# Making Person Search Enjoy the Merits of Person Re-identification

Chuang Liu, Hua Yang, Qin Zhou, and Shibao Zheng

*Abstract*—**Person search is an extended task of person re-identification (Re-ID). However, most existing one-step person search works have not studied how to employ existing advanced Re-ID models to boost the one-step person search performance due to the integration of person detection and Re-ID. To address this issue, we propose a faster and stronger one-step person search framework, the Teacher-guided Disentangling Networks (TDN), to make the one-step person search enjoy the merits of the existing Re-ID researches. The proposed TDN can significantly boost the person search performance by transferring the advanced person Re-ID knowledge to the person search model. In the proposed TDN, for better knowledge transfer from the Re-ID teacher model to the one-step person search model, we design a strong one-step person search base framework by partially disentangling the two subtasks. Besides, we propose a Knowledge Transfer Bridge module to bridge the scale gap caused by different input formats between the Re-ID model and one-step person search model. During testing, we further propose the Ranking with Context Persons strategy to exploit the context information in panoramic images for better retrieval. Experiments on two public person search datasets demonstrate the favorable performance of the proposed method.**

*Index Terms*—**Person search, person re-identification, knowledge distillation, teacher-guided disentangling networks, ranking with context persons.**
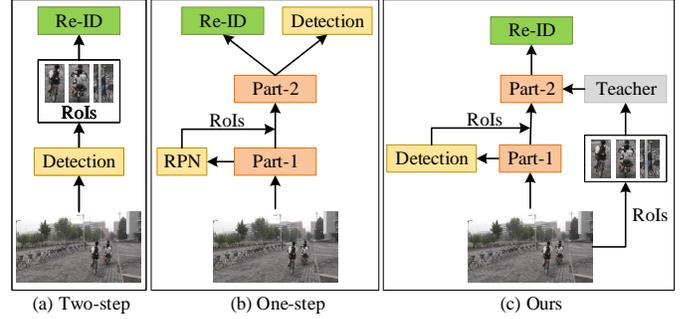
Fig. 1. Person search frameworks. (a) The two-step framework. It separates person search into two independent subtasks, namely person detection and person Re-ID. (b) The widely-used one-step framework. It is based on the Faster R-CNN framework and shares the whole backbone between the two subtasks. (c) The proposed framework. It partially disentangles the two subtasks to solve the conflict between them and is guided by a Re-ID teacher model to learn more discriminative Re-ID features.

## I. INTRODUCTION

**P**ERSON search task aims to locate the given query persons in the panoramic images. Different from person re-identification (Re-ID), person search unifies the person detection and re-identification tasks, making it more suitable for real-world applications. Existing CNN-based person search methods can be divided into two categories: two-step framework and one-step framework. As shown in Fig. 1, the two-step framework separates person search into two independent subtasks, person detection and person Re-ID, while the one-step framework integrates the detection and Re-ID subtasks into a unified and end-to-end trainable framework with shared networks. In the two-step framework, the Re-ID and detection models are trained separately, which ignores the correlations between them, leading to sub-optimal results. Therefore, we focus on the one-step framework to exploit the correlations between the two subtasks.

Chuang Liu, Hua Yang, and Shibao Zheng are with the Institute of Image Communication and Network Engineering, Department of Electronic Engineering, Shanghai Jiao Tong University, Shanghai 200240, China and Shanghai Key laboratory of Digital Media Processing and Transmission, Shanghai Jiao Tong University, China (Corresponding authors: Hua Yang and Shibao Zheng).

Qin Zhou are with the Institute of Medical Robotics, Shanghai Jiao Tong University, Shanghai 200240, China.

As an extended task of person Re-ID, it is intuitive that the one-step person search can benefit from the existing person Re-ID researches. However, due to model structure differences between the one-step person search model and independent Re-ID model, most previous one-step person search works have not studied how to make the advanced Re-ID models benefit the one-step person search model. Inspired by the thought of Knowledge Distillation [1], we propose a new one-step person search framework, the Teacher-guided Disentangling Networks (TDN), to make the one-step person search enjoy the merits of the powerful state-of-the-art Re-ID models. Specifically, we introduce a teacher branch into the TDN to guide the learning of the one-step person search model. The teacher branch is an independent Re-ID model. To generate discriminative enough Re-ID features as the guidance of the one-step person search model, a powerful state-of-the-art Re-ID model is chosen and pre-trained following classical Re-ID protocols.

However, there still remain challenges in transferring knowledge from the Re-ID teacher model to the one-step person search model. Firstly, as pointed out by [2], person detection focuses on the commonness of various persons, while person Re-ID pays attention to the uniqueness among different persons, leading to conflicting optimization targets in the one-step framework depicted in Fig. 1(b). This conflict prevents the one-step person search model from learning discriminative Re-ID features. Secondly, different from the traditional KD pipeline where the student model and teacher models have the same input, the Re-ID teacher model takes fixed-scale

person patches as input, while the one-step person search model takes the panoramic images as input, which leads to huge scale variations. As analyzed in [3], the CNN model is not robust against scale variations and can not learn scale-invariant features, making it difficult to directly transferring the strong prior knowledge learned on fixed-scale inputs to the one-step person search model.

To tackle these issues for effective knowledge transfer, firstly, we design a new strong one-step person search base framework for the proposed TDN, as shown in Fig. 1(c). The proposed new base framework partially disentangles the two subtasks by only sharing part convolutional networks between them. In this way, the conflict of optimization targets between two subtasks can be reduced to allow the one-step person search model to learn discriminative Re-ID features from the teacher branch more easily. Secondly, we propose a Knowledge Transfer Bridge (KTB) module to solve the scale variations for knowledge transfer. The KTB processes the same inputs as the teacher branch, and fuses its output feature maps into the one-step person search. The KTB functions as a bridge to help transfer knowledge from the teacher model to the one-step person search model. Altogether, the proposed TDN can effectively learn from the powerful Re-ID teacher model to generate discriminative Re-ID features. As a result, the proposed TDN demonstrates faster and stronger than the previous one-step framework.

Besides the proposed TDN, we also improve the ranking process. The widely-used ranking process only considers the individual similarity between persons, ignoring the context persons provided by the panoramic scene images. To exploit the context information in the panoramic scene images in the ranking phase, we design a Ranking with Context Persons (RCP) strategy which takes context persons as additional cues and generates better ranking results. Specifically, in the proposed RCP, the co-occurrence index score is proposed to measure the correlation between the target person and a context person. The total co-occurrence index scores are taken as a supplement to the individual similarity score to generate the final similarity score for the ranking process.

To summarize, the contributions of this paper are as follows:

- We propose the Teach-guided Disentangling Networks (TDN) to make the one-step person search model enjoy the merits of the powerful state-of-the-art Re-ID models. The proposed TDN can significantly boost the person search performance by integrating the knowledge of powerful Re-ID models;
- We design the Ranking with Context Person (RCP) strategy to exploit the context persons for better retrieval performance.
- Experiments on two person search benchmarks demonstrate that the proposed method achieves much higher performance compared with state-of-the-art methods.

## II. RELATED WORK

**Person Re-ID.** Person Re-ID has been studied for many years and achieved great progress. Some works [4], [5], [6], [7], [8], [9], [10] focus on designing effective metric principles for person Re-ID. Some works [11], [12], [13], [14], [15], [16], [17], [18] pay attention to developing novel CNN models to learn discriminative features. As an extended task of person Re-ID, person search task should benefit from the existing person Re-ID researches.

