# Sequential Transformer for End-to-End Person Search

Long Chen, Jinhua Xu East China Normal University, Shanghai, China

longchen@stu.ecnu.edu.cn, jhxu@cs.ecnu.edu.cn

# Abstract

Person Search aims to simultaneously localize and recognize a target person from realistic and uncropped gallery images. One major challenge of person search comes from the contradictory goals of the two sub-tasks, i.e., person detection focuses on finding the commonness of all persons so as to distinguish persons from the background, while person re-identification (re-ID) focuses on the differences among different persons. In this paper, we propose a novel Sequential Transformer (SeqTR) for end-to-end person search to deal with this challenge. Our SeqTR contains a detection transformer and a novel re-ID transformer that sequentially addresses detection and re-ID tasks. The re-ID transformer comprises the self-attention layer that utilizes contextual information and the cross-attention layer that learns local fine-grained discriminative features of the human body. Moreover, the re-ID transformer is shared and supervised by multi-scale features to improve the robustness of learned person representations. Extensive experiments on two widely-used person search benchmarks, CUHK-SYSU and PRW, show that our proposed SeqTR not only outperforms all existing person search methods with a 59.3% mAP on PRW but also achieves comparable performance to the state-of-the-art results with an mAP of 94.8% on CUHK-SYSU.

#### **1. Introduction**

Practical applications of person search, such as searching for suspects and missing people in intelligent surveillance, require separating people from complex background and discriminating target identities (IDs) from other IDs. It involves two fundamental tasks in computer vision, *i.e.*, pedestrian detection and person re-identification (re-ID). Pedestrian detection aims at detecting the bounding boxes (Bboxes) of all candidates in the image. Person re-ID aims at retrieving a person of interest across multiple nonoverlapping cameras. Person search has recently attracted tremendous interest of researchers in the computer vision community for its importance in building smart cities. How-



Figure 1. Comparison of person search frameworks. (a) The twostep framework. (b) The one-step framework. (c) Our proposed SeqTR adopts the sequential framework to perform detection and re-ID in order.

ever, it remains a difficult task that suffers from many challenges, such as jointly optimizing contradictory objectives of two sub-tasks in a unified framework, scale/pose variations, background clutter and occlusions and so on.

According to training manners, existing person search methods can be generally grouped into two categories: twostep frameworks and one-step frameworks. Two-step methods typically perform detection and re-ID with two separate independent models. As shown in Fig. 1(a), pedestrians are first detected by an off-the-shelf detection model. After non-maximum suppression (NMS), the person patches are cropped and resized (C&R) into a fixed size. Then the person re-ID model is applied to produce ID feature embeddings, which will be used to calculate the similarity between the query persons and the candidates. The two-step frameworks can achieve satisfactory performance since each step focuses on one task and no contradictory is involved. However, this pipeline is time-consuming and resource-consuming. In contrast, one-step methods simultaneously optimize two sub-tasks in a joint framework (Fig. 1(b)). The two sub-tasks first share a common backbone for features extraction and then detection head and re-ID head are applied in parallel.

In terms of architecture, the sequential framework combines the merits of two-step and one-step frameworks. It not only inherits the better performance of two-stage frameworks via providing accurate bounding boxes (Bboxes) for the re-ID stage but also preserves the efficiency of the endto-end training manner of one-step frameworks. However, as Li *et al.* [13] has pointed out, the performance bottleneck of this architecture lies in the design of the re-ID subnetwork. In addition, we find that NMS, commonly used in the detection models, primarily hinders the inference speed of this architecture, especially in crowded scenes.

As transformers [17] become popular in vision tasks, transformers-based person search frameworks [1, 24] also show advantages over CNN-based models, such as no NMS needed and powerful capability of learning fine-grained features.

Motivated by the above observations, we propose a novel Sequential transformer (SeqTR) for end-to-end Person Search (Fig. 1(c)). It is a sequential framework, in which two transformers are integrated seamlessly to address the detection and re-ID tasks. Meanwhile, the two transformers are decoupled with different features for the two contradictory tasks.

In summary, we make the following contributions:

- We propose a novel Sequential Transformer (SeqTR) model for end-to-end person search, which utilizes two transformers to sequentially perform pedestrian detection and re-ID without NMS post-processing.
- We propose a novel re-ID transformer to generate discriminative re-ID feature embeddings. To make full use of context information, we introduce the selfattention mechanism in our re-ID transformer. Meanwhile, we employ multiple cross-attention layers to learn local fine-grained features. To obtain scaleinvariant person representations, our re-ID transformer is shared by multi-scale features.
- We achieve a state-of-the-art result on two datasets. Comprehensive experiments show the merits of our proposed modules. Furthermore, with PVTv2-B2 [19] backbone, SeqTR achieves 59.3% mAP that outperforms all existing person search models on PRW [25].

