

Towards Fully Decoupled End-to-End Person Search

Pengcheng Zhang

School of Computer Science and Engineering,
State Key Laboratory of Software Development Environment,
Jiangxi Research Institute, Beihang University
Beijing, China
pengchengz@buaa.edu.cn

Jin Zheng

School of Computer Science and Engineering,
State Key Laboratory of Software Development Environment,
Jiangxi Research Institute, Beihang University
Beijing, China
jinzheng@buaa.edu.cn

Xiao Bai*

School of Computer Science and Engineering,
State Key Laboratory of Software Development Environment,
Jiangxi Research Institute, Beihang University
Beijing, China
baixiao@buaa.edu.cn

Xin Ning

Institute of Semiconductors, Chinese Academy of Sciences,
Cognitive Computing Technology Joint Laboratory,
Wave Group,
Beijing, China
ningxin@semi.ac.cn

Abstract—End-to-end person search aims to jointly detect and re-identify a target person in raw scene images with a unified model. The detection task unifies all persons while the re-id task discriminates different identities, resulting in conflict optimal objectives. Existing works proposed to decouple end-to-end person search to alleviate such conflict. Yet these methods are still sub-optimal on one or two of the sub-tasks due to their partially decoupled models, which limits the overall person search performance. In this paper, we propose to fully decouple person search towards optimal person search. A task-incremental person search network is proposed to incrementally construct an end-to-end model for the detection and re-id sub-task, which decouples the model architecture for the two sub-tasks. The proposed task-incremental network allows task-incremental training for different objectives thus fully decoupled the model for persons search. Comprehensive experimental evaluations demonstrate the effectiveness of the proposed fully decoupled models for end-to-end person search.

Index Terms—person search, decoupling, task-incremental learning

I. INTRODUCTION

Person search [1] jointly performs two sub-tasks, *i.e.* person detection [2]–[4] and re-id [12], [23], [24], [36], to locate a query person across a gallery of uncropped scene images. To guarantee the overall person search accuracy, it requires optimal detection performance to integrally contain all true positive persons, and optimal re-id performance to discriminate persons of different identities.

Recent researches focus more on end-to-end methods [16], [18]–[20], [35], [37], [38] that complete person search with a unified model. Although this paradigm shows better efficiency than the two-step ones [14], [22], [29], [32]–[34], it suffers from the conflicting objectives of the two sub-tasks as the detection task unifies all persons while the re-id task discriminates different identities. For illustration purposes, we

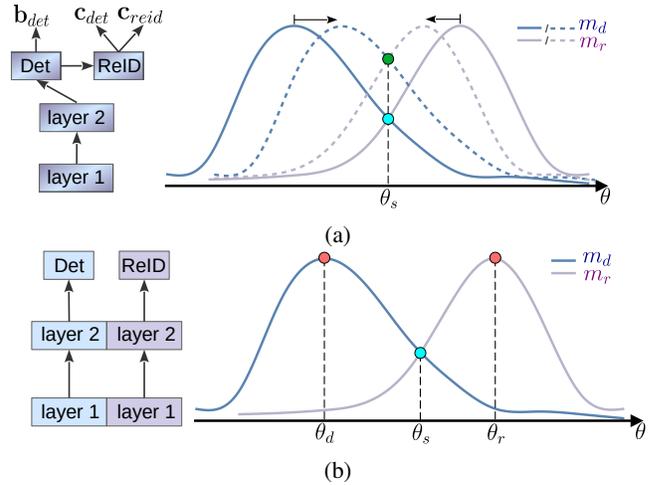


Fig. 1: Comparison of fully decoupled person search (b) with previous decoupled models [16], [19] (a). We employ cyan points to indicate the performance of the vanilla end-to-end model [18]. **(a) Left:** Partially decoupled model in [16], [19]. **Right:** The green point indicates the performance of the partially decoupled model. By closing the two task-specific feature spaces, the upper bound of performance upon shared parameters θ_s is boosted. **(b) Left:** The proposed fully decoupled person search network. **Right:** The pink points illustrate the performance of the proposed model. It eliminates the coupled parameters and achieves the optimum for both sub-tasks towards optimal person search.

denote by θ_s the model parameters shared the two tasks, and show the overall performance of a vanilla end-to-end model [18] in Figure 1a where m_d and m_r denote the individual performance (larger is better) of the detection and re-id sub-task upon model parameters θ . Due to that θ_s is the majority of model parameters, it obtains only a compromised average

*Corresponding author.

solution for person search. To alleviate this problem, previous works explored partially decoupled person search by closing the two task-specific feature spaces [16], [19] and boosting the performance upper-bound upon θ_s as in Figure 1a.

Another solution is to separate the respective prediction branches [17], [20] for different tasks, which constraints the impact of θ_s and decouples the prediction parameters to improve sub-task performances.

Despite the progress made for decoupled person search, these methods are still sub-optimal on one or two of the sub-tasks due to the models are still partially coupled. This further limits the overall person search performances. To this end, we propose to fully decouple end-to-end models, as is shown in Figure 1b, to achieve the optimum for both tasks towards the optimal solution for person search. The advantages of fully decoupled end-to-end person search are threefold: (1) The detection sub-network is comparable with standalone detectors, which maximumly reduces under-detected target persons in variant scene images; (2) Person retrieval features are learned without compromising to the detection sub-task, which guarantees the feature discrimination; (3) The architecture of the re-id sub-network is less dependent on the detection sub-network, which unleashes more space to design specialized re-id modules.

