momentum parameter used to adjust the update rate. While the features of the unlabeled persons are stored in a circular queue with dimension  $Q$ , and  $u_k$  indicates the features for the  $k$ -th unlabeled person. The objective of OIM loss is to maximize the expected log-likelihood, and the identification loss is then defined as:

$$\mathcal{L}_{identification} = E_x[\log p_{id}]. \quad (5)$$

Most importantly, in order to improve the discriminability of  $f_{id}$ , we establish a parallel mask branch on top of the shared features  $\mathcal{F}_p$ . Specifically, we pick out the feature maps associated with labeled persons from  $\mathcal{F}_p$ , and use these feature maps  $\mathcal{F}_{pm} \in \mathcal{R}^{c \times m \times m}$  to predict segmentation masks with the size of  $2m \times 2m$  for each proposal with mask ( $m$  equals 7 in our case), and the predicted masks

are then computed into binary cross entropy loss [48], which can be written as:

$$\mathcal{L}_{mask} = \sum_{s=1}^S \sum_{t=1}^T (-Y_t \log(x_t) - (1 - Y_t) \log(1 - x_t)), \quad (6)$$

where  $x_t$  denotes the probability of  $t$ -th pixel in the predicted mask being recognized as foreground,  $Y_t$  is its related ground truth,  $T$  is the number of pixels in the predicted pedestrian mask ( $T = 2m \times 2m$ ), and  $S$  is the number of proposals with mask.

Finally, we adopt a multi-task loss to train our person search framework in an end-to-end manner. The total loss is defined as:

$$\mathcal{L}_{total} = \mathcal{L}_{RPN} + \mathcal{L}_{RCNN} + \mathcal{L}_{identification} + \lambda \mathcal{L}_{mask}, \quad (7)$$

specifically,  $\lambda = 1$  when input image contains labeled segmentation masks, otherwise,  $\lambda = 0$ .

With the assist of the partially labeled segmentation masks, our proposed person search framework can generate more discriminative features invariant to background clutters, compared with the previous segmentation mask based state-of-the-art approach [15].

## 4. Experimental Results

In this section, the dataset and evaluation metrics we used are first introduced, followed by implementation details. We also compare our proposed method with previous state-of-the-art results. At last, our proposed person search framework is verified in the ablation study.

### 4.1. Dataset and Evaluation Metrics

We use our newly labeled CHUK-SYSU dataset, as introduced in Sec. 3.1, in our experiments. We adopt both mean average precision (mAP) and top-1 matching rate to evaluate all our experiment performance, similar to [17]. A matching is accepted only if the overlap between the detected pedestrian

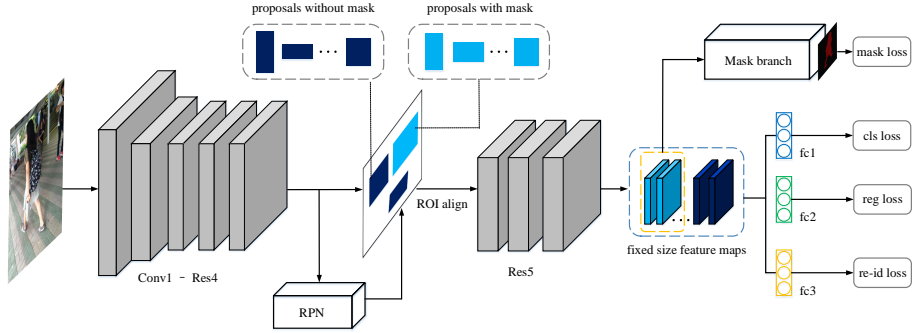


Figure 3: The schematic of our proposed segmentation masks guided person search framework. The model is trained end-to-end with multi-task loss. We adopt RPN to generate candidate boxes, and we denote the proposals associated with labeled person as proposals with mask, while the proposals associated with unlabeled person as proposals without mask since we only partially label the segmentation masks for the labeled persons in a raw image. Feature vectors of all candidate boxes go into the regression loss, classification loss and identification loss, whereas only the feature maps of the proposals with mask are fed into the mask branch.

bounding boxes and the ground truth bounding boxes is larger than pre-defined intersection over union (IOU) threshold, which equals to 0.5.