#### 2. Related Work

#### 2.1. CNN-based Person Search

Person search has attracted a lot of attention from the computer vision community. A large number of meth-

ods have been proposed and achieved remarkable results. According to the training manner, existing person search frameworks can be divided into two-step and one-step methods. Two-step person search models first perform pedestrian detection and subsequently crop the detected people for re-ID. Zheng et al. [25] first exhaustively evaluate the combinations of different detectors and re-ID models. Chen et al. [4] propose a mask-guided two-stream network to obtain enhanced feature representation. Lan et al. [12] analyze the multi-scale misalignment caused by the detector and exploit knowledge distillation to address it. Wang *et al.* [18] utilize an identity-guided query detector to extract the query-like proposals and employ a detection-adapted model for re-ID. One-step person search models integrate detection and re-ID into a joint framework, which enables end-to-end training of two sub-tasks. Xiao et al. [21] propose the first one-step person search model by introducing a re-ID branch and Online Instance Matching (OIM) loss in the Faster R-CNN detector. Liu et al. [14] and Chang *et al.* [3] discard the proposal generation operation and search the query person directly on the uncropped images by sequential decision making or reinforcement learning. Xiao et al. [20] use Center Loss to enhance feature discrimination. Yan et al. [23] enrich the features with surrounding persons. Munjal et al. [15] build the relationship between the query image and gallery image by integrating a query-guided Siamese squeeze-and-excitation block into the backbone. Han et al. [8] develop an RoI transform layer that enables gradient flow from the re-identifier to the detector for localization refinement. Chen et al. [5] propose a norm-aware embedding (NAE) to improve re-ID performance. Dong et al. [6] employ a Siamese network that takes both the entire image and cropped persons to better guide the feature learning of the person. Yan et al. [22] introduce the first anchor-free approach for person search. Li et al. [13] propose a Sequential End-to-end Network (SeqNet) to obtain accurate Bboxes for the re-ID stage, in which detection and re-ID are considered as a progressive process and tackled with two sub-networks sequentially. SeqNet inherits the sequential process of two-stage methods and the end-to-end training fashion and efficiency of the one-step methods. Our work is inspired by SeqNet, and we use the sequential framework and replace the CNN sub-networks for detection and re-ID with two transformers. Employing the structure advantage of the transformer, no NMS is needed during training and inference, and the two subnetworks are integrated with deformable attention seamlessly rather than the ROI-align in SeqNet.

### 2.2. Transformer-based Person Search

Recently, transformers-based person search frameworks [1, 24] have been also proposed. The COAT model [24] is a cascaded one-step method, in which an occluded atten-



Figure 2. Architecture of our proposed SeqTR, which comprises a backbone, a detection transformer and a re-ID transformer.

tion transformer is used for feature enhancement before the parallel detection head and re-ID head. In PSTR [1], a detection decoder and a re-ID decoder are designed for the two tasks. The output features of the detection decoder are fed into the re-ID decoder, therefore the two decoders with contradictory goals are coupled. Considering the advantages of the transformer, we aim to utilize the transformer to design a robust re-ID sub-network to alleviate the performance bottleneck of the sequential framework.

### 3. Method

In this section, we introduce our proposed SeqTR in detail. Firstly, we give an overview architecture of SeqTR in Sec. 3.1. Secondly, the details of our designed re-ID transformer are elaborated in Sec. 3.2. Finally, we introduce the training and inference process in Sec. 3.3.

### 3.1. SeqTR Architecture

The overall architecture of our SeqTR is depicted in Fig. 2. It contains three main components: a backbone to extract multi-scale feature maps of the input image, a detection transformer to predict Bboxes, and a novel re-ID transformer to learn robust person feature embeddings.

**Backbone.** Starting from the initial image  $x_{img} \in \mathbb{R}^{3 \times H_0 \times W_0}$  (with 3 color channels). The backbone extracts original multi-scale feature maps  $\{x^l\}_{l=1}^3$  from stages  $P_2$  through  $P_4$  in PVTv2-B2 [19] (or from stages  $C_3$  through  $C_5$  in RestNet [10]). The resolution of  $x^l$  is  $2^{l+2}$  lower than the input image.

**Detection Transformer.** We introduce the transformerbased detector, deformable DETR [27], into our framework to predict the pedestrian bounding boxes. However, The difference with the original deformable DETR is the input features. First, the channel dimensions of all feature maps  $\{x^l\}_{l=1}^3$  from the backbone are mapped to a smaller dimension d = 256 by  $1 \times 1$  convolution. Then, a  $3 \times 3$  deformable convolution is used to generate more accurate feature maps. Finally,  $\{F_{bi} \in \mathbb{R}^{d \times H \times W}\}_{i=2}^4$  are transfomed from original feature maps  $\{x^l\}_{l=1}^3$  by the above two steps and fed into a standard deformable DETR.

re-ID Transformer. Our re-ID transformer aims to



Figure 3. Architecture of our proposed re-ID transformer.

adaptively learn discriminative re-ID features around the human body center. Motivated by object queries in DETR [2], we set a fixed number of learnable re-ID queries  $Q_r$  to reconcile the relationship between detection and re-ID and obtain re-ID feature embeddings.