Inspired by the task-incremental learning (TIL) [7], [8], [40], [43], [49] mechanism, we propose a task-incremental network to enable the aforementioned fully decoupled end-to-end person search. With the assumption that the task identity is known during inference, TIL accumulates task-dependent parameters during training and selects proper ones at test time to mitigate catastrophic forgetting for incremental learning [39]. In contrast, for task-incremental person search, we construct an expandable model which is derived from a standard person detector and expanded for an incremental task, *i.e.* the re-id sub-task. The overall model architecture is thus decoupled for the two sub-tasks.

The task-incremental person search network also degrades to a partially decoupled model when jointly trained by the two sub-tasks. We thus propose to conduct task-incremental training to decouple the training procedure. We first train a standard person detector and then expand the model for the re-id task. The whole model is then trained only by the re-id sub-task while the detection sub-network is frozen. The training for the re-id sub-task introduces spatial noises on GT boxes to simulate overlapping detection bounding boxes in previous works [16], [18], [19], [21], [35], [37] without performing computationally expensive person detection. Therefore, the model is fully decoupled for the two conflicting sub-tasks. This paradigm introduces an extra training phase, yet it keeps the end-to-end efficiency for inference. And a portion of the added re-id modules runs in parallel with the detection sub-network, which seldom increases the time cost.

In summary, this paper makes the following contributions:

- We propose the first fully decoupled end-to-end person search model. By designing a novel task-incremental per-

son search network, we decouple the model architecture for the two conflicting sub-tasks.

- The proposed task-incremental person search network allow task-incremental training for the two sub-tasks. This enables independent learning for the conflicting task objectives. Thus the proposed method achieves the optimum for the two conflicting sub-tasks towards optimal end-to-end person search.
- Comprehensive experiments demonstrate that the proposed model significantly outperforms previous decoupled models on PRW [22] and achieves the best on CUHK-SYSU [18].

II. RELATED WORK

Person Search. Person search methods can be categorized into two types: two-step methods and end-to-end methods. Two-step methods typically employ a standalone detector to detect and crop person images, and a re-id model to retrieve a target across the cropped images. [22] first explores combining popular person detector and re-id models for person search. To obtain more representative features for each identity, [29] proposes to enhance foreground image regions and [28] performs multi-scale matching for person search. [14], [33], [34] instead design target-guided person detectors to suppress retrieval distractors. To improve the efficiency of the two-step paradigms, end-to-end methods propose to perform person search by a unified model. [18] constructs the first end-to-end model and proposes the Online Instance Matching (OIM) [18] loss for model training. Following [18], [30], [31], [35], [50] significantly improve the performances of end-to-end models. [37], [38] draw inspiration from the Transformers [6], [51] and obtain more discriminative person features with well-designed models.

One key challenge for end-to-end person search is the contradictory objectives of person detection and re-id. To deal with the problem, [16] proposes hierarchical classification score calculation and [19] disentangles the output embedding for the two tasks. These close the gap between two marginal task-specific feature spaces. Another solution is to construct separate prediction branches as in [17], which decouples the predictions for different tasks. [20] proposes a "re-id first" design and can be viewed as employing a non-parametric identity re-id prediction branch for this. Though achieved significant progress, these models are still sub-optimal on one or two of the sub-tasks due to that a subset of model parameters are left coupled for the two sub-tasks. This limits the overall person search performances of the models. To this end, this paper proposes to fully decouple end-to-end models towards optimal person search.

Task-incremental Learning. TIL [39] is one of the major scenarios for incremental learning. A key assumption for TIL is that the task identity is known during inference. The catastrophic forgetting problem is thus can be solved by incrementally learning task-specific sub-networks while preserving all learned ones. At test time, the task identity is then employed to select the proper sub-network. Specifically,

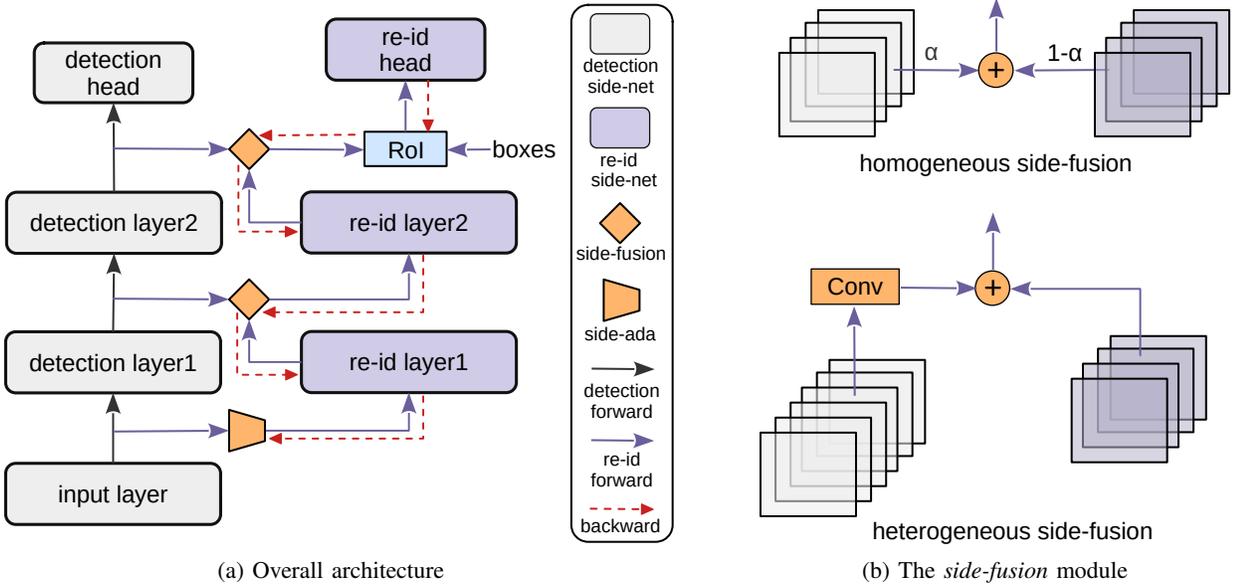


Fig. 2: (a) The proposed fully decoupled person search network which consists of a *detection side-net*, a *re-id side-net*, and the modules that bridge them. The model is trained incrementally by person detection and re-id tasks, which fully decouples the parameters for the two conflicting sub-tasks. (b) The architectures of homogeneous and heterogeneous *side-fusions* for their corresponding *re-id side-net*. These modules transfer knowledge from the trained *detection side-net* to the *re-id side-net*.