#### 4.2. Implementation Details

**Training Phase.** During training, we initialize our residual backbone with the ImageNet pretrained ResNet-50 model and adopt SGD as optimizer. The initial learning rate sets to 0.0004 and decayed by a factor of 0.1 at every 4 epochs. Because of the large memory consumption of the Faster R-CNN framework [12], we set the batch size to 1 during the 12 training epochs. All our experiments are implemented by Pytorch on Titan X Pascal GPU.

It should be noticed that, the shorter side of the input images is resized to 600 pixels, and we also augment the training data by horizontal flipping the training images and their related ground truth bounding boxes as well as the ground truth masks. In particular, the ground truth masks are resized to  $14 \times 14$  to match the masks generated from the mask branch. For the implementation of the mask branch, we adopt a similar architecture as in the Mask R-CNN [48].

6 losses are used jointly to supervise the training process.

**Inference Phase.** At test time, the shorter side for both query and gallery images is resized to 600 pixels as in the training process. We use the features generated from the last residual block (Res5) to represent each pedestrian, either probe person or the persons detected from the gallery set. Euclidean distance is then computed for each probe and gallery pair to assess the similarity level.

### 4.3. Comparison with State-of-the-Art Approaches

In this subsection, we report the person search performance of our model on our newly labeled person search dataset CUHK-SYSU, and we also give the comparison to several state-of-the-art approaches, including methods that optimize pedestrian detection and identification jointly (OIM [17], IAN [18], NPSM [42], QEEPS [43] and GCNPS [49]), as well as the methods solving pedestrian detection and person identification separately (DSIFT+Euclidean[50], DSIFT+KISS-ME[51], BoW[52]+Cosine similarity, LO-MO+XQDA, and MGTS [15]).

#### 4.3.1. Overall Person Search Performance on CUHK-SYSU

The comparative results with gallery size 100 are summarized in Table 2. We follow the annotations defined in [15] and [49], where “CNN” denotes the Faster R-CNN detector with ResNet-50 backbone, and “CNN<sub>v</sub>” denotes the VGG-based detector.

The methods above the dash line handle pedestrian detection and person identification separately. It can be observed that the deep CNN based pedestrian features [15] achieved better performance than hand-crafted features [50][51][52]. CNN<sub>v</sub>+MGTS [15] also utilizes segmentation mask to produce more discriminative features by filtering out the background, and achieves the best performance among those methods addressing pedestrian detection and person identification separately. Our proposed method uses segmentation mask to guide the network to extract discriminative pedestrian features by specifying the foreground regions. Meanwhile, pedestrian detection and person identification are optimized

jointly. Our framework achieved 3% gain compared with [15] on both mAP and top-1 matching rate.

All the joint methods (below the dash line) are built upon the Faster R-CNN [12] framework where OIM [17] can be regarded as the benchmark. The major distinction between our method and OIM [17] is that a new pedestrian segmentation mask branch is added. We achieve a significant performance improvement, with 10.8% mAP and 7.8% top-1 higher compared with [17]. It demonstrates the importance of the pedestrian segmentation mask and the newly labeled dataset. Other methods [18][42][43][49] are state-of-the-art person search approaches with good performance. IAN [18] improves the person search performance by introducing center loss to reduce the intra-class variations. NPSM [42] designs a person search approach by recursively shrinking the search area. QEEPS [43] proposes a strong person search framework by learning query-guided global context. [49] utilizes GCN to explore the impact of context persons on the person search task. Nevertheless, we still achieve 2% gain on both mAP and top-1 accuracy compared with [43], and 2% improvement on mAP compared with [49], all of which prove the effectiveness of our method.