#### 3.2. re-ID transformer

The architecture of the re-ID transformer is shown in Fig. 3. Each re-ID transformer layer is composed of a self-attention layer and K cross-attention layers. The self-attention layer comprises a multi-head attention module and a layer normalization. The cross-attention layer contains a deformable attention module and a layer normalization. Suppose that the detection transformer decodes N objects in each image. The re-ID query number is also set as N. Taking the enhanced backbone features  $F_{bi}$ ,  $i \in [2, 4]$ , N

reference points  $P_q$  from the detection transformer and N re-ID queries  $Q_r$  as input, the re-ID transformer outputs N instance-level re-ID embeddings  $F_{ri}$  that have the same dimension as the pixel features. These instance-level re-ID feature embeddings are highly associated with pedestrian locations. Furthermore, to aggregate multi-scale features, multi-scale feature maps  $\{F_{bi}\}_{i=2}^4$  are used to generate multi-scale re-ID embeddings  $\{F_{ri} \in \mathbb{R}^{d \times H \times W}\}_{i=2}^4$  by the re-ID transformer. During inference, all multi-scale re-ID embeddings  $\{F_{ri}\}_{i=2}^4$  are concatenated to perform matching.

**Re-ID Queries.** To mitigate the objective contradictory problem, we set re-ID queries  $Q_r$ , like object queries, to obtain re-ID features. Specifically, re-ID queries guarantees that the final re-ID embeddings  $\{F_{ri}\}_{i=2}^4$  are instance-level fine-grained features learned from the augmented multiscale backbone features  $\{F_{bi}\}_{i=2}^4$ . Through this design, the final learned re-ID feature embeddings are highly correlated with the detected pedestrian locations, but not affected by the detection features. This is different from the re-ID decoder in PSTR [1], in which the re-ID queries come from the output features of the detection decoder.

Self-Attention Layer. To produce discriminative re-ID feature embeddings, we introduce the self-attention layer into the re-ID transformer to learn contextual information. This is different from the re-ID decoder in PSTR [1], in which no self-attention layer is used. From the ablation study (Table 3) in experiments, the performance is improved with the self-attention layer. Specifically, we adopt a standard multi-head self-attention (with H heads) in the Transformer [17]. We denote the input of the self-attention layer as  $Y_q$ . The initial input  $Y_q = Q_r$ .  $Y_q$  are transformed into query vectors  $Q \in \mathbb{R}^{N \times d_k}$ , key vectors  $K \in \mathbb{R}^{N \times d_k}$  and value vectors  $V \in \mathbb{R}^{N \times d_v}$  by three different linear projections. The output embeddings then are generated by performing the multi-head self-attention module.

$$nead_i = Attention(QW_i^Q, KW_i^K, VW_i^V), \quad (1)$$

where  $W_i^Q \in \mathbb{R}^{d \times d_k}$ ,  $W_i^K \in \mathbb{R}^{d \times d_k}$ ,  $W_i^V \in \mathbb{R}^{d \times d_v}$ ,  $d_k = d_v = d/H$ . The self-attention module use Scaled Dot-Product Attention in each head:

Attention
$$(Q, K, V) = \operatorname{softmax}\left(\frac{QK^{\top}}{\sqrt{d_k}}\right) V.$$
 (2)

The embeddings from all heads are concatenated and projected to yield *d*-demensional embeddings:

$$MultiHead(Q, K, V) = Concat(head_1, ..., head_H)W^O,$$
(3)

where  $W^O \in \mathbb{R}^{Hd_k \times d}$ . At last, we use a layer normalization to get the final embeddings  $\hat{Y}_q$ .

$$\hat{Y}_q = \text{layernorm}(Y_q + \text{dropout}(\text{MultiHead}(Q, K, V)).$$
(4)

The self-attention layer in the first re-ID transformer layer can be skipped. After passing through the first re-ID transformer layer, the output features are correlated with reference points. N feature embeddings correspond to N locations respectively. These embeddings interact with each other for learning spatial relationship by the self-attention layer in the  $m^{th}$  ( $m \in [2, M]$ ) re-ID transformer layer, resulting to enhance feature embeddings by instances in the same scene.

**Cross-Attention Layer.** Different from the previous works that use the RoI-Align layer on detection features, we employ and stack several cross-attention layers to address the region misalignment. In the cross-attention layer, there is a deformable attention module and a layer normalization. The deformable attention module proposed by deformable DETR [27], only attends to a small set of key sampling points around a reference point. It is useful for learning fine-grained features. Given an input feature map  $F_{bi} \in \mathbb{R}^{C \times H \times W}$ , a set of detected bounding boxes, *i.e.*, reference points (denoted  $P_q$ ), and query features (denoted  $Z_q$ ), the output feature embeddings  $\hat{Z}_q$  can be calculated:

$$\hat{Z}_q = layernorm(Z_q + dropout(DeformAttn(Z_q, P_q, F_{bi})),$$
(5)

$$DeformAttn(Z_q, P_q, F_{bi}) = \sum_{h=1}^{H} W_h \left[ \sum_{s=1}^{S} A_{hs} \cdot W'_h F_{bi}(P_q + \Delta P_{hs}) \right]$$
(6)

where H is the total attention heads, S is the total sampled key number.  $A_{hs}$  and  $\Delta P_{hs}$  denote attention weight of the  $s^{th}$  sampling point in the  $h^{th}$  attention head and the sampling offset, respectively. Both are obtained via linear projection over the query feature  $Z_q$ , respectively. In this way, each query feature corresponds to one detected bounding boxes and integrates the features of the surround sampling points. In PSTR [1], features at sampling points are averaged rather than using the attention weight  $A_{hs}$  as in Eq. 6 because it was observed that the attention weights from the query struggle to effectively capture the features of a person instance. We think it may be caused by the coupling of the two decoders since the re-ID queries in PSTR [1] are from the detection decoder.