[40]–[43] directly expands the network by adding new layers or branches for new tasks. Yet this can be limited in practice due to unbounded parameter growth. Other works thus propose to freeze partial network with masks for old tasks [44]–[47] and adapt to the new task with left trainable parameters, which may suffer from running out of model capacities. To mitigate the limitation that the task identity is not always available, [48], [49] additionally trains a task classifier to infer task identity or designs adaptive parameter selection mechanisms at test time. With similar techniques, [7], [8] also suits TIL by incrementally expandable side-networks.

In this work, we design a task-incremental network to enable the aforementioned fully decoupled end-to-end person search. This achieves the optimum for both the two sub-tasks to facilitate optimal person search performance.

III. METHOD

A. Task-incremental Person Search Network

Detection side-net. The *detection side-net* f_d can be various modern detectors, *e.g.* Faster R-CNN [2], RetinaNet [4] and FCOS [3], given the simple and expandable overall architecture. For illustration purposes, we employ the Faster R-CNN [2] detector as f_d in this section. The *input layer* in Figure 2a is composed of the ‘conv1’ and ‘conv2’ blocks of the ResNet [9] backbone. The *detection layer1* and *detection layer2* are the ‘conv3’ and ‘conv4’ blocks, respectively. The *detection head* consists of the RPN [2] and ‘conv5’ block to predict probable person locations and corresponding classification scores. Following the common practice in object detection [2]–[4], [6], the *input layer* is initialized with ImageNet [10] pre-trained parameters and frozen during training, which makes it

independent neither on the detection sub-task nor on the re-id sub-task.

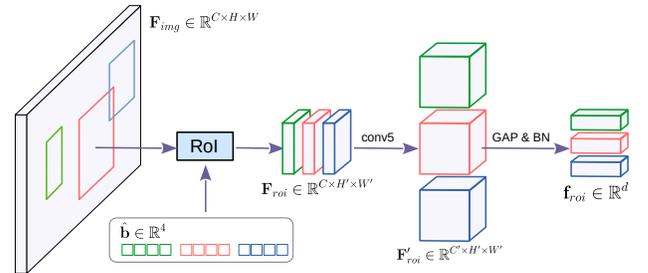


Fig. 3: Illustration of the *re-id head*. Person feature maps are drawn from the output of ‘conv4’ and refined by the ‘conv5’ block. By consecutive global average pooling and batch normalization, this module produces 1-D person feature vectors.

Re-id side-net. As is in Figure 2a, the *re-id side-net* f_r shares the *input layer* with f_d . Similar to f_d , another ‘conv3’ block and ‘conv4’ block are employed as the consecutive *re-id layers* to extract image feature maps $F_{img} \in \mathbb{R}^{C \times H \times W}$. Given probable person bounding boxes $\hat{\mathbf{b}} \in \mathbb{R}^d$, as is in Figure 3, RoIAlign [11] is performed to obtain corresponding person feature maps $F_{roi} \in \mathbb{R}^{C \times H' \times W'}$. We further employ the ‘conv5’ block (stride is set to 1) of the backbone on top of F_{roi} to produce $F'_{roi} \in \mathbb{R}^{C' \times H' \times W'}$. Subsequent global average pooling (GAP) and batch normalization (BN) [13] layers then aggregate F'_{roi} to 1-D person features $\mathbf{f}_{roi} \in \mathbb{R}^d$ for person retrieval.

Side-ada. The network architecture of f_r can be homo-

geneous (e.g., ResNet50 for both) or heterogeneous (e.g., ResNet50 for f_d and ResNet34 for f_r) to f_d . The *side-ada* module is thus inserted to transform the output of the *input layer* to match the input size required by f_r . For homogeneous *re-id side-net*, the *side-ada* module is implemented by an identity function $I(\cdot)$ where

$$I(\mathbf{x}) = \mathbf{x}, \mathbf{x} \in \mathbb{R}^{c' \times h' \times w'}. \quad (1)$$

And for heterogeneous *re-id side-net*, we instead employ a single convolution layer

$$\text{Conv}(\mathbf{x}) = \mathbf{x}', \quad (2)$$

where $\mathbf{x} \in \mathbb{R}^{c \times h \times w}$ and $\mathbf{x}' \in \mathbb{R}^{c' \times h' \times w'}$, to transform the input to be spatially compatible with f_r .