**Person Search.** In recent years, person search task has already received a lot of attention due to the integration of person detection and Re-ID. Many two-step structure methods are proposed in [19], [20], [2], [21], [22]. Zheng et al. [19] first propose the two-step person search framework and evaluate the impact of combinations of various person detector and person Re-ID models. Lan et al. [20] propose the Cross-Level Semantic Alignment to solve the multi-scale challenge by combining cross-level feature maps with the Faster R-CNN [23] as person detector. Wang et al. [22] propose a Task-Consist Two-Stage person search framework including an identity-guided query detector to generate query-like person detections and a Detection Results Adapted Re-ID model to make the Re-ID model adapted to the detections.

Different from the above two-step methods, some researchers propose to tackle person search using the one-step structure methods [24], [25], [26], [27], [26], [28], [29]. Xiao et al. [24] first propose the one-step framework based on the Faster R-CNN and design the OIM loss function to solve the ill-conditioned training problem. Munjal et al. [26] propose a query-guided one-step person search framework which employs the query to generate query-relevant proposals. Chen et al. [29] propose the Norm-Aware Embedding method which uses the norm and angle of the person embedding to conduct person detection and person re-identification, respectively.

**Knowledge Distillation.** Hinton et al. [1] propose the Knowledge Distillation (KD) to compress complex large models into a small model which can keep comparative performance. Romero et al. [30] introduce the intermediate representations from the teacher model as hints to improve the distillation training procedure. Zagoruyko et al. [31] propose to transfer knowledge by attention maps rather than feature maps. Yim et al. [32] define the flow of solution procedure (FSP) matrix to represent the flow between two layers and propose to transfer distilled knowledge by minimizing the distance between the teacher and student FPS matrices. Generally, these knowledge distillation works focus on how to effectively transfer knowledge from the teacher model to the student model, where both the teacher and student models solve the same task and have the same input.

Munjal et al. [33] also propose to apply knowledge distillation to person search. However, their method distills the knowledge of a pre-trained person detector to the one-step OIM person search model to improve the person detection. Differently, in this paper, inspired by the thought of knowledge distillation, we propose a new one-step person search framework (the TDN) to integrate the state-of-the-art person Re-ID models to promote the development of one-step person search. Different from the above-mentioned knowledge distillation works, the proposed TDN framework aims to transfer the knowledge from a single-task teacher model taking cropped person patches as inputs to a multi-task student model taking panoramic images as inputs, which brings some new chal-
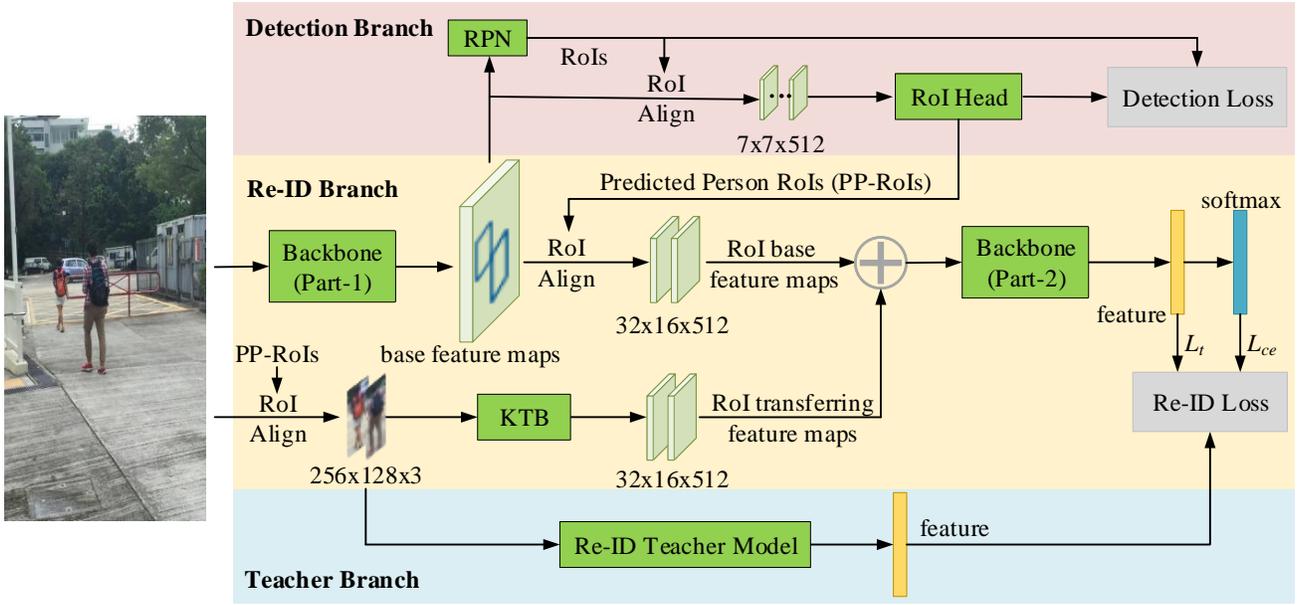
Fig. 2. Illustration of the proposed TDN. The TDN contains three branches, a detection branch, a Re-ID branch, and a teacher branch. The Re-ID branch is utilized to extract Re-ID features for the person RoIs predicted by the detection branch. The teach branch is employed to supervise the learning of the Re-ID branch and is removed during the reference.

lenges to the knowledge distillation process. For effective knowledge distillation, we take some measures to overcome the new challenges in the proposed TDN framework.

## III. METHOD

In this section, we will elaborate on the proposed TDN framework as well as the RCP ranking strategy. As illustrated in Fig. 2, the TDN framework includes a detection branch, a Re-ID branch, and a teacher branch. To perform person search, the detection branch is used to locate persons in the panoramic images based on the outputs of the Part-1, then Re-ID features are extracted from the Re-ID branch to represent each predicted person. The teacher branch is used to guide the Re-ID branch to learn discriminative features in the training phase. Finally, the RCP ranking method provides the ranking list for each query image.

### A. Teacher-guided Disentangling Networks

In the proposed TDN, to make the one-step person search model learn more discriminative Re-ID features from the powerful Re-ID teacher model, we design a new one-step base framework, namely the Partially Disentangled Framework, and propose a Knowledge Transfer Bridge (KTB) module. As shown in Fig. 2, the proposed TDN consists of a Re-ID branch, a detection branch and a teacher branch. The Re-ID branch and detection branch form the Partially Disentangled Framework which carries out the person detection and re-identification subtasks. The KTB module is introduced between the Re-ID branch and the teacher branch to help the one-step person search model learn from the teacher model more easily. The details are introduced in the following parts.

**Partially Disentangled Framework.** For the Re-ID branch, we adopt the ResNet50-IBN-a [34] as the backbone. The

ResNet50-IBN-a is divided into two parts, Part-1 and Part-2. The Part-1 is composed of layers from the conv1 to conv3_x, while the Part-2 includes layers from conv4_x to conv5_x. Following the Part-2, a Global Average Pooling (AVG) layer and a Batch Normalization (BN) layer are used to generate the final Re-ID features. Finally, a softmax layer classifies person identities in the training phase.

For the detection branch, we do not adopt the widely-used one-step framework [24] shown in Fig. 1(b). As analyzed in the introduction, the detection and Re-ID subtasks have a conflict of optimization targets in this widely-used one-step framework [24] which shares all convolutional networks between two subtasks. The conflict prevents the Re-ID subtask from learning discriminative Re-ID features for each detected person RoI. To address this problem, we propose to partially disentangle the two branches to retain both relevance and independence between them. Specifically, in the proposed architecture, only the Part-1 network is shared between the detection and Re-ID branches (as illustrated in Fig. 2), while Part-2 is only kept for the Re-ID branch to learn discriminative features for person RoI detected by the detection branch.