The visualization of person search results on the CUHK-SYSU dataset are shown in Fig. 5. The upper images in each group are the searching results of OIM [17], and the lower images in each group are that of our model. It is observed that the persons in the bounding boxes of the third and the fourth images in the upper rows of group (a) and (b), as well as the person in the bounding box of the third image in the upper row of group (c), are different from their probe images. However, these persons are ranked before the persons who have the same identities as the probe images, simply because their background is more similar to the probe images. Nevertheless, with the assist of partially labeled segmentation masks, our model focus on the foreground and can distinguish persons based on the detailed textural information rather than the background-noise.

Method	mAP(%)	top-1(%)
CNN + DSIFT + Euclidean [50]	34.5	39.4
CNN + DSIFT + KISSME [50][51]	47.8	53.6
CNN + BoW + Cosine [52]	56.9	62.3
CNN + LOMO + XQDA [3]	68.9	74.1
CNN <sub>v</sub> + MGTS [15]	83.0	83.7
OIM [17]	75.5	78.7
IAN(ResNet-50) [18]	76.3	80.1
NPSM [42]	77.9	81.2
QEEPS [43]	84.4	84.4
GCNPS [49]	84.1	<b>86.5</b>
Ours	<b>86.3</b>	<b>86.5</b>

Table 2: Comparison with the state-of-the-art on CUHK-SYSU dataset with gallery size equals to 100.

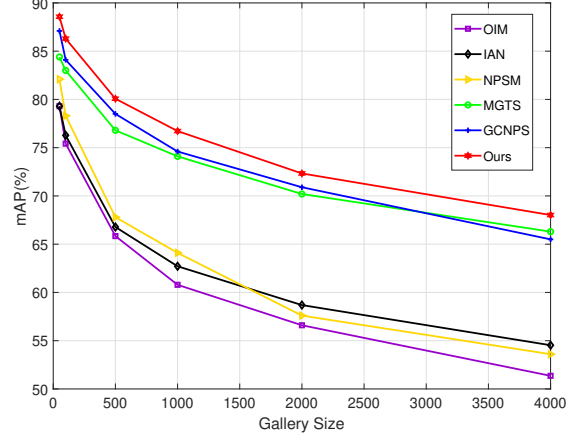
#### 4.3.2. Impact of Gallery Size

Each gallery image in CUHK-SYSU dataset contains around 6 pedestrians on average. With gallery size 100, person search aims to retrieve each target person from about 600 pedestrians. The person search is more challenging with the increasing number of gallery size. We also report the performance of our model with various gallery size, including [50, 100, 500, 1,000, 2,000, 4,000]. The results are demonstrated in Fig. 4. As expected, the person search performance of all methods drops with the increasing gallery size. While our person search framework remains superior than other approaches with various gallery sizes.

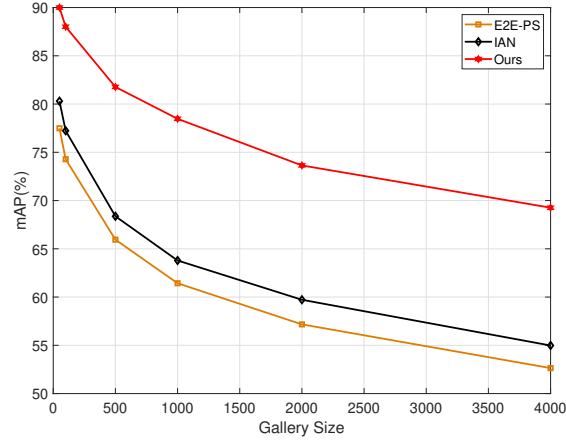
#### 4.3.3. Impact of Occlusion and Low Resolution

Person search becomes even harder when pedestrians are occluded or the resolution of the captured images is low. Therefore, to prove the robustness of our method, we further evaluate our model on two subsets. One subset contains 187 target persons with occlusion, and the other subset contains 290 target





(a)



(b)

Figure 4: Person search performance comparison on CUHK-SYSU dataset with different gallery size, [50, 100, 500, 1,000, 2,000, 4,000]. (a) Model with ResNet-50 backbone. (b) Model with ResNet-101 backbone.