Schemes of Employing Multi-scale Features. Much previous work has demonstrated that employing multi-scale feature maps is useful for addressing scale variation in person search. To obtain scale-invariant re-ID features, we propose several schemes of employing multi-scale features. First, A straightforward way is to concatenate the augmented backbone features  $\{F_{bi}\}_{i=2}^{4}$  and feed to the re-ID transformer to produce  $F_{rm} \in \mathbb{R}^{N \times d}$ , as shown in Fig.



Figure 4. Comparison of different re-ID transformer schemes. (a) Multi-scale re-ID transformer.  $\{F_{bi}\}_{i=2}^{4}$  are concatenated as input features. (b) Parallel re-ID transformer. Each independent re-ID transformer is responsible for a single-scale input feature. (c) Shared re-ID transformer.  $\{F_{bi}\}_{i=2}^{4}$  respectively go through a common shared re-ID transformer.

4(a). To align the dimension of the final matching embeddings with other schemes in Fig. 4, we also design "Multiscale re-ID transformer-3d", whose deformable attention modules are replaced by multi-scale deformable attention modules [27]. Correspondingly, the re-ID queries are adjusted to  $Q_r \in \mathbb{R}^{N \times 3d}$ , resulting in  $F_{rm} \in \mathbb{R}^{N \times 3d}$ . We also build three independent re-ID transformers for threelevel features  $\{F_{bi}\}_{i=2}^{4}$  to obtain three-level re-ID feature embeddings  $\{F_{ri} \in \mathbb{R}^{N \times d}\}_{i=2}^{4}$ , respectively. we call it parallel re-ID transformer (Fig. 4(b)). As opposed to parallel re-ID transformer, the shared re-ID transformer (Fig. 4(c)) means that three-scale input features are respectively fed to a common re-ID transformer to generate re-ID feature embeddings. The following ablation studies (Table 4) verify that the shared re-ID transformer achieves the best performance.

#### 3.3. Training and Inference

For each image, our SeqTR predicts N classification scores, bounding boxes and re-ID feature embeddings

 $\{F_{ri}\}_{i=2}^4$ . In the training phase,  $\{F_{ri}\}_{i=2}^4$  are supervised separately. They are concatenated during inference.

During training, our SeqTR is trained end-to-end for detection and re-ID. Specifically, detection transformer is supervised with loss functions of deformable DETR [27] for classification  $(L_{cls})$ , bounding-box IoU loss  $(L_{iou})$ , bounding-box Smooth-L1 loss  $(L_{cls})$ . While the re-ID transformer is supervised by the Focal OIM loss  $(L_{oim})$  [22].

The overall loss is given by:

$$L = \lambda_1 L_{cls} + \lambda_2 L_{iou} + \lambda_3 L_{l1} + \lambda_4 L_{oim} \tag{7}$$

where  $\lambda_1, \lambda_2, \lambda_3, \lambda_4$ , responsible for the relative loss importance, are set as 2.0, 5.0, 2.0, 0.5, respectively.

During inference, our SeqTR predicts Bboxes and corresponding re-ID feature embeddings for gallery images. For the query person, we get predictions of the query image in the same way and then choose the one that has maximum overlap with its annotated bounding box.

#### 4. Experiments

In this section, we conduct experiments on two widely utilized person search datasets. We first introduce two large datasets and evaluation metrics. Then we describe some implementation details. Afterwards, we compare the overall performance of our methods with state-of-the-art methods. Finally, we perform ablation studies to validate the effectiveness of our methods on the PRW [25] dataset.

#### 4.1. Datasets and Settings

**CUHK-SYSU.** Scene images in the CUHK-SYSU [21] are collected from real street snaps and movies. There are a total of 18,184 realistic and uncropped images, 96,143 annotated bounding boxes and 8,432 different identities. The dataset is partitioned into two parts without overlap. The training set includes 11,206 images, 55,272 pedestrians, and 5,532 identities. The test set contains 6,978 images, 40,871 pedestrians, and 2,900 identities. During inference, for each query, the dataset defines a gallery set with different sizes from 50 to 4,000 to evaluate the performance scalability of models. Following the previous works, we report the results with the gallery size of 100 if not specified.