Based on the feature maps generated by the Part-1, a standard RPN [23] is built to generate possible person RoIs. And the RoI Align [35] operation is adopted to crop the feature maps of the RoIs into fixed-size feature maps ($7 \times 7 \times 512$). These 3-D feature maps are then flattened into 1-D feature vectors before fed into the RoI Head. The RoI Head is composed of two fully connected (FC) layers with 1024 units. Finally, a classification layer and a coordinate regression layer are employed to predict the probability that the RoIs are persons and refine their box coordinates, respectively. The Non Maximum Suppression (NMS) is applied to the proposals generated from the RoI Head to output the final predicted

person RoIs. The Part-2 only extracts Re-ID features for the final predicted person RoIs rather than all the RoIs from the RPN, which improves the running speed.

**Teacher-guided Learning with KTB.** For the teacher branch, we choose a state-of-the-art Re-ID model, the strong Re-ID baseline model [15], as the default teacher model. To pre-train the Re-ID teacher model, we construct a Re-ID style training set by cropping all labeled person patches from the panoramic images and resizing them to $256 \times 128 \times 3$ ones with the bilinear interpolation. Then, we pre-train the strong Re-ID teacher model following its training settings [15]. During the training of the TDN, the pretrained teacher model is utilized to extract features for each predicted person RoI to supervise the learning of the one-step person search model. Please note that the teacher branch is only used in the TDN training phase, which will not bring an extra computational burden during inference.

In the one-step person search framework, different from traditional fixed-size Re-ID inputs, features are extracted from predicted person RoIs with arbitrary sizes. As analyzed in [3], the CNN model is not robust against scale variations, making it difficult to directly transfer the strong prior knowledge learned on fixed-size inputs to the Re-ID branch in the one-step person search model. To ensure effective knowledge transfer from the strong Re-ID teacher model to the Re-ID branch, we further propose the KTB module for the Re-ID branch to reduce scale variations between the Re-ID teacher model and the Re-ID branch. Specifically, denote a predicted person RoI with varying sizes as $x$, we introduce $\hat{x}$ as a supplementary input of the Re-ID branch, where $\hat{x}$ is a fixed-size ($256 \times 128 \times 3$) image interpolated from $x$ by the RoI Align operation. Then, the KTB module processes the supplementary input $\hat{x}$ to generate the $32 \times 16 \times 512$ RoI transferring feature maps. The KTB has the same structure as Part-1 in the Re-ID branch, but does not share parameters with Part-1 to avoid entanglement with the detection task.

Next, we combine the RoI base feature maps extracted from the Part-1 with the RoI transferring feature maps generated from the KTB module. As the RoIs predicted by the detection branch are of variable sizes, to allow fusion with the fixed-size RoI transferring feature maps from the KTB, the RoI base feature maps are interpolated to $32 \times 16 \times 512$ ones using the RoI Align operation (as shown in Fig. 2). Then, the RoI base feature maps and the RoI transferring feature maps can be fused by pixel-wise addition. The fused feature maps contain feature maps extracted from the same inputs as the teacher model, and consequently can help the Re-ID branch to learn from the teacher model more easily by reducing scale variations.

Finally, in the Re-ID branch, the Part-2 together with the "GAP+BN" layer is utilized to further process the fusion feature maps to generate the final 2048-D Re-ID features, followed by a softmax classifier. Additionally, as shown in Fig. 2, there are two RoI Align operations in the Re-ID branch. One is performed on the base feature maps, and the other is conducted on the input scene image. Please kindly note that these two RoI Align operations use the same person RoIs predicted by the detection branch in both the training and inference stages.

**Training Loss.** Following the Faster R-CNN [23], we employ the RPN training losses ($L_{cls}^{rpn}$ and $L_{reg}^{rpn}$) and RoI Head training losses ($L_{cls}$ and $L_{reg}$) to train the detection branch. The total detection loss $L_{det}$ is defined as:

$$L_{det} = L_{cls}^{rpn} + L_{reg}^{rpn} + L_{cls} + L_{reg}. \qquad (1)$$

For the Re-ID branch, we adopt the following loss function to perform knowledge transfer.

$$L_t = \frac{1}{K} \sum_{x \in \text{RoIs}} \| f_s(x, \hat{x}) - f_t(\hat{x}) \|^2, \qquad (2)$$

where $K$ is the number of the final predicted person RoIs in a mini-batch, $f_t(\hat{x})$ is the L2-normalized Re-ID supervision feature from the teacher model, and $f_s(x, \hat{x})$ is the L2-normalized Re-ID feature from the Re-ID branch. Besides, we also introduce the cross-entropy loss $L_{ce}$. Overall, the total loss for the Re-ID branch is as follows:

$$L_{reid} = w L_t + L_{ce}, \qquad (3)$$

where $w$ is a weight factor to balance between the two components of $L_{reid}$. In the training phase, to let the teacher model lead the learning of Re-ID branch, we empirically set $w$ as follows:

$$w = \begin{cases} 5, & epoch < 15 \\ 11 - 0.4 \times epoch, & 15 \le epoch < 25 \\ 1, & \text{otherwise} \end{cases}. \qquad (4)$$

There are labeled and unlabeled person RoIs in person search datasets. Please kindly note that the $L_t$ loss is applied to both the labeled and unlabeled person RoIs, while the $L_{ce}$ loss is only applied to the labeled person RoIs.

The total training loss for the proposed TDN is the combination of the detection loss $L_{det}$ and the Re-ID loss $L_{reid}$, which is mathematically formulated as follows,

$$L = L_{reid} + L_{det}. \qquad (5)$$

### B. Ranking with Context Persons



Fig. 3. Illustration of the fact that context persons tend to simultaneously appear in more than one scenes.

Different from the traditional Re-ID task, the person search task provides the panoramic scene images which contain

context persons (co-travelers). Mazzon et al. [36] point out that persons are likely to walk in groups, which means that context persons of the query person are likely to simultaneously appear with him in more than one gallery scenes (as shown in Fig. 3). Thus, the context persons of the query person can be regarded as additional cues to help search for the query. In this paper, we propose an effective ranking method to take advantage of the context cues for better ranking results.

Given a query scene image containing the query person $q$ and context persons $q_1, \cdots, q_N$ and a gallery scene including gallery persons $g_1, g_2, \cdots, g_M$, it is a conventional practice to take the most similar gallery person (denoted as $g$) to the query person $q$ as the candidate for this gallery scene image and rank all candidates from all gallery scene images based on the individual similarity. This common procedure ignores the valuable information of context persons. It is intuitive that if more context persons of the query person $q$ appear in the gallery scene image with him, it is more likely that the $g$ is a right candidate of $q$. However, it is still unknown how to utilize the context persons to rectify the matching degree of the target query-gallery pair $(q, g)$.

In fact, one one hand, when the target query-gallery pair $(q, g)$ is not the same person, any context person should make few contributions to the rectification, even if this context person appears in both the query and gallery scenes. On the other hand, when the target pair $(q, g)$ is the same person, the contributions a context person makes to the rectification should depend on the possibility that this context person also appears in the gallery scene. To model this relationships between the query person $q$ and a context person $q_i$, we propose the co-occurrence index score based on the target query-gallery pair $(q, g)$ and a context query-gallery pair $(q_i, g_j)$. The co-occurrence index score $S_{co}(q, q_i)$ for the query person $q$ and a context person $q_i$ is defined as follows:

$$S_{co}(q, q_i) = \begin{cases} S_I(q, g) \cdot S_I(q_i, g_j), & S_I(q_i, g_j) \geq b \geq 0 \\ 0, & \text{otherwise} \end{cases} \quad (6)$$

where $S_I(q, g)$ is the cosine similarity between person $q$ and $g$, the gallery person $g_j$ is the most similar person to the context person $q_i$ in the gallery scene image, and $b$ is a positive threshold to keep confident context query-gallery pair $(q_i, g_j)$.