Side-fusion. Inspired by [7], [8], we further add the *side-fusion* module to fuse \mathbf{x}_d^i and \mathbf{x}_r^i , the outputs of two i th parallel side-network layers, as input \mathbf{x}_s^i to the subsequent *re-id side-net* block. This transfers valuable knowledge from the pre-trained f_d to f_r . Analogously, we design homogeneous *side-fusion* and heterogeneous *side-fusion* modules for homogeneous f_r and heterogeneous f_r , respectively. As Figure 2b shows, the homogeneous *side-fusion* performs alpha blending

$$\mathbf{x}_s^i = \alpha_i \mathbf{x}_d^i + (1 - \alpha_i) \mathbf{x}_r^i, \quad (3)$$

where \mathbf{x}_s^i , \mathbf{x}_d^i and \mathbf{x}_r^i are all in $\mathbb{R}^{c'_i \times h'_i \times w'_i}$, to fuse the outputs. $\alpha_i \in [0, 1]$ is the only learnable parameter of this module. Similar to the *side-ada* module, the heterogeneous *side-fusion* module employs a single convolution layer to transform \mathbf{x}_d^i to be spatially compatible with \mathbf{x}_r^i . The fused output is thus given by

$$\mathbf{x}_s^i = \text{Conv}(\mathbf{x}_d^i) + \mathbf{x}_r^i \quad (4)$$

where \mathbf{x}_s^i , $\mathbf{x}_r^i \in \mathbb{R}^{c'_i \times h'_i \times w'_i}$ and $\mathbf{x}_d^i \in \mathbb{R}^{c_i \times h_i \times w_i}$.

B. Task-incremental Model Training

The proposed task-incremental person search network in Section III-A is still partially decoupled when conducting simple joint training for person search. To this end, we propose task-incremental training for the task-incremental person search network. Specifically, we first train a standard person detector f_d on the training set with detection losses according to the specific settings of the detector. Afterwards, we freeze the detection side network f_d and expand the model with f_r and the bridge modules. The whole model is then trained only by the re-id task with the OIM loss [18] \mathcal{L}_{oim} and the triplet loss \mathcal{L}_{tri} similar to [20], which makes f_r learn discriminative person representations without competing with f_d . The overall model is thus fully decoupled in a task-incremental learning manner for person search.

Note that previous works [18]–[20], [35], [37] detect multiple overlapping bounding boxes of the same person and extract feature embeddings upon the boxes to train the re-id modules, which implicitly performs spatial augmentation

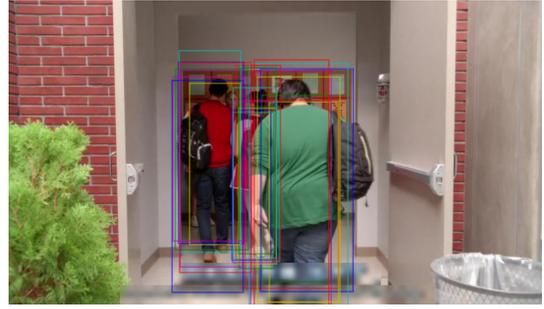


Fig. 4: Illustration of augmented person bounding boxes from the GT boxes. Center shifting and box scaling are employed to add spatial noises.

for representation learning. To decouple the training of the re-id side-network from the detection side-network without sacrificing this advantage, we propose to add spatial noises, namely Spatial-noise Augmentation (SnA), to the GT boxes to obtain augmented training samples for f_r . Similar to [15], we combine two types of spatial noises: center shifting and box scaling. Denoted by (c_x, c_y) and (h_b, w_b) the center point and spatial size of GT box \mathbf{b} , center shifting adds random noise $(\Delta c_x, \Delta c_y)$ to the center, where $|\Delta c_x| < \frac{\lambda_1 w_b}{2}$ and $|\Delta c_y| < \frac{\lambda_1 h_b}{2}$. And box scaling randomly samples height and width from $[(1 - \lambda_2) h_b, (1 + \lambda_2) h_b]$ and $[(1 - \lambda_2) w_b, (1 + \lambda_2) w_b]$. $\lambda_1 \in (0, 1)$ and $\lambda_2 \in (0, 1)$ are hyper-parameters that control the noise scales of the augmentations. For training, as in Figure 4, we randomly generate n augmented versions of each GT box to enhance the robustness of f_r .

IV. EXPERIMENTS

A. Datasets

CUHK-SYSU [18] presents 18,184 images from both movies and street snapshots. A total number of 96,143 pedestrian bounding boxes and 8,432 labeled identities are manually annotated for person search. The training subset provides 11,206 frames with 5532 identities from both sources. The testing subset selects 2900 query identities with gallery sizes varying from 50 to 4,000. The training and testing sets have no overlaps on images or identities.

PRW [22] deployed 6 cameras at a campus to record multi-view videos of pedestrians. 11,816 frames are selected and densely annotated as 43,110 pedestrian bounding boxes with 932 identities. The training subset contains 5,134 frames with 432 identities, while the rest 6,112 frames make up the test subset. The evaluation protocol of PRW takes the full test set as the gallery for evaluation by default, which tends to be more challenging than CUHK-SYSU.

Evaluation Metrics. Similar to the evaluation metrics in person re-identification [25], [26], the mAP and top-1 accuracy are utilized in person search. During the evaluation, the gallery of probable persons is dynamically built upon detection results. A retrieved person is considered positive if it shares the same identity with the query and the detected bounding box has the

Intersection over Union (IOU) larger than 0.5 with the GT, which makes the mAP and top-1 accuracy also be affected by person detection results. For the evaluation of person detection, the Average Precision (AP) and Recall [27] are the mostly employed metrics in recent person search works.

B. Implementation Details

For training, we first train f_d by person detection task for 18 epochs with batch size of 8. Then we freeze the parameters of f_d and train the whole model by the re-id sub-task for 18 epochs with batch size of 5. The output size of RoIAlign is set to 16×8 for both f_d and f_r . The input image is randomly scaled for training and resized to 1500×900 for testing. We use SGD optimizer with initial learning rate of 0.003. The learning rate is linearly warmed up during the first epoch. For f_d , we decrease the learning rate by 10 at the 12th epoch. And for f_r , the learning rate is decreased by 10 at the 10th epoch. For CUHK-SYSU/PRW, the circular queue size of OIM is set to 5000/500 and the softmax temperature is set to 1/30.