**PRW.** Images in the PRW [25] dataset are collected by 6 static cameras at Tsinghua university. There are 11,816 video frames and 43,110 annotated bounding boxes. 34,304 of these boxes are annotated with 932 labelled identities and the rest are marked as unknown identities. It is also divided into two groups. The training set contains 5,704 images, 18,048 pedestrians, and 482 identities. The test set has 6,112 images and 2,057 query persons with 450 identities. During inference, for each query person, the gallery set is the whole test set, *i.e.*, the gallery size is 6,112.

**Evaluation Metrics.** Following the previous works [21], we employ Mean Average Precision (mAP) and Cumulative Matching Characteristics (CMC top-K) to evaluate the performance of the person search.

#### **4.2. Implementation Details**

We adopt ResNet50 [10] and transformer-based PVTv2-B2 [19] that are pre-trained on ImageNet [16] as backbone. To train our model, we adopt the AdamW optimizer with a weight decay rate of 0.0001. The initial learning rate is set to 0.0001 that is warmed up during the first epoch and decreased by a factor of 10 at 19th and 23th epoch, with a total of 24 epochs. For CUHK-SYSU/PRW, the circular queue size of OIM is set to 5000/500. During training, we employ a multi-scale training strategy, where the longer side of the image is randomly resized from 400 to 1666. For inference, we rescale the test images to a fixed size of  $1500 \times 900$  pixels. For our SeqTR with ResNet50 [10] backbone, we use one NVIDIA GeForce RTX 3090 to run all experiments and batch size set to 2. Our SeqTR with PVTv2-B2 [19] backbone is trained on two RTX 3090 GPUs with batch size set to 1 because of the limitation of GPU memory.

#### **4.3.** Comparison to the State-of-the-arts

We compare our SeqTR with the state-of-the-arts, including both two-step models [4, 7, 8, 12, 18] and one-step models [1, 3, 5, 6, 9, 11, 13-15, 20-24, 26], on two datasets.

**Results on CUHK-SYSU.** As shown in Table 1, our SeqTR outperforms most one-step methods and achieves comparable performance to two-step methods on the CUHK-SYSU test set [21].

The best two-step method TCTS [18] achieves mAP scores of 93.9%. Among one-step methods with the ResNet50 [10], COAT [24] achieves the best mAP score of 94.2%. Our SeqTR with the same ResNet50 backbone, which achieves comparable 93.4% mAP and 94.1% top-1 accuracy, outperforms AlignPS [22] by 0.3% and 0.7% in mAP and top-1 accuracy, respectively. Our results are slightly worse than the transformer-based COAT [24] and PSTR [1].

Then, based on PVTv2-B2 [19] backbone, the performance of our SeqTR is significantly improved to 94.8% mAP and 95.5% top-1 accuracy. For a fair comparison, we reproduce the performance of PSTR [1] with the same PVTv2-B2 [19] backbone (named PSTR\*) to eliminate the effects of different training strategies, *i.e.*, single-GPU training and distributed training. Specifically, we set batch size from 2 to 1 and use two RTX 3090 GPUs for distributed training three times. The average of the three reproduced results is then calculated and reported in Table 1. Our method outperforms the reproduced results of PSTR by 0.2% in mAP. Moreover, the post-processing strategy Context Bipartite Graph Matching(CBGM) [13] is widely used to im-

Mathad	Paalthono	CUHK-SYSU		PRW				
wiethou	Dackbolle	mAP(%)	Top-1(%)	mAP(%)	Top-1(%)			
Two-step methods								
MGTS [4]	4] VGG16 83.0 83.7		32.6	72.1				
CLSA [12]	ResNet50	87.2	88.5	38.7	65.0			
RDLR [8]	ResNet50	93.0	94.2	42.9	70.2			
IGPN [7]	ResNet50	90.3	91.4	47.2	87.0			
TCTS [18]	ResNet50	93.9	95.1	46.8	87.5			
One-step methods with	CNNs							
OIM [21]	ResNet50	75.5	78.7	21.3	49.4			
NPSM [14]	ResNet50	77.9	81.2	24.2	53.1			
RCAA [3]	ResNet50	79.3	81.3	-	-			
IAN [20]	ResNet50	76.3	80.1	23.0	61.9			
CTXGraph [23]	ResNet50	84.1	86.5	33.4	73.6			
QEEPS [15]	ResNet50	88.9	89.1	37.1	76.7			
BI-Net [6]	ResNet50	90.0	90.7	45.3	81.7			
APNet [26]	ResNet50	88.9	89.3	41.9	81.4			
NAE [5]	ResNet50	91.5	92.4	43.3	80.9			
NAE+ [5]	ResNet50	92.1	92.9	44.0	81.1			
PGSFL [11]	ResNet50	90.2	91.8	42.5	83.5			
SeqNet [13]	ResNet50	93.8	94.6	46.7	83.4			
DMRN [9]	ResNet50	93.2	94.2	46.9	83.3			
AlignPS [22]	ResNet50	93.1	93.4	45.9	81.9			
One-step methods with	transformers							
COAT [24]	ResNet50	94.2	94.7	53.3	87.4			
PSTR [1]	ResNet50	93.5	95.0	49.5	87.8			
SeqTR(Ours)	ResNet50	93.4	94.1	52.0	86.5			
PSTR [1]	PVTv2-B2	95.2	96.2	56.5	89.7			
PSTR* [1]	PVTv2-B2	94.6	95.6	57.6	90.1			
SeqTR(Ours)	PVTv2-B2	94.8	95.5	59.3	89.4			
COAT [24]+CBGM	ResNet50	94.8	95.2	54.0	89.1			
PSTR [1]+CBGM	PVTv2-B2	95.8	96.8	58.1	92.0			
PSTR* [1]+CBGM	PVTv2-B2	95.2	96.1	58.2	91.5			
SeqTR(Ours)+CBGM	PVTv2-B2	95.4	96.3	59.8	90.6			