persons with low resolution. The results are demonstrated in Table 3. We follow the notations defined in [18], where “whole” denote the full set which contains 2,900 probe images. We observe that the performance degenerate under these two extreme conditions compared with full set. However, our person search

Method	Low-Res		Occlusion		Whole	
	mAP(%)	top-1(%)	mAP(%)	top-1(%)	mAP(%)	top-1(%)
E2E-PS(VGGNet)	46.1	51.0	44.3	45.4	69.6	72.9
E2E-PS(Res-101)	47.9	52.0	47.7	48.1	74.2	78.1
IAN(Res-101)	52.6	54.4	53.0	54.5	77.2	80.4
Ours(Res-50)	<b>66.7</b>	<b>66.8</b>	<b>70.8</b>	<b>71.3</b>	<b>86.3</b>	<b>86.5</b>

Table 3: Person search performance on low resolution and occlusion subsets.

framework still outperforms the other approaches [16][18].

#### 4.4. Ablation Study

With the assist of the newly labeled dataset, our proposed person search framework produces more discriminative features by utilizing partially labeled segmentation mask. To evaluate the effectiveness of our approach, we report the person search performance when we progressively increase the number of images with segmentation mask. The results are shown in Table 4, where we denote the proportion of the images with segmentation mask as  $\alpha$ . In total, 1,833 images are labeled with segmentation mask, which accounts for around 16% of the 11,206 training images. When all those 1,833 images are used for training, we denote as “Full”. It can be observed that there is an obvious gain when 12% images with segmentation mask are used for training, and tend to be stable until 15% images are used. That is why we only label 16% of all the images.

Value of $\alpha$	3%	6%	9%	12%	15%	Full
mAP(%)	85.1	85.3	85.3	86.1	<b>86.3</b>	<b>86.3</b>
top-1(%)	85.2	85.4	85.7	86.1	<b>86.5</b>	<b>86.5</b>

Table 4: Person search performance on CUHK-SYSU dataset with various proportion of images with segmentation masks



(a)



(b)



(c)

Figure 5: Three groups of top-3 comparison results for person search on CHUK-SYSU dataset. The upper row in each group are the searching results of OIM [17], and the lower row in each group are the searching results of our model (both models adopt ResNet-50 as backbone). The blue boxes in the first column indicate the probe images, and the green boxes in other columns indicate their top-3 searching results. Best viewed in color.

## 5. Conclusion

Person search handles the challenges from both pedestrian detection and person identification, and inevitably introduces background clutters into the detected candidate boxes. To address this problem, with the assist of our new created dataset which contains the labeled segmentation masks for a portion of images in the existing CUHK-SYSU dataset, we propose a novel segmentation mask guided person search framework to extract more discriminative and robust features invariant to background clutters for each human individual. Moreover, our person search framework is trained end-to-end, which proves that joint optimization of pedestrian detection, person re-identification, and pedestrian segmentation is an effective solution for person search. Finally, extensive experiments show that our proposed method achieves state-of-the-art performance on CUHK-SYSU dataset.

## References

- [1] L. Zheng, Y. Yang, A. G. Hauptmann, Person re-identification: Past, present and future, arXiv preprint arXiv:1610.02984.
- [2] A. Hermans, L. Beyer, B. Leibe, In defense of the triplet loss for person re-identification, arXiv preprint arXiv:1703.07737.
- [3] S. Liao, Y. Hu, X. Zhu, S. Z. Li, Person re-identification by local maximal occurrence representation and metric learning, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2015, pp. 2197–2206.
- [4] D. Gray, H. Tao, Viewpoint invariant pedestrian recognition with an ensemble of localized features, in: European conference on computer vision, Springer, 2008, pp. 262–275.
- [5] L. Wu, C. Shen, A. v. d. Hengel, Personnet: Person re-identification with deep convolutional neural networks, arXiv preprint arXiv:1601.07255.