Based on the co-occurrence index score, we define the context similarity score $S_C(q, g)$ for the target query-gallery pair $(q, g)$ as follows:

$$S_C(q, g) = \sum_{i=1}^{N} S_{co}(q, q_i). \quad (7)$$

During the ranking phase, the context similarity score $S_C(q, g)$ is regarded as an additional cue to rank the gallery persons. Specifically, we combine it with $S_I(q, g)$ to define the final similarity score $S(q, g)$.

$$S(q, g) = S_I(q, g) + \lambda \cdot S_C(q, g), \quad (8)$$

where $\lambda \in (0, 1)$ balances the relative importance between $S_I(q, g)$ and $S_C(q, g)$.

## IV. EXPERIMENTS

In this section, we conduct experiments on two public person search datasets, the PRW [19] and CUHK-SYSU [24] datasets, and compare the proposed method with state-of-the-art methods. Afterward, thorough ablation studies are conducted to validate the effectiveness of each component.

### A. Datasets

**PRW** dataset [19] is collected in a university with six cameras. It provides 11,816 video frames with 43,110 bounding boxes. Among them, 5,704 frames are split into the training set, 2,057 person bounding boxes are taken as the query set, and 6,112 frames are taken as the gallery set. The training set contains 15,575 bounding boxes from 482 labeled identities and many unlabeled bounding boxes. Both the training set and gallery sets include plenty of unlabeled bounding boxes which have different identities from those labeled ones. Different from the CUHK-SYSU, the search scope is the whole gallery set, making it more challenging than the CUHK-SYSU.

**CUHK-SYSU** dataset [24] contains video frames from the street snap and movies. It provides 18,184 frames with total 96,143 bounding boxes containing both labeled and unlabeled identities. The training set consists of 11,206 frames containing 15,080 bounding boxes from 5,532 labeled identities and many unlabeled bounding boxes, while the testing set is composed of 2,900 query persons and 6,987 gallery frames. For each query person, it provides several gallery subsets with various gallery sizes.

### B. Evaluation Protocol

The widely-used Cumulative Matching Characteristic (CMC) and mean Average Precision (mAP) in the Re-ID task are adopted to evaluate the performance of person search. For each query, its AP is scaled by its recall rate. Then the mAP is computed as the average of all APs across all query persons.

### C. Implementation Details

Following the settings proposed in [15], we pre-train the strong Re-ID teacher model. When training the TDN, the Re-ID teacher model is frozen. The TDN is trained for total 40 epochs with batch size 16 using the SGD optimizer. The initial learning rate is 0.01, decayed to 0.001 and 0.0001 in the 25-th and 30-th epochs, respectively. For the RPN in the detection branch, the anchor sizes are set to 4, 8, 16, and 32 for each location on feature maps. Since the height of a person is not less than the width in a scene image, the anchor aspect ratios are empirically set to 1, 2, and 3. The height and width of an input scene image are scaled by the same factor to make the shorter side not less than 640 pixels or the longer side not more than 960 pixels. During reference, the predicted person RoIs with foreground scores lower than 0.5 are removed, and only person RoIs whose Intersection over Union (IoU) with ground truth bounding boxes larger than 0.5 are regarded as true detection results. In addition, the parameters $b$ in (6) and $\lambda$ in (8) are 0.3 and 0.2, respectively.

TABLE I
COMPARISON WITH STATE-OF-THE-ART METHODS. THE GALLERY SIZE IS
100 FOR THE CUHK-SYSU DATASET, WHILE THE WHOLE GALLERY SET
IS UTILIZED FOR THE PRW DATASET. THE STRONG Re-ID MODEL IS USED
AS THE TEACHER MODEL FOR OUR METHOD HERE.

| | Method | PRW | | CUHK-SYSU | |
|---|---|---|---|---|---|
| | | mAP (%) | top-1 (%) | mAP (%) | top-1 (%) |
| two-step | DPM+IDE [19] | 20.5 | 48.3 | - | - |
| | MGTS [2] | 32.6 | 72.1 | 83.0 | 83.7 |
| | CLSA [20] | 38.7 | 65.0 | 87.2 | 88.5 |
| | RDLR [21] | 42.9 | 70.2 | 93.0 | 94.2 |
| | TCTS [22] | 46.8 | 87.5 | 93.9 | 95.1 |
| one-step | OIM [24] | - | - | 75.7 | 78.7 |
| | IAN [37] | 23.0 | 61.9 | 76.3 | 80.1 |
| | NPSM[38] | 24.2 | 53.1 | 77.9 | 81.2 |
| | RCAA [39] | - | - | 79.3 | 81.3 |
| | LCGPS [40] | 33.4 | 73.6 | 84.1 | 86.5 |
| | QEEPS [26] | 37.1 | 76.7 | 88.9 | 89.1 |
| | NAE+ [29] | 44.0 | 81.1 | 92.1 | 92.9 |
| | APNet [41] | 41.9 | 81.4 | 88.9 | 89.3 |
| | IGPN [42] | 47.2 | 87.0 | 90.3 | 91.4 |
| | BINet [28] | 45.3 | 81.7 | 90.0 | 90.7 |
| | PSFL [43] | 44.2 | 85.2 | 92.3 | 94.7 |
| | Ours | **70.2** | **93.5** | **94.9** | **96.3** |

### D. Comparison with State-of-the-art Methods

In this section, we compare the proposed method with some state-of-the-art methods. Experimental results are reported in Table I. In this section, the strong Re-ID model [15] is taken as the default teacher model in our method.

**Results on PRW.** As shown in Table I, the proposed method achieves 70.2% mAP and 93.5% top-1 recognition rate, outperforming all the compared state-of-the-art methods. Among one-step state-of-the-art methods, the IGPN obtains the highest performance (47.2% mAP and 87.0% top-1). Compared to the IGPN, the proposed method achieves 23.0% mAP and 6.5% top-1 improvement. We also make comparisons with some two-step methods among which the TCTS method achieves the best performance. According to Table I, our proposed method outperforms TCTS by 23.4% in mAP and 6.0% in top-1 recognition rate respectively. It is observed that the proposed method surpasses all the listed state-of-the-art methods by a large margin, especially in mAP, which shows that the proposed method can retrieve more positive gallery persons at the top of the ranking list.

**Results on CUHK-SYSU.** As shown in Table I, the proposed method achieves 94.9% mAP and 96.3% top-1 recognition rate, surpassing all the compared state-of-the-art methods. For example, compared to TCTS which achieves the best performance (93.9% mAP and 95.1% top-1) among all the compared state-of-the-art methods, the proposed method achieves 1.0% gains in mAP and 1.2% gains in top-1 performance.

Besides, we also conduct experiments to evaluate the influence of various gallery sizes. The gallery sizes vary from 50 to 4,000. Detailed comparison results are presented in Fig. 4. On one hand, we observe that the performance of all methods decreases with gallery size increasing. On the other hand, the proposed method generates much higher and more stable performance compared with other state-of-the-
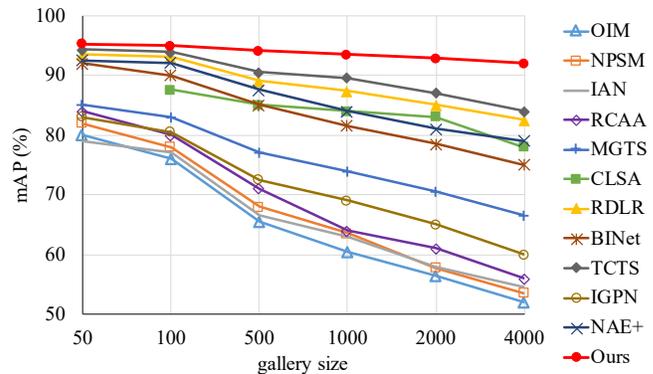


Fig. 4. Performance comparison on the CUHK-SYSU dataset with gallery size increasing from 50 to 4000.

art methods, especially when the gallery size is large. The robustness against various gallery sizes is valuable in practical scenarios.