In the following experiments, we employ Faster R-CNN with ResNet50 backbone as f_d and ResNet50 blocks initialized with ImageNet pre-trained parameters as f_r unless otherwise specified. We use \dagger to mark models evaluated with Class Weighted Similarity (CWS) [22]. The subscript $_{34}$ indicates that ResNet34 blocks are employed for f_d and f_r . Evaluations on CUHK-SYSU are with gallery size of 100 by default. The models are all trained and tested on a single RTX 3090 GPU. Extensive experimental results and implementation details are further presented in the Supplementary Material.

C. Analytical Studies

Detectors	PRW		CUHK-SYSU	
	mAP	top-1	mAP	top-1
RetinaNet [4]	47.6	86.3	92.8	93.8
FCOS [3]	47.8	86.7	93.0	94.3
Faster R-CNN [2]	47.9	86.5	92.1	93.4
RetinaNet $^\nabla$ [4]	22.1	55.8	73.6	75.2
FCOS $^\nabla$ [3]	22.4	56.4	76.9	77.8
Faster R-CNN $^\nabla$ [2]	21.0	47.2	74.2	77.6
RetinaNet † [4]	46.8	85.9	93.1	93.8
FCOS † [3]	46.1	85.5	92.6	92.9
Faster R-CNN † [2]	49.0	87.1	93.4	<u>94.2</u>

TABLE I: The person search evaluation results of our proposed model with different detectors as f_d . We use $^\nabla$ to denote the coupled models.

Different choices of f_d . The proposed fully decoupled person search network specifies no concrete architecture of f_d . To verify the effect of different detector architectures, we employ RetinaNet [4], FCOS [3] and Faster R-CNN [2], as f_d for our proposed model and present the evaluation results in Table I. On PRW, the proposed model achieves similar performance with different detectors. And on CUHK-SYSU, employing RetinaNet or FCOS as f_d slightly outperforms that with Faster R-CNN. These results demonstrate that the proposed fully decoupled network is well-compatible with

various detector architectures. We additionally test their respective coupled versions, *i.e.*, models jointly trained by the two sub-tasks and directly share the features between the detection prediction layers and the re-id losses. The proposed fully decoupled models also outperform the coupled versions by a large margin. We further conduct the evaluation with CWS [22] on the models as in [19], [35]. This significantly boosts the person search performance when employing Faster R-CNN as f_d yet harms that with other detectors. We present a more comprehensive analysis upon this in the Supplementary Material.

Combination of f_d and f_r	PRW		CUHK-SYSU	
	mAP	top-1	mAP	top-1
R-50 & R-50	47.9	86.5	92.1	93.4
R-34 & R-34	<u>47.1</u>	<u>85.5</u>	<u>92.0</u>	93.3
R-50 & R-34	41.6	83.7	91.0	92.1
R-50 & OSNet [36]	40.9	84.7	88.6	90.2

TABLE II: Comparison between different combinations of f_d and f_r . The first two rows are homogeneous f_r s while others are heterogeneous f_r s. We use R-50 and R-34 to denote ResNet50 and ResNet34, respectively.

Comparison between homogeneous f_r and heterogeneous f_r . As is described in Section III-A, f_r can be with homogeneous or heterogeneous architecture to f_d . To investigate the effects of different combinations of f_d and f_r , we conduct experiments as in Table II. We initialize the f_r s with their respective ImageNet pre-trained parameters for fair comparison. It can be observed that the performances of the heterogeneous combinations are marginally inferior to the homogeneous counterparts. The main reason for this can be the limited capacity of f_r or the misalignment between heterogeneous blocks. We thus additionally test a homogeneous combination of ResNet34 for both f_d and f_r . Although the capacity of this combination is also limited, it only performs moderately inferior to the ResNet50 combination, which demonstrates that the misalignment between heterogeneous blocks mainly impedes the performance. We thus employ homogeneous side-networks in the proposed model by default.

Models	PRW		CUHK-SYSU	
	mAP	top-1	mAP	top-1
ours w/ <i>side-fusion</i>	47.9	86.5	92.1	<u>93.4</u>
ours w/o <i>side-fusion</i>	46.1	85.5	91.1	92.4
ours $_{34}$ w/ <i>side-fusion</i>	<u>47.1</u>	<u>85.5</u>	<u>92.0</u>	93.3
ours $_{34}$ w/o <i>side-fusion</i>	44.4	84.7	91.0	92.1

TABLE III: Person search performances of the proposed models with and without *side-fusion*. Note that the tested *side-fusion* modules are homogeneous *side-fusion* modules.

The effect of *side-fusion*. We propose the *side-fusion* modules to transfer useful knowledge, *e.g.* suppressing the background and enhancing the foreground, from the trained f_d to f_r . The proposed model is also capable of performing person search without the *side-fusion* modules. We thus

evaluate the person search performances of models with and without *side-fusion*, as in Table III, to verify the impact of the modules. It can be observed that the *side-fusion* modules consistently boost the person search accuracy, suggesting that the knowledge learned from the detection sub-task can be implicitly employed to strengthen the re-id side-network.