- [6] F. Wang, W. Zuo, L. Lin, D. Zhang, L. Zhang, Joint learning of single-image and cross-image representations for person re-identification, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 1288–1296.
- [7] H. Zhao, M. Tian, S. Sun, J. Shao, J. Yan, S. Yi, X. Wang, X. Tang, Spindle net: Person re-identification with human body region guided feature decomposition and fusion, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 1077–1085.
- [8] M. Tian, S. Yi, H. Li, S. Li, X. Zhang, J. Shi, J. Yan, X. Wang, Eliminating background-bias for robust person re-identification, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 5794–5803.
- [9] C.-P. Tay, S. Roy, K.-H. Yap, Aanet: Attribute attention network for person re-identifications, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 7134–7143.
- [10] S. Li, S. Bak, P. Carr, X. Wang, Diversity regularized spatiotemporal attention for video-based person re-identification, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 369–378.
- [11] W. Li, X. Zhu, S. Gong, Harmonious attention network for person re-identification, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 2285–2294.
- [12] S. Ren, K. He, R. Girshick, J. Sun, Faster r-cnn: Towards real-time object detection with region proposal networks, in: Advances in neural information processing systems, 2015, pp. 91–99.
- [13] W. Ouyang, X. Wang, Joint deep learning for pedestrian detection, in: Proceedings of the IEEE International Conference on Computer Vision, 2013, pp. 2056–2063.

- [14] W. Ouyang, X. Wang, A discriminative deep model for pedestrian detection with occlusion handling, in: 2012 IEEE Conference on Computer Vision and Pattern Recognition, IEEE, 2012, pp. 3258–3265.
- [15] D. Chen, S. Zhang, W. Ouyang, J. Yang, Y. Tai, Person search via a mask-guided two-stream cnn model, in: Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 734–750.
- [16] T. Xiao, S. Li, B. Wang, L. Lin, X. Wang, End-to-end deep learning for person search, arXiv preprint arXiv:1604.01850 2 (2016) 2.
- [17] T. Xiao, S. Li, B. Wang, L. Lin, X. Wang, Joint detection and identification feature learning for person search, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 3415–3424.
- [18] J. Xiao, Y. Xie, T. Tillo, K. Huang, Y. Wei, J. Feng, Ian: the individual aggregation network for person search, Pattern Recognition 87 (2019) 332–340.
- [19] N. Dalal, B. Triggs, Histograms of oriented gradients for human detection, 2005.
- [20] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, D. Ramanan, Object detection with discriminatively trained part-based models, IEEE transactions on pattern analysis and machine intelligence 32 (9) (2009) 1627–1645.
- [21] P. Dollár, Z. Tu, P. Perona, S. Belongie, Integral channel features.
- [22] S. Zhang, C. Bauckhage, A. B. Cremers, Informed haar-like features improve pedestrian detection, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2014, pp. 947–954.
- [23] W. Nam, P. Dollár, J. H. Han, Local decorrelation for improved pedestrian detection, in: Advances in Neural Information Processing Systems, 2014, pp. 424–432.

- [24] I. Filali, M. S. Allili, N. Benblidia, Multi-scale salient object detection using graph ranking and global–local saliency refinement, *Signal Processing: Image Communication* 47 (2016) 380–401.
- [25] H. Cholakkal, J. Johnson, D. Rajan, A classifier-guided approach for top-down salient object detection, *Signal Processing: Image Communication* 45 (2016) 24–40.
- [26] J. Dai, Y. Li, K. He, J. Sun, R-fcn: Object detection via region-based fully convolutional networks, in: *Advances in neural information processing systems*, 2016, pp. 379–387.
- [27] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, S. Belongie, Feature pyramid networks for object detection, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2117–2125.
- [28] R. Girshick, Fast r-cnn, in: *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1440–1448.
- [29] L. Zhang, L. Lin, X. Liang, K. He, Is faster r-cnn doing well for pedestrian detection?, in: *European conference on computer vision*, Springer, 2016, pp. 443–457.
- [30] L. Zheng, H. Zhang, S. Sun, M. Chandraker, Y. Yang, Q. Tian, Person re-identification in the wild, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1367–1376.
- [31] D. Cheng, Y. Gong, S. Zhou, J. Wang, N. Zheng, Person re-identification by multi-channel parts-based cnn with improved triplet loss function, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 1335–1344.
- [32] R. R. Varior, M. Haloi, G. Wang, Gated siamese convolutional neural network architecture for human re-identification, in: *European conference on computer vision*, Springer, 2016, pp. 791–808.