### E. Ablation Study

In this section, experiments are first conducted to validate the effectiveness of each proposed component. Then, the comparison results between the one-step and two-step frameworks are reported. Next, the running speed of the previous one-step framework and the proposed TDN is compared. Finally, the parameter analysis is conducted to study the influence of hyper-parameters in the proposed RCP.

TABLE II
EFFECTIVENESS OF THE PROPOSED TDN AND RCP. "OIM*" REPRESENTS
THE MODIFIED OIM METHOD, THE PREVIOUS WIDELY-USED ONE-STEP
FRAMEWORK.

| Method | PRW | | CUHK-SYSU | |
|---|---|---|---|---|
| | mAP (%) | top-1 (%) | mAP (%) | top-1 (%) |
| OIM* | 22.8 | 49.1 | 81.5 | 83.2 |
| TDN | 68.7 | 91.9 | 93.8 | 94.9 |
| TDN + RCP | **70.2** | **93.5** | **94.9** | **96.3** |

**Effectiveness of TDN.** To validate the effectiveness of the proposed TDN, we re-implement the widely-used one-step person search framework (denoted as OIM*) proposed in [24]. For a fair comparison, we replace the original ResNet50 backbone with the ResNet50-IBN-a, and divides the backbone into two parts from the conv3_4. As shown in Table II, on the PRW dataset, the proposed TDN surpasses the OIM* by 45.9% in mAP and 42.8% in top-1. On the CUHK-SYSU dataset, the proposed TDN outstrips the OIM* by 12.3% in mAP and 11.7% in top-1. This demonstrates that the proposed TDN one-step person search framework can significantly boost the one-step person search performance by integrating the advanced Re-ID model.

In the proposed TDN, to help the one-step person search model learn discriminative Re-ID features from the powerful Re-ID teacher model, we propose to partially disentangle the detection and Re-ID subtasks and introduce the KTB

TABLE III
IMPACT OF COMPONENTS IN THE PROPOSED TDN.

| Method | PRW | | CUHK-SYSU | |
|---|---|---|---|---|
| | mAP (%) | top-1 (%) | mAP (%) | top-1 (%) |
| OIM* | 22.8 | 49.1 | 81.5 | 83.2 |
| Baseline | **38.0** | **77.4** | **82.4** | **84.8** |
| OIM* only w/ teacher | 51.0 | 82.7 | 88.5 | 90.5 |
| Baseline only w/ teacher | **63.0** | **90.2** | **93.1** | **94.5** |
| TDN (w/o KTB) | 63.0 | 90.2 | 93.1 | 94.5 |
| TDN (w/o teacher) | 55.2 | 89.5 | 90.0 | 91.8 |
| TDN | **68.7** | **91.9** | **93.8** | **94.9** |

TABLE IV
RESULTS OF DIFFERENT STATE-OF-THE-ART RE-ID MODELS AS THE
TEACHER MODEL IN THE PROPOSED TDN FRAMEWORK. THE "NONE"
MEANS THE BASELINE MODEL WHERE NO TEACHER MODEL IS APPLIED.
THE PROPOSED RCP RANKING METHOD IS NOT USED HERE.

| Teacher model | PRW | | CUHK-SYSU | |
|---|---|---|---|---|
| | mAP (%) | top-1 (%) | mAP (%) | top-1 (%) |
| None (baseline) | 38.0 | 77.4 | 82.4 | 84.8 |
| Strong Re-ID [15] | 68.7 | 91.9 | 93.8 | 94.9 |
| MGN [44] | 71.4 | 92.7 | 93.5 | 95.1 |
| AGW [45] | 68.9 | 93.1 | 92.9 | 94.5 |

module. To validate their effectiveness, we construct a baseline model by removing the teacher branch and KTB module in the proposed TDN. The baseline is also trained end-to-end. Experimental results are reported in Table III. Compared to the OIM*, the proposed baseline achieves better results on both datasets, especially on the more challenging PRW dataset (15.2% mAP and 28.3% top-1 improvement). This demonstrates that the proposed Partially Disentangled Framework can help the one-step person search model learn more discriminative Re-ID features, and can be used as a new strong baseline framework for the one-step person search model in place of the OIM one-step framework. Besides, we also employ the Re-ID teacher model to guide the learning of the OIM* model. As shown in Table III, compared to the "OIM* only w/ teacher", the "Baseline only w/ teacher" achieves much higher performance on both datasets. For example, on the PRW dataset, the "Baseline only w/ teacher" surpasses the "OIM* only w/ teacher" by 12.0% in mAP and 7.5% in top-1. This further validates that the proposed Partially Disentangled Framework is effective to help the one-step person search model learn from the powerful Re-ID teacher model.

Without the KTB module, the performance of TDN drops from 68.7% to 63.0% in mAP on the PRW dataset and from 93.8% to 93.1% in mAP on the CUHK-SYSU dataset. This validates that the KTB module is crucial to help the one-step person search model learn from the teacher model. When the Re-ID teacher branch is removed, the performance of TDN suffers heavy losses. On the PRW dataset, the mAP drops from 68.7% to 55.2% and the top-1 drops from 91.9% to 89.5%. On the CUHK-SYSU dataset, the mAP drops from 93.8% to 90.0% and the top-1 drops from 94.9% to 91.8%. This shows that it is effective to improve the one-step person search performance by learning from the powerful Re-ID model.

**Effectiveness of RCP.** As shown in Table II, by introducing RCP in the ranking phase, the person search performance further witnesses a gain of 1.5% in mAP and 1.6% in top-1 recognition rate on the PRW dataset, and 1.1% in mAP and 1.4% in top-1 recognition rate on the CUHK-SYSU dataset. It demonstrates that the proposed RCP can effectively improve performance by exploiting context information.

**Enjoying the merits of person Re-ID.** In this paper, the TDN is proposed to make the end-to-end person search enjoy the merits of person Re-ID researches. To demonstrate that the proposed TDN is able to realize this purpose, besides employing the strong Re-ID baseline model proposed by [15]

as the teacher model, we also adopt some other state-of-the-art Re-ID models as the teacher model, the MGN [44] model and AGW model [45]. Experimental results are reported in Table IV. It is observed that the proposed TDN achieves much higher performance on both person search datasets with any one of the three state-of-the-art Re-ID models as the teacher model when compared to the baseline model.

The MGN model is designed to extract multiple granularity locals features as well as a global feature for a person patch. The final feature of a person patch is obtained by concatenating the global feature and multiple granularity global features. Although the TDN has no local branches to learn local features for person RoIs, it can still learn discriminative Re-ID features with the guidance of the MGN teacher model. The AGW model integrates the Non-local Attention Block into the backbone networks to capture non-local relations and replaces the widely-used max-pooling or average pooling with a learnable generalized-mean pooling layer to capture the fine-grained discriminative features. Even if the TDN has no Non-local Attention Blocks and the generalized-mean pooling layer to learn the non-local relations between pixels and fine-grained features, it can also learn discriminative Re-ID features from the AGW teacher model. These experimental results demonstrate that the proposed TDN can significantly boost the performance of the end-to-end person search by enjoying the merits of state-of-the-art person Re-ID models.