Comparison between joint training and task-incremental training. To enable the fully decoupled person search, we propose a task-incremental person search network and train the model in a task-incremental manner. It is worth noting that the proposed task-incremental person search network becomes a partially decoupled model when jointly trained for the two sub-tasks. Specifically, the shared *input layer* is fixed with ImageNet pre-trained parameters and thus independent of the two sub-tasks. The *side fusion* modules are extremely lightweight, which constraints the coupled optimization for the contradictory objectives. And the prediction modules are separated for their respective tasks similar to [17]. To validate the effect of task-incremental training, we conduct performance comparisons between joint training and task-incremental training. As the results in Table IV and Table V show, the model also obtains comparable performances to previous decoupled methods when jointly trained by the two sub-tasks. And task-incremental training further boosts the person search performance as well as the person detection performance. We also test a hybrid training scheme, *i.e.* jointly training f_d with a very small $lr = 1.0 \times 10^{-5}$ when training f_r . The results in Table IV and Table V suggest that this slightly improves the person search performance yet significantly impedes the person detection capability. And the memory cost is also increased when jointly training the two side-networks.

Training	PRW		CUHK-SYSU	
	mAP	top-1	mAP	top-1
joint	47.3	86.1	91.3	92.7
task-incremental	47.9	86.5	92.1	93.4
hybrid	48.0	86.2	92.2	93.5

TABLE IV: Person search performance comparison between joint training and task-incremental training.

Training	PRW		CUHK-SYSU	
	AP	Recall	AP	Recall
joint	92.8	96.8	87.0	93.3
task-incremental	93.4	97.6	87.8	94.0
hybrid	90.4	95.1	82.7	86.1

TABLE V: Person detection performance comparison between joint training and task-incremental training.

D. Comparison with State-of-the-art

To demonstrate the advantages of fully decoupled person search networks, we first compare our proposed model with existing decoupled person search models in terms of person search (Table VI) and person detection (Table VII) performances. We then present the comparison results with recent

Method	PRW		CUHK-SYSU	
	mAP	top-1	mAP	top-1
HOIM [16]	39.8	80.4	89.7	90.8
NAE [†] [19]	43.3	80.9	91.5	92.4
NAE+ [†] [19]	44.0	81.1	92.1	92.9
ours ₃₄ [†]	47.3	85.9	92.6	93.4
ours [†]	49.0	87.1	93.4	94.2
DMRNet [17]	46.1	83.2	91.6	93.0
ours ₃₄ (RetinaNet)	46.5	84.9	92.0	93.1
ours (RetinaNet)	47.6	86.3	92.8	93.8
AlignPS [20]	45.9	81.9	93.1	93.4
AlignPS+ [20]	46.1	82.1	94.0	94.5
ours ₃₄ (FCOS)	46.3	85.4	92.6	93.4
ours(FCOS)	47.8	86.7	93.0	94.3

TABLE VI: Person search performance comparison with existing decoupled person search methods. We employ the results with RetinaNet of [17] for fair comparison.

well-established state-of-the-art methods (Table VIII). And qualitative visualization is shown at last.

Results on CUHK-SYSU. As is shown in Table VI, with similar configurations of detectors, our proposed models significantly outperform that in [16], [17], [19], [20] on CUHK-SYSU. By applying decoupled model initialization, the proposed model achieves the best top-1 accuracy and mAP score. Note that AlignPS+ [20] utilizes deformable convolution backbone [53] and multi-scale features which are not included in our models. We further evaluate the detection performances of previous decoupled models and our proposed ones. By fully decoupling for the detection and re-id sub-tasks, the *detection side-net* f_d inherits the capability of independently trained detectors. The detection sub-network of our proposed model thus guarantees the optimum for person detection and significantly outperforms previous methods. When compared with recent well-established models, as in Table VIII, the proposed model achieves competitive mAP and top-1 scores without bells and whistles. Compared with the best model COAT [37] that employs self-attention blocks [51] and cascaded refinement [54] of multi-scale features, the proposed model achieves comparable performance with simple model architecture and single-scale features. We also test the combination of our proposed full decoupled person search framework with the modules in [37]. This achieves superior performances on CUHK-SYSU, demonstrating the effectiveness of our proposed method.

Results on PRW. On the PRW dataset, as Table VI shows, our proposed models surpass the best of previous decoupled methods by a large margin even with lightweight ResNet34 blocks. We also test the detection performances of models in [16], [19], [20] by their released checkpoints and present the results in Table VII. It can be observed that our proposed models consistently achieve better performances than these methods. When compared with recent state-of-the-art end-to-end methods, the proposed model achieves the second-best mAP and the second-best top-1 accuracy. On the more challenging multi-view gallery of PRW, the proposed model

Detector	PRW		CUHK-SYSU	
	AP	Recall	AP	Recall
Faster R-CNN by [16]	87.8	95.7	85.7	91.8
Faster R-CNN by [19]	88.8	93.3	86.8	92.6
ours Faster R-CNN	93.4	97.6	87.8	94.0
RetinaNet by [17]	-	-	91.3	-
ours RetinaNet	93.0	95.6	91.7	97.5
FCOS by [20]	88.4	90.5	86.9	89.1
FCOS ⁺ by [20]	89.1	91.1	86.0	88.8
ours FCOS	93.4	95.8	92.2	95.5

TABLE VII: Person detection performance comparison with existing decoupled person search methods. ⁺ indicates the model with DCN [53] backbone.

