# Cascade Transformers for End-to-End Person Search

Rui Yu[1,2][*],     Dawei Du[1],     Rodney LaLonde[1],     Daniel Davila[1],
Christopher Funk[1],     Anthony Hoogs[1],     Brian Clipp[1]
[1]Kitware, Inc., NY & NC, USA,     [2]Pennsylvania State University, PA, USA

https://github.com/Kitware/COAT

## Abstract

*The goal of person search is to localize a target person from a gallery set of scene images, which is extremely challenging due to large scale variations, pose/viewpoint changes, and occlusions. In this paper, we propose the Cascade Occluded Attention Transformer (COAT) for end-to-end person search. Our three-stage cascade design focuses on detecting people in the first stage, while later stages simultaneously and progressively refine the representation for person detection and re-identification. At each stage the occluded attention transformer applies tighter intersection over union thresholds, forcing the network to learn coarse-