**One-step vs. two-step.** A straightforward solution to person search is the two-step methods which employ a person detector to detect person RoIs and use a person Re-ID model to extract feature representations for the predicted person RoIs. To validate that the proposed one-step person search model TDN is a better solution to person search task compared to the two-step ones, we construct several two-step methods by training a Faster R-CNN person detector and employing different Re-ID teacher model to extract features for the predicted person RoIs. Comparison results are reported in Table V. All the three two-step methods achieve pretty high performance on both datasets with the advanced Re-ID models as the feature extractor. It is also observed that these two-step methods outperforms most state-of-the-art person search methods, which shows that the state-of-the-art Re-ID models are very powerful when applied to person search task. However, when adopting any one of these advanced Re-ID models as the teacher model, the proposed TDN realizes higher performance than the corresponding two-step method. It validates that the

TABLE V
COMPARISON OF ONE-STEP AND TWO-STEP METHODS. THE "TDN (MGN)" REPRESENTS THE TDN MODEL WITH THE MGN RE-ID MODEL AS THE TEACHER. THE PROPOSED RCP RANKING METHOD IS NOT USED HERE.

| Method | PRW | | CUHK-SYSU | |
|---|---|---|---|---|
| | mAP (%) | top-1 (%) | mAP (%) | top-1 (%) |
| Faster R-CNN + Strong Re-ID | 66.9 | 91.3 | 92.1 | 94.1 |
| TDN (Strong Re-ID) | **68.7** | **91.9** | **93.8** | **94.9** |
| Faster R-CNN + MGN | 66.8 | 92.3 | 90.5 | 92.6 |
| TDN (MGN) | **71.4** | **92.7** | **93.5** | **95.1** |
| Faster R-CNN + AGW | 67.1 | 91.3 | 91.4 | 93.4 |
| TDN (AGW) | **68.9** | **93.1** | **92.9** | **94.5** |

proposed TDN can not only learn from the powerful Re-ID teacher models but also can be more excellent than the teacher by exploiting the correlations between two subtasks in the end-to-end one-step framework.

**Running Speed.** Most of the previous one-step methods adopt the Faster R-CNN-based one-step framework (denoted as OIM*) and design more additional modules based on it. This means that the running speed of these methods is slower than the OIM*. Thus, in Table VI, we only compare the running speed of the OIM* and the proposed TDN. It is observed that the proposed TDN runs three times as fast as the OIM*. In the OIM* framework, the second part has to process the RoI feature maps from all the RoIs generated by the RPN (e.g. 128 RoIs per image). However, the RoIs generated by the RPN contain a lot of invalid proposals which do not contain persons or are removed after the NMS. The computation on invalid RoIs greatly lowers the running speed. In contrast, in the proposed TDN, the Part-2 in the Re-ID branch only processes the final detections generated by the lightweight detection branch for each scene image. Assuming 10 final predicted person bounding boxes in each frame on average, we quantitatively compare the computation complexity of main parts for two one-step frameworks. As shown in Table VII, the computation of the proposed TDN is much lower than the previous one-step framework OIM*. Thus, the proposed TDN can run much faster than the OIM*.

TABLE VI
RUNNING SPEED COMPARISON OF THE OIM* AND TDN.

| Method | GPU | PRW | CUHK-SYSU |
|---|---|---|---|
| OIM* | TITAN RTX | 11.5 fps | 12.4 fps |
| TDN | TITAN RTX | 34.3 fps | 37.7 fps |

**Parameter Analysis.** We conduct experiments to evaluate the impact of parameters $b$ in (5) and $\lambda$ in (6).

As shown in Fig. 5, on the PRW dataset, both the mAP and top-1 can obtain stable gains when the threshold $b \leq 0.6$, and when $b > 0.6$, the mAP and top-1 begins to decrease slightly. On the CUHK-SYSU dataset, when $b \leq 0.3$, the performance achieves the highest improvement, and when $b \geq 0.3$, the improvement becomes smaller. Generally, the proposed RCP ranking method is not sensitive to the threshold $b$. In a large reasonable range, any value of $b$ can bring performance gains.

TABLE VII
COMPUTATION COMPLEXITY COMPARISON OF TWO ONE-STEP FRAMEWORKS. THE NUMBERS 128 AND 10 ARE THE NUMBERS OF RoIs FROM THE RPN AND ASSUMED FINAL PREDICTED PERSON BOUNDING BOXES, RESPECTIVELY.

| Computation (GFLOPs) | Part-1 | RoI Head | KTB | Part-2 | Total |
|---|---|---|---|---|---|
| input size | 960x600x3 | 7x7x512 | 256x128x3 | 32x16x512 | |
| OIM* | 42.2 | - | - | 128x2.9=371.2 | 413.4 |
| TDN | 42.2 | 128x0.05=6.4 | 10x2.4=24.0 | 10x2.9=29.0 | 101.6 |

It is recommended to set $b$ to 0.3, considering the performance on both datasets.

Fig. 6 shows the impact of $\lambda$. On the PRW dataset, when $\lambda \in [0.1, 0.3]$, the proposed RCP brings performance improvement. On the CUHK-SYSU dataset, any value of $\lambda \in [0.1, 0.9]$ can improve the performance. Experimental results show that the proposed RCP is not sensitive to $\lambda \in [0.1, 0.3]$ on two datasets. In this paper, we set $\lambda$ to 0.2 for both datasets.
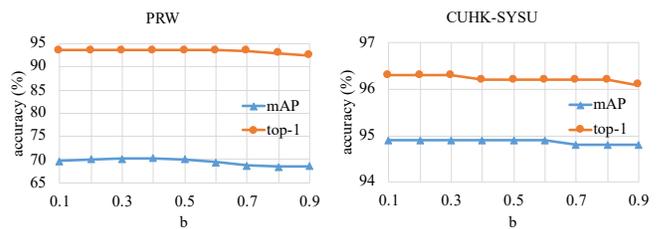


Fig. 5. Impact of parameter $b$. $\lambda$ is set to 0.2 for both datasets.
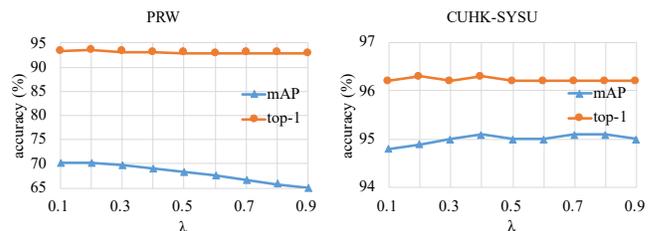


Fig. 6. Impact of parameter $\lambda$. $b$ is set to 0.3 for both datasets.

### F. Visualization

To further validate the effectiveness of the proposed TDN model as well as the RCP ranking method, some visualization results are reported in Fig. 7 and Fig. 8. Compared to the previous widely-used OIM*, the proposed TDN can retrieve more true candidates at the top of the ranking list. When the proposed RCP ranking method is applied, the quality of the ranking list is further improved. These visualization results demonstrate the effectiveness of the proposed TDN and RCP.