Methods		PRW		CUHK-SYSU	
		mAP	top-1	mAP	top-1
two-step	IDE [†] [22]	20.5	48.3	-	-
	MGTS [29]	32.6	72.1	83.0	83.7
	CLSA [28]	38.7	65.0	87.2	88.5
	RDLR [32]	42.9	70.2	93.0	94.2
	IGPN [33]	<u>47.2</u>	87.0	90.3	91.4
	TCTS [34]	46.8	<u>87.5</u>	<u>93.9</u>	<u>95.1</u>
end-to-end	OIM [18]	21.3	49.4	75.5	78.7
	CTXG [31]	33.4	73.6	86.5	84.1
	BiNet [52]	45.3	81.7	90.0	90.7
	PGA [50]	44.2	85.2	92.3	94.7
	SeqNet [†] [35]	45.8	81.7	93.4	94.1
	OIMNet++ [21]	47.7	84.8	93.1	94.1
	PSTR [38](R-50)	49.5	87.8	93.5	<u>95.0</u>
	COAT [†] [37]	<u>53.3</u>	87.4	<u>94.2</u>	94.7
	ours [†]	49.0	87.1	93.4	94.2
	ours + COAT [†] [37]	54.0	<u>87.5</u>	94.3	95.2

TABLE VIII: Comparison with other state-of-the-art methods. We also evaluate the models on the multi-view gallery of PRW as in [37] and present the results in the last 4 rows.

also obtain the second-best mAP and the best top-1 accuracy. We observe that [37], [38] construct sophisticated attention blocks [6], [51] and multi-scale feature representations, which can be further incorporated to boost the performance of our proposed model. When combining our proposed method with [37], our proposed model achieves the best mAP and the second best top-1 accuracy.

Efficiency comparison. To verify the efficiency of the fully decoupled person search, we compare the proposed models with previous decoupled models by training time (h), number of parameters (M), and run time (ms), in Table IX. All results are conducted on the PRW dataset by a single RTX 3090 GPU. It can be observed that the proposed models are with more parameters. Yet the inference time is barely increased due to the highly parallel model architecture. Although the task-incremental person search network introduces another training stage, the fully decoupled mechanism makes the training of both sub-tasks converge faster, resulting in an acceptable total training time.

Methods	Training Time (h)	Params (M)	Run time (ms)
HOIM [16]	5.9	33.5	73
NAE [19]	5.2	33.5	71
ours ₃₄	2.1 + 1.7	42.9	35
ours	3.3 + 2.7	56.3	64
DMRNet [17]	-	49.1	-
ours ₃₄ (Retina)	3.0 + 1.7	50.9	42
ours (RetinaNet)	4.6 + 2.7	59.6	58
AlignPS [20]	19.0	42.2	51
AlignPS+ [20]	20.7	43.1	54
ours ₃₄ (FCOS)	2.9 + 1.7	50.2	41
ours(FCOS)	4.4 + 2.7	55.3	56

TABLE IX: Efficiency comparison between the proposed fully decoupled models and previous decoupled models. The result of [17] is estimated according to the paper as the source code is not released yet.

V. CONCLUSION AND LIMITATIONS

Decoupling end-to-end person search has been effectively explored for end-to-end person search to date. By analyzing their respective limitations, this paper takes a further step to enable fully decoupled end-to-end person search which achieves the optimum for both sub-tasks towards optimal person search performance. Notably, we design a task-incremental person search framework that decouples the end-to-end model architecture for the contradictory sub-tasks. The proposed task-incremental person search network further enables task-incremental training for the two sub-tasks, leading to fully decoupled end-to-end person search. Experimental results demonstrate the advantages of our proposed fully decoupled person search models. Moreover, this paper presents only baseline models for fully decoupled person search. We shall explore improving the parameter-efficiency and effectiveness of the proposed fully decoupled person search mechanism in future research.

REFERENCES

- [1] Xu, Yuanlu, et al. "Person search in a scene by jointly modeling people commonness and person uniqueness." Proceedings of the 22nd ACM international conference on Multimedia. 2014.
- [2] Ren, Shaoqing, et al. "Faster r-cnn: Towards real-time object detection with region proposal networks." Advances in neural information processing systems 28 (2015).
- [3] Tian, Zhi, et al. "Fcos: Fully convolutional one-stage object detection." Proceedings of the IEEE/CVF international conference on computer vision. 2019.
- [4] Lin, Tsung-Yi, et al. "Focal loss for dense object detection." Proceedings of the IEEE international conference on computer vision. 2017.
- [5] Carion, Nicolas, et al. "End-to-end object detection with transformers." European conference on computer vision. Cham: Springer International Publishing, 2020.
- [6] Zhu, Xizhou, et al. "Deformable detr: Deformable transformers for end-to-end object detection." arXiv preprint arXiv:2010.04159 (2020).
- [7] Zhang, Jeffrey O., et al. "Side-tuning: a baseline for network adaptation via additive side networks." Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part III 16. Springer International Publishing, 2020.
- [8] Sung, Yi-Lin, Jaemin Cho, and Mohit Bansal. "Lst: Ladder side-tuning for parameter and memory efficient transfer learning." Advances in Neural Information Processing Systems 35 (2022): 12991-13005.