- [33] T. Xiao, H. Li, W. Ouyang, X. Wang, Learning deep feature representations with domain guided dropout for person re-identification, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 1249–1258.
- [34] N. McLaughlin, J. Martinez del Rincon, P. Miller, Recurrent convolutional network for video-based person re-identification, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 1325–1334.
- [35] D. Yi, Z. Lei, S. Liao, S. Z. Li, Deep metric learning for person re-identification, in: 2014 22nd International Conference on Pattern Recognition, IEEE, 2014, pp. 34–39.
- [36] W. Li, R. Zhao, T. Xiao, X. Wang, Deepreid: Deep filter pairing neural network for person re-identification, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2014, pp. 152–159.
- [37] H. Liu, J. Feng, M. Qi, J. Jiang, S. Yan, End-to-end comparative attention networks for person re-identification, *IEEE Transactions on Image Processing* 26 (7) (2017) 3492–3506.
- [38] W. Chen, X. Chen, J. Zhang, K. Huang, Beyond triplet loss: a deep quadruplet network for person re-identification, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 403–412.
- [39] L. Zheng, Z. Bie, Y. Sun, J. Wang, C. Su, S. Wang, Q. Tian, Mars: A video benchmark for large-scale person re-identification, in: European Conference on Computer Vision, Springer, 2016, pp. 868–884.
- [40] M. Geng, Y. Wang, T. Xiang, Y. Tian, Deep transfer learning for person re-identification, *arXiv preprint arXiv:1611.05244*.
- [41] Y. Xu, B. Ma, R. Huang, L. Lin, Person search in a scene by jointly modeling people commonness and person uniqueness, in: Proceedings of the 22nd ACM international conference on Multimedia, ACM, 2014, pp. 937–940.



- [42] H. Liu, J. Feng, Z. Jie, K. Jayashree, B. Zhao, M. Qi, J. Jiang, S. Yan, Neural person search machines, in: Proceedings of the IEEE International Conference on Computer Vision, 2017, pp. 493–501.
- [43] B. Munjal, S. Amin, F. Tombari, F. Galasso, Query-guided end-to-end person search, arXiv preprint arXiv:1905.01203.
- [44] J. Long, E. Shelhamer, T. Darrell, Fully convolutional networks for semantic segmentation, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2015, pp. 3431–3440.
- [45] Y. Li, H. Qi, J. Dai, X. Ji, Y. Wei, Fully convolutional instance-aware semantic segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 2359–2367.
- [46] K. Wada, labelme: Image Polygonal Annotation with Python, <https://github.com/wkentaro/labelme> (2016).
- [47] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778.
- [48] K. He, G. Gkioxari, P. Dollár, R. Girshick, Mask r-cnn, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 2961–2969.
- [49] Y. Yan, Q. Zhang, B. Ni, W. Zhang, M. Xu, X. Yang, Learning context graph for person search, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 2158–2167.
- [50] R. Zhao, W. Ouyang, X. Wang, Unsupervised salience learning for person re-identification, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2013, pp. 3586–3593.
- [51] M. Koestinger, M. Hirzer, P. Wohlhart, P. M. Roth, H. Bischof, Large scale metric learning from equivalence constraints, in: 2012 IEEE conference on computer vision and pattern recognition, IEEE, 2012, pp. 2288–2295.

- [52] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, Q. Tian, Scalable person re-identification: A benchmark, in: Proceedings of the IEEE international conference on computer vision, 2015, pp. 1116–1124.