## V. CONCLUSION

In this paper, aiming to make the one-step person search model enjoy the merits of the state-of-the-art Re-ID models,
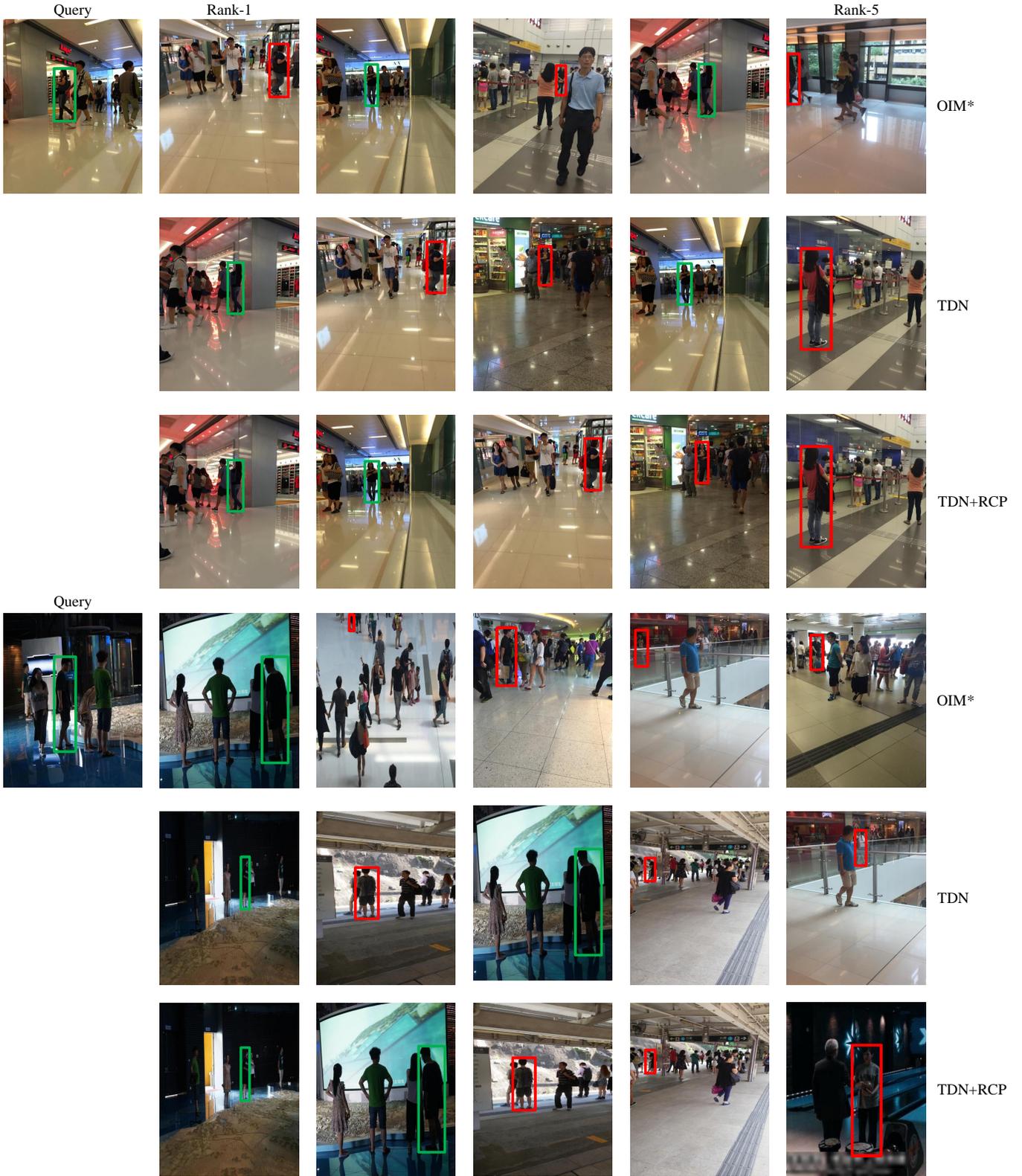
Fig. 7. Example results of two query persons. The query persons are marked with green bounding boxes in the first column scene images. For each query person, the first row, the second row and the third row ranking results correspond to the baseline method OIM*, the TDN method and the "TDN + RCP" method, respectively. In each ranking scene image, the most similar gallery person to the query person is marked within a bounding box. The green bounding box is for the right match of the query person, and the red one is for the wrong match.
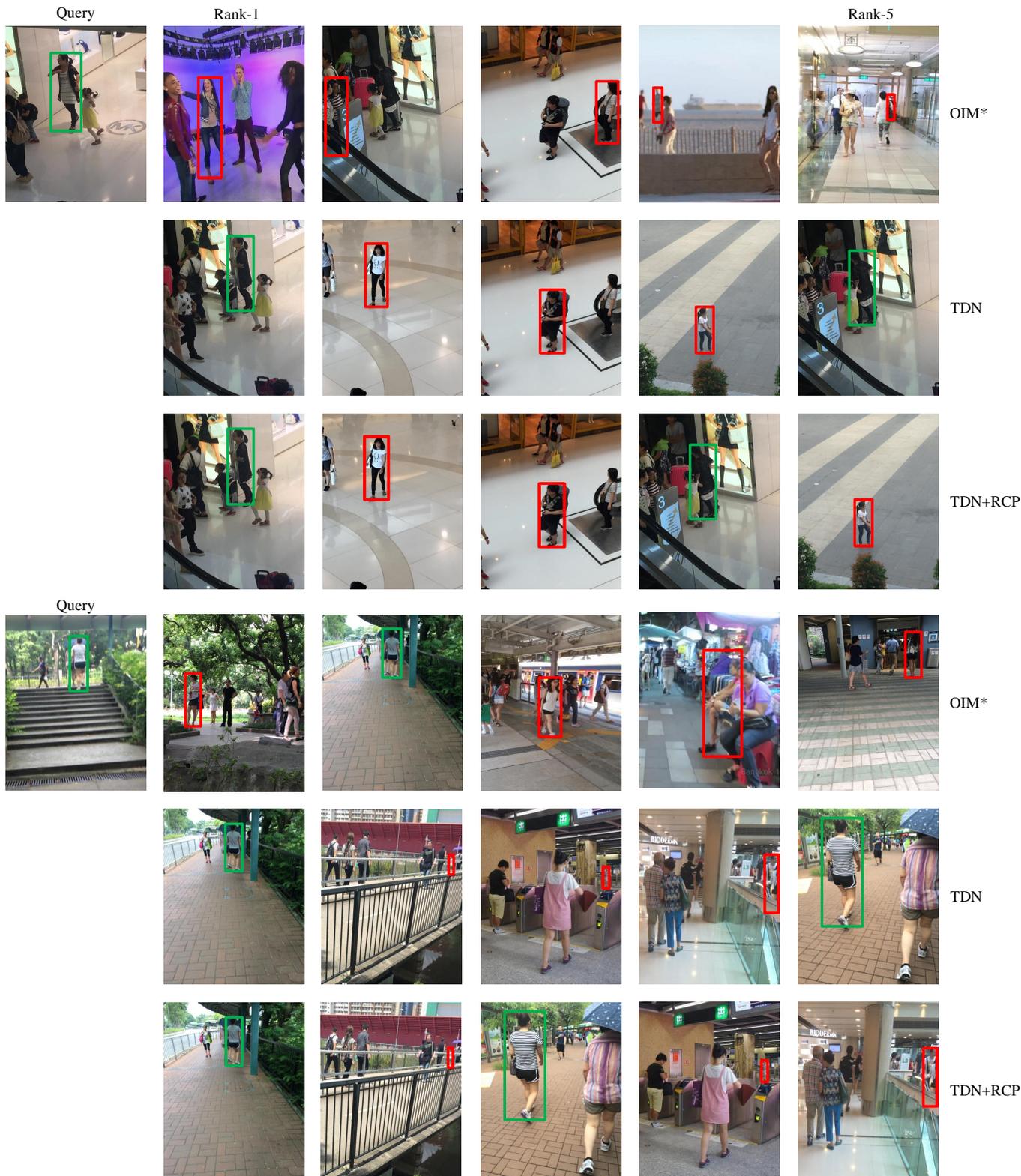
Fig. 8. Example results of another two query persons. The query persons are marked with green bounding boxes in the first column scene images. For each query person, the first row, the second row and the third row ranking results correspond to the baseline method OIM*, the TDN method and the "TDN + RCP" method, respectively. In each ranking scene image, the most similar gallery person to the query person is marked within a bounding box. The green bounding box is for the right match of the query person, and the red one is for the wrong match.

we propose the Teacher-guided Disentangling Networks which is a faster and stronger one-step person search framework. The proposed TDN presents a pipeline to integrate the prior knowledge from the powerful Re-ID model to boost the person search performance. In the proposed TDN, to make the one-step person search model learn discriminative Re-ID features from the powerful Re-ID teacher model better, we propose to partially disentangle the Partially Disentangled Framework and introduce the Knowledge Transfer Bridge module. Besides, we propose the Ranking with Context Persons strategy to exploit the context persons provided in person search task. The proposed RCP ranking method can further improve the person search performance. Comparison results with the previous state-of-the-art methods, as well as thorough ablation studies, demonstrate the favorable performance of the proposed method.