- [9] He, Kaiming, et al. "Deep residual learning for image recognition." Proceedings of the IEEE conference on computer vision and pattern recognition. 2016.
- [10] Deng, Jia, et al. "Imagenet: A large-scale hierarchical image database." 2009 IEEE conference on computer vision and pattern recognition. Ieee, 2009.
- [11] He, Kaiming, et al. "Mask r-cnn." Proceedings of the IEEE international conference on computer vision. 2017.
- [12] Luo, Hao, et al. "Bag of tricks and a strong baseline for deep person re-identification." Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops. 2019.
- [13] Ioffe, Sergey, and Christian Szegedy. "Batch normalization: Accelerating deep network training by reducing internal covariate shift." International conference on machine learning. pmlr, 2015.
- [14] Wang, Cheng, et al. "Person search by a bi-directional task-consistent learning model." IEEE Transactions on Multimedia 25 (2022): 1190-1203.
- [15] Li, Feng, et al. "Dn-detr: Accelerate detr training by introducing query denoising." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2022.
- [16] Chen, Di, et al. "Hierarchical online instance matching for person search." Proceedings of the AAAI Conference on Artificial Intelligence. Vol. 34. No. 07. 2020.
- [17] Han, Chuchu, et al. "Decoupled and memory-reinforced networks: Towards effective feature learning for one-step person search." Proceedings of the AAAI Conference on Artificial Intelligence. Vol. 35. No. 2. 2021.
- [18] Xiao, Tong, et al. "Joint detection and identification feature learning for person search." Proceedings of the IEEE conference on computer vision and pattern recognition. 2017.
- [19] Chen, Di, et al. "Norm-aware embedding for efficient person search." Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2020.
- [20] Yan, Yichao, et al. "Anchor-free person search." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2021.
- [21] Lee, Sanghoon, et al. "Oimnet++: Prototypical normalization and localization-aware learning for person search." European Conference on Computer Vision. Cham: Springer Nature Switzerland, 2022.
- [22] Zheng, Liang, et al. "Person re-identification in the wild." Proceedings of the IEEE conference on computer vision and pattern recognition. 2017.
- [23] Ye, Mang, et al. "Deep learning for person re-identification: A survey and outlook." IEEE transactions on pattern analysis and machine intelligence 44.6 (2021): 2872-2893.
- [24] Wang, Guanshuo, et al. "Learning discriminative features with multiple granularities for person re-identification." Proceedings of the 26th ACM international conference on Multimedia. 2018.
- [25] Zheng, Liang, et al. "Scalable person re-identification: A benchmark." Proceedings of the IEEE international conference on computer vision. 2015.
- [26] Ristani, Ergys, et al. "Performance measures and a data set for multi-target, multi-camera tracking." European conference on computer vision. Cham: Springer International Publishing, 2016.
- [27] Everingham, Mark, et al. "The pascal visual object classes (voc) challenge." International journal of computer vision 88 (2010): 303-338.
- [28] Lan, Xu, Xiatian Zhu, and Shaogang Gong. "Person search by multi-scale matching." Proceedings of the European conference on computer vision (ECCV). 2018.
- [29] Chen, Di, et al. "Person search via a mask-guided two-stream cnn model." Proceedings of the european conference on computer vision (ECCV). 2018.
- [30] Munjal, Bharti, et al. "Query-guided end-to-end person search." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2019.
- [31] Yan, Yichao, et al. "Learning context graph for person search." Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2019.
- [32] Han, Chuchu, et al. "Re-id driven localization refinement for person search." Proceedings of the IEEE/CVF International Conference on Computer Vision. 2019.
- [33] Dong, Wenkai, et al. "Instance guided proposal network for person search." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2020.
- [34] Wang, Cheng, et al. "Tcts: A task-consistent two-stage framework for person search." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2020.
- [35] Li, Zhengjia, and Duoqian Miao. "Sequential end-to-end network for efficient person search." Proceedings of the AAAI Conference on Artificial Intelligence. Vol. 35. No. 3. 2021.
- [36] Zhou, Kaiyang, et al. "Omni-scale feature learning for person re-identification." Proceedings of the IEEE/CVF international conference on computer vision. 2019.
- [37] Yu, Rui, et al. "Cascade transformers for end-to-end person search." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2022.
- [38] Cao, Jiale, et al. "Pstr: End-to-end one-step person search with transformers." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2022.
- [39] Van de Ven, Gido M., and Andreas S. Tolias. "Three scenarios for continual learning." arXiv preprint arXiv:1904.07734 (2019).
- [40] Rusu, Andrei A., et al. "Progressive neural networks." arXiv preprint arXiv:1606.04671 (2016).
- [41] Li, Xilai, et al. "Learn to grow: A continual structure learning framework for overcoming catastrophic forgetting." International Conference on Machine Learning. PMLR, 2019.
- [42] Wang, Yu-Xiong, Deva Ramanan, and Martial Hebert. "Growing a brain: Fine-tuning by increasing model capacity." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2017.
- [43] Yoon, Jaehong, et al. "Lifelong learning with dynamically expandable networks." arXiv preprint arXiv:1708.01547 (2017).
- [44] Golkar, Siavash, Michael Kagan, and Kyunghyun Cho. "Continual learning via neural pruning." arXiv preprint arXiv:1903.04476 (2019).
- [45] Hung, Ching-Yi, et al. "Compacting, picking and growing for unforgetting continual learning." Advances in Neural Information Processing Systems 32 (2019).
- [46] Mallya, Arun, and Svetlana Lazebnik. "Packnet: Adding multiple tasks to a single network by iterative pruning." Proceedings of the IEEE conference on Computer Vision and Pattern Recognition. 2018.
- [47] Serra, Joan, et al. "Overcoming catastrophic forgetting with hard attention to the task." International conference on machine learning. PMLR, 2018.
- [48] Abati, Davide, et al. "Conditional channel gated networks for task-aware continual learning." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2020.
- [49] Wallingford, Matthew, et al. "Task adaptive parameter sharing for multi-task learning." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2022.
- [50] Kim, Hanjae, et al. "Prototype-guided saliency feature learning for person search." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2021.
- [51] Dosovitskiy, Alexey, et al. "An image is worth 16x16 words: Transformers for image recognition at scale." arXiv preprint arXiv:2010.11929 (2020).
- [52] Dong, Wenkai, et al. "Bi-directional interaction network for person search." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2020.
- [53] Dai, Jifeng, et al. "Deformable convolutional networks." Proceedings of the IEEE international conference on computer vision. 2017.
- [54] Cai, Zhaowei, and Nuno Vasconcelos. "Cascade r-cnn: Delving into high quality object detection." Proceedings of the IEEE conference on computer vision and pattern recognition. 2018.
- [55] Lin, Tsung-Yi, et al. "Feature pyramid networks for object detection." Proceedings of the IEEE conference on computer vision and pattern recognition. 2017.