Table 1. Comparison with the state-of-the-art methods on CUHK-SYSU and PRW test sets. \* denotes our reproduced result. The highest scores in each group are highlighted in bold.

prove mAP and top-1 accuracy. By employing CBGM, our SeqTR achieves 95.4% mAP and 96.3% top-1 accuracy, which outperforms the reproduced results of PSTR\* with CBGM.

We also evaluate the performance scalability of these models with different gallery sizes. Fig. 5 shows that the mAP of all methods decreases monotonically as the gallery size increases, which illustrates the fact that more distracting persons introduced in the larger gallery make searching much more difficult. As shown in Fig. 5, our SeqTR outperforms most models.

**Results on PRW.** The PRW dataset [25] is more challenging than the CUHK-SYSU dataset [21] for less training data and larger gallery size. Furthermore, there is a large number of people wearing similar uniforms and there are more scale variations, pose/viewpoint changes and occlusions. Nevertheless, as can be observed from Table 1, our method achieves strong performance.

With ResNet50 [10] backbone, our SeqTR achieves 52.0% mAP and 86.5% top-1 accuracy, outperforming all two-step methods and with a significant gain of 2.5% mAP than PSTR [1] with the same backbone. The performance of our method is slightly lower than COAT [24] by 1.3% mAP and 0.9% top-1 accuracy.

With PVTv2-B2 backbone [19], our SeqTR achieves 59.3% mAP and 89.4% top-1 accuracy, outperforming all



Figure 5. Comparison with (a) two-step models and (b) one-step models on CUHK-SYSU with different gallery sizes.

Transformer layers M	Cross-attention layers K	mAP(%)	<b>Top-1</b> (%)
	2	50.8	86.5
2	3	50.4	86.0
	4	49.9	85.5
	2	50.4	86.7
3	3	52.0	86.5
	4	50.7	86.8
	2	50.5	86.7
4	3	50.3	86.1
	4	50.6	86.7

Table 2. Ablation study for different shared re-ID transformer structures on PRW dataset.

existing methods with a clear margin on mAP. We attribute it to our designed re-ID transformer which alleviates some challenges, such as recognizing the query person from cotravellers wearing the same uniform. Finally, our SeqTR is improved to the best 59.8% mAP and comparable 90.6% top-1 accuracy with CBGM.

### 4.4. Ablation Study

We perform a series of ablation studies on the PRW [25] dataset to analyze our design decisions. Limited by the memory of the RTX 3090 GPU, we choose the ResNet50 [10] backbone to void the impact of the distributed training.

**Re-ID Transformer Structure.** Setting the number of the self-attention layer in each transformer layer to 1, we evaluate the impact of the number of transformer layers M and the number of cross-attention layers K. As shown in Table 2, when the number of transformer layers is greater than 2, different combinations of M and K have a slight impact on performance. Among these configurations, when M = 3 and K = 3, our SeqTR achieves the best performance of 52.0% mAP and 86.5% top-1 accuracy.

Method	mAP(%)	<b>Top-1</b> (%)
re-ID transformer	52.0	86.5
re-ID transformer w/o self-attention layers	49.6	85.5

Table 3. Comparative results of adding and removing selfattention layer on PRW dataset.

Re-ID transformer scheme	mAP(%)	<b>Top-1</b> (%)
Multi-scale re-ID transformer-d	44.1	80.9
Multi-scale re-ID transformer-3d	44.1	83.1
Parallel re-ID transformer	51.1	85.2
Shared re-ID transformer	52.0	86.5

Table 4. Comparative results with different variants of the re-ID transformer on PRW dataset.

Input feature	Е	$F_{b4}$	$F_{b3}$	$F_{b2}$	mAP(%)	<b>Top-1</b> (%)
Single-scale feature	$\checkmark$				26.5	66.4
		$\checkmark$			41.6	79.5
			$\checkmark$		45.6	82.9
				$\checkmark$	41.7	82.6
Multi-scale feature		$\checkmark$			41.6	79.5
		$\checkmark$	$\checkmark$		47.8	82.7
		$\checkmark$	$\checkmark$	$\checkmark$	52.0	86.5

Table 5. Comparative results by employing different input features on PRW dataset. " $\checkmark$ " means using the corresponding feature. "E" denotes the output feature of the encoder in the detection transformer.