## REFERENCES

[1] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, 2015.
[2] D. Chen, S. Zhang, W. Ouyang, J. Yang, and Y. Tai, "Person search via a mask-guided two-stream cnn model," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 734–750.
[3] I. Goodfellow, Y. Bengio, A. Courville, and Y. Bengio, *Deep learning*. MIT press Cambridge, 2016, vol. 1.
[4] S. Bak and P. Carr, "One-shot metric learning for person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2990–2999.
[5] Z. Liu, D. Wang, and H. Lu, "Stepwise metric promotion for unsupervised video person re-identification," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2429–2438.
[6] A. Hermans, L. Beyer, and B. Leibe, "In defense of the triplet loss for person re-identification," *arXiv preprint arXiv:1703.07737*, 2017.
[7] W. Chen, X. Chen, J. Zhang, and 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.
[8] M. Yu, Z. Chang, Q. Zhou, S. Zheng, and T. P. Wu, "Reference-oriented loss for person re-identification," in *2019 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2019, pp. 1–8.
[9] Z. Chang, Q. Zhou, M. Yu, S. Zheng, H. Yang, and T.-P. Wu, "Distribution context aware loss for person re-identification," *arXiv preprint arXiv:1911.07273*, 2019.
[10] C. Yan, G. Pang, X. Bai, C. Liu, N. Xin, L. Gu, and J. Zhou, "Beyond triplet loss: person re-identification with fine-grained difference-aware pairwise loss," *IEEE Transactions on Multimedia*, 2021.
[11] D. Chen, D. Xu, H. Li, N. Sebe, and X. Wang, "Group consistent similarity learning via deep crf for person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8649–8658.
[12] D. Chen, H. Li, T. Xiao, S. Yi, and X. Wang, "Video person re-identification with competitive snippet-similarity aggregation and co-attentive snippet embedding," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1169–1178.
[13] G. Song, B. Leng, Y. Liu, C. Hetang, and S. Cai, "Region-based quality estimation network for large-scale person re-identification," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.
[14] Z. Chang, Z. Qin, H. Fan, H. Su, H. Yang, S. Zheng, and H. Ling, "Weighted bilinear coding over salient body parts for person re-identification," *Neurocomputing*, vol. 407, pp. 454–464, 2020.
[15] H. Luo, W. Jiang, Y. Gu, F. Liu, X. Liao, S. Lai, and J. Gu, "A strong baseline and batch normalization neck for deep person re-identification," *IEEE Transactions on Multimedia*, vol. 22, no. 10, pp. 2597–2609, 2019.
[16] X. Chen, C. Fu, Y. Zhao, F. Zheng, J. Song, R. Ji, and Y. Yang, "Salience-guided cascaded suppression network for person re-identification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 3300–3310.

[17] Z. Zhang, C. Lan, W. Zeng, X. Jin, and Z. Chen, "Relation-aware global attention for person re-identification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 3186–3195.
[18] X. Gong, Z. Yao, X. Li, Y. Fan, B. Luo, J. Fan, and B. Lao, "Lag-net: Multi-granularity network for person re-identification via local attention system," *IEEE Transactions on Multimedia*, 2021.
[19] L. Zheng, H. Zhang, S. Sun, M. Chandraker, Y. Yang, and Q. Tian, "Person re-identification in the wild," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1367–1376.
[20] X. Lan, X. Zhu, and S. Gong, "Person search by multi-scale matching," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 536–552.
[21] C. Han, J. Ye, Y. Zhong, X. Tan, C. Zhang, C. Gao, and N. Sang, "Re-id driven localization refinement for person search," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 9814–9823.
[22] C. Wang, B. Ma, H. Chang, S. Shan, and X. Chen, "Tcts: A task-consistent two-stage framework for person search," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11 952–11 961.
[23] S. Ren, K. He, R. Girshick, and 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.
[24] T. Xiao, S. Li, B. Wang, L. Lin, and 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.
[25] L. Li, H. Yang, and L. Chen, "Spatial invariant person search network," in *Chinese Conference on Pattern Recognition and Computer Vision (PRCV)*. Springer, 2018, pp. 122–133.
[26] B. Munjal, S. Amin, F. Tombari, and F. Galasso, "Query-guided end-to-end person search," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 811–820.
[27] Y. Hong, H. Yang, L. Li, L. Chen, and C. Liu, "A cascaded multitask network with deformable spatial transform on person search," *International Journal of Advanced Robotic Systems*, vol. 16, no. 3, p. 1729881419858162, 2019.
[28] W. Dong, Z. Zhang, C. Song, and T. Tan, "Bi-directional interaction network for person search," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 2839–2848.
[29] D. Chen, S. Zhang, J. Yang, and B. Schiele, "Norm-aware embedding for efficient person search," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 12 615–12 624.
[30] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio, "Fitnets: Hints for thin deep nets," *arXiv preprint arXiv:1412.6550*, 2014.
[31] S. Zagoruyko and N. Komodakis, "Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer," *arXiv preprint arXiv:1612.03928*, 2016.
[32] J. Yim, D. Joo, J. Bae, and J. Kim, "A gift from knowledge distillation: Fast optimization, network minimization and transfer learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 4133–4141.
[33] B. Munjal, F. Galasso, and S. Amin, "Knowledge distillation for end-to-end person search," *arXiv preprint arXiv:1909.01058*, 2019.
[34] X. Pan, P. Luo, J. Shi, and X. Tang, "Two at once: Enhancing learning and generalization capacities via ibn-net," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 464–479.
[35] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2961–2969.
[36] R. Mazzon, F. Poiesi, and A. Cavallaro, "Detection and tracking of groups in crowd," in *2013 10th IEEE International Conference on Advanced Video and Signal Based Surveillance*. IEEE, 2013, pp. 202–207.
[37] J. Xiao, Y. Xie, T. Tillo, K. Huang, Y. Wei, and J. Feng, "Ian: the individual aggregation network for person search," *Pattern Recognition*, vol. 87, pp. 332–340, 2019.
[38] H. Liu, J. Feng, Z. Jie, K. Jayashree, B. Zhao, M. Qi, J. Jiang, and S. Yan, "Neural person search machines," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 493–501.
[39] X. Chang, P.-Y. Huang, Y.-D. Shen, X. Liang, Y. Yang, and A. G. Hauptmann, "Rcaa: Relational context-aware agents for person search," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 84–100.

[40] Y. Yan, Q. Zhang, B. Ni, W. Zhang, M. Xu, and X. Yang, "Learning context graph for person search," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2158–2167.

[41] Y. Zhong, X. Wang, and S. Zhang, "Robust partial matching for person search in the wild," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 6827–6835.

[42] W. Dong, Z. Zhang, C. Song, and T. Tan, "Instance guided proposal network for person search," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 2585–2594.

[43] H. Kim, S. Joung, I.-J. Kim, and K. Sohn, "Prototype-guided saliency feature learning for person search," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 4865–4874.

[44] G. Wang, Y. Yuan, X. Chen, J. Li, and X. Zhou, "Learning discriminative features with multiple granularities for person re-identification," in *Proceedings of the 26th ACM international conference on Multimedia*, 2018, pp. 274–282.

[45] M. Ye, J. Shen, G. Lin, T. Xiang, L. Shao, and S. C. Hoi, "Deep learning for person re-identification: A survey and outlook," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.