**Importance of Self-Attention Layer.** We also evaluate the importance of the self-attention layer. In Table 3, we find that adding self-attention layers yields improvements of 2.4% on mAP and 1% on top-1 accuracy respectively.

Schemes of Employing Multi-scale Features. To evaluate the effect of different re-ID transformer schemes, we design three different variants as illustrated in Fig. 4 and report the results in Table 4. First, For a multi-scale re-ID transformer-*d*, it outputs *d* dimensional re-ID feature embeddings for matching. We obtain 44.1% on mAP and 80.9% on top-1 accuracy. To align with the 3*d* dimensional matching embeddings of other schemes (Fig. 4(b) and Fig. 4(c)), we also design a multi-scale re-ID transformer-3*d*. However, it has no improvement on mAP. Compared to the multi-scale re-ID transformer-3*d*, the parallel re-ID transformer (Fig. 4(b)) has absolute gains of 7.0% on mAP and 2.1% on top-1 accuracy. Then, the shared re-ID transformer (Fig. 4(c)) achieves the best performance with 52.0% on mAP and 86.5% on top-1 accuracy.

**Choices of Input Features to re-ID transformer.** We conduct experiments on employing different input features to the shared re-ID transformer, including single-level and multi-scale features. The results are reported in Table 5. Specifically, we first evaluate the single-level feature respectively. Among these single-level features, the output feature of the encoder in the detection transformer provides



Figure 6. Qualitative comparison with PSTR [1]. The yellow bounding boxes denote the queries, while the green and red bounding boxes denote correct and incorrect top-1 matches, respectively. Row (a) are two cases to illustrate the strength of the sequential framework. Row (b) are two cases to show the importance of self-attention layers in our re-ID transformer. Row (c) are two cases to show the advantages of cross-attention layers in our re-ID transformer.

Method	Backbone	GPU	Time(ms)	mAP(%)	<b>Top-1</b> (%)
NAE [5]	ResNet50	RTX 3090	80	43.3	80.9
SeqNet [13]	ResNet50	RTX 3090	106	46.7	83.4
AlignPS [22]	ResNet50	RTX 3090	44	45.9	81.9
COAT [24]	ResNet50	RTX 3090	130	53.3	87.4
PSTR [1]	ResNet50	RTX 3090	52	49.5	87.8
SeqTR(Ours)	ResNet50	RTX 3090	86	52.0	86.5
PSTR [1]	PVTv2-B2	RTX 3090	88	56.5	89.7
SeqTR(Ours)	PVTv2-B2	RTX 3090	130	59.3	89.4

Table 6. Comparative results of person search efficiency on the PRW dataset.

less information for the re-ID task and is discarded in later experiments. Relatively, C4 yields the best performance. Furthermore, we also show the performance of utilizing multi-scale features. As can be observed, the best performance is achieved by using three-level features.

Efficiency Comparison. Generally, there are more pedestrians in every scene image in the PRW [25] dataset. To evaluate our contributions in the sequential framework, we conduct runtime efficiency analysis on PRW [25] dataset. As shown in Table 6, our SeqTR with ResNet50 [10] backbone takes 86 milliseconds to process an image, which is faster than SeqNet [13] and COAT [24]. It is attributed to the design without requiring an NMS. For using PVTv2-B2 backbone [19], our SeqTR is slower than PSTR [1], but achieves an absolute of 2.8% mAP over PSTR [1]. Our SeqTR with PVTv2-B2 backbone has the same speed of 130 milliseconds with COAT [24] with ResNet50 backbone, however our method achieves +6.0% and +2.0% gains of mAP and top-1 accuracy respectively.

Qualitative Results. To demonstrate the performance of our SeqTR, we show some qualitative comparisons between our SeqTR with PSTR [1] on PRW [25] dataset. As shown in Fig. 6(a), our SeqTR achieves more accurate pedestrian localizations in both examples, because the sequential framework produces high-quality detection results first that then benefit for the re-ID stage. In both cases of Fig. 6(b), compared to PSTR [1], our SeqTR accurately identifies the query persons, whose co-travellers wear similar uniforms. It is attributed to the self-attention layer that employs contextual information. In addition, the cross-attention layers in our re-ID transformer contribute to focusing on meaningful regions, although occlusions occur in the given query person in Fig. 6(c). The above examples also illustrate that our SeqTR further alleviates some challenges, such as occlusions and distinguishing similar appearances.

# 5. Conclusion

In this paper, we propose a novel Sequential Transformer (SeqTR) for end-to-end person search. Within our SeqTR, a detection transformer and a re-ID transformer are integrated to solve the two contradictory tasks sequentially. We design a re-ID transformer that contains self-attention layers and cross-attention layers to generate discriminative re-ID feature embeddings. Furthermore, our re-ID transformer adopts a share strategy for employing multi-scale features. Extensive experiments demonstrate the performance of our proposed framework, which achieves state-of-the-art results on PRW [25] dataset.

## References

- Jiale Cao, Yanwei Pang, Rao Muhammad Anwer, Hisham Cholakkal, Jin Xie, Mubarak Shah, and Fahad Shahbaz Khan. Pstr: End-to-end one-step person search with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9458– 9467, 2022. 2, 3, 4, 6, 8
- [2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-toend object detection with transformers. In *European Conference on Computer Vision*, pages 213–229. Springer, 2020.
   3
- [3] Xiaojun Chang, Po-Yao Huang, Yi-Dong Shen, Xiaodan Liang, Yi Yang, and Alexander G Hauptmann. Rcaa: Relational context-aware agents for person search. In *Proceedings of the European Conference on Computer Vision* (ECCV), pages 84–100, 2018. 2, 6
- [4] Di Chen, Shanshan Zhang, Wanli Ouyang, Jian Yang, and Ying Tai. Person search via a mask-guided two-stream cnn model. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 734–750, 2018. 2, 6
- [5] Di Chen, Shanshan Zhang, Jian Yang, and Bernt Schiele. Norm-aware embedding for efficient person search. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 12615–12624, 2020. 2, 6, 8
- [6] Wenkai Dong, Zhaoxiang Zhang, Chunfeng Song, and Tieniu Tan. Bi-directional interaction network for person search. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2839–2848, 2020. 2, 6
- [7] Wenkai Dong, Zhaoxiang Zhang, Chunfeng Song, and Tieniu Tan. Instance guided proposal network for person search. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2585–2594, 2020.
- [8] Chuchu Han, Jiacheng Ye, Yunshan Zhong, Xin Tan, Chi Zhang, Changxin Gao, and Nong Sang. Re-id driven localization refinement for person search. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9814–9823, 2019. 2, 6
- [9] Chuchu Han, Zhedong Zheng, Changxin Gao, Nong Sang, and Yi Yang. Decoupled and memory-reinforced networks: Towards effective feature learning for one-step person search. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 1505–1512, 2021. 6
- [10] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceed-ings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 3, 6, 7, 8
- [11] Hanjae Kim, Sunghun Joung, Ig-Jae Kim, and Kwanghoon Sohn. Prototype-guided saliency feature learning for person search. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pages 4865– 4874, 2021. 6
- [12] Xu Lan, Xiatian Zhu, and Shaogang Gong. Person search by multi-scale matching. In *Proceedings of the European conference on computer vision (ECCV)*, pages 536–552, 2018.
   2, 6

- [13] Zhengjia Li and Duoqian Miao. Sequential end-to-end network for efficient person search. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 2011–2019, 2021. 2, 6, 8
- [14] Hao Liu, Jiashi Feng, Zequn Jie, Karlekar Jayashree, Bo Zhao, Meibin Qi, Jianguo Jiang, and Shuicheng Yan. Neural person search machines. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 493–501, 2017. 2, 6
- [15] Bharti Munjal, Sikandar Amin, Federico Tombari, and Fabio Galasso. Query-guided end-to-end person search. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 811–820, 2019. 2, 6
- [16] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015. 6
- [17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017. 2, 4
- [18] Cheng Wang, Bingpeng Ma, Hong Chang, Shiguang Shan, and Xilin Chen. Tcts: A task-consistent two-stage framework for person search. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11952–11961, 2020. 2, 6
- [19] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pvt v2: Improved baselines with pyramid vision transformer. *Computational Visual Media*, 8(3):415–424, 2022. 2, 3, 6, 8
- [20] Jimin Xiao, Yanchun Xie, Tammam Tillo, Kaizhu Huang, Yunchao Wei, and Jiashi Feng. Ian: the individual aggregation network for person search. *Pattern Recognition*, 87:332– 340, 2019. 2, 6
- [21] Tong Xiao, Shuang Li, Bochao Wang, Liang Lin, and Xiaogang Wang. Joint detection and identification feature learning for person search. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3415–3424, 2017. 2, 5, 6
- [22] Yichao Yan, Jinpeng Li, Jie Qin, Song Bai, Shengcai Liao, Li Liu, Fan Zhu, and Ling Shao. Anchor-free person search. arXiv preprint arXiv:2103.11617, 2021. 2, 5, 6, 8
- [23] Yichao Yan, Qiang Zhang, Bingbing Ni, Wendong Zhang, Minghao Xu, and Xiaokang Yang. Learning context graph for person search. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2158–2167, 2019. 2, 6
- [24] Rui Yu, Dawei Du, Rodney LaLonde, Daniel Davila, Christopher Funk, Anthony Hoogs, and Brian Clipp. Cascade transformers for end-to-end person search. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 7267–7276, 2022. 2, 6, 8
- [25] Liang Zheng, Hengheng Zhang, Shaoyan Sun, Manmohan Chandraker, Yi Yang, and Qi Tian. Person re-identification in the wild. In *Proceedings of the IEEE Conference on Com-*

*puter Vision and Pattern Recognition*, pages 1367–1376, 2017. 2, 5, 6, 7, 8

- [26] Yingji Zhong, Xiaoyu Wang, and Shiliang Zhang. Robust partial matching for person search in the wild. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 6827–6835, 2020. 6
- [27] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. 3, 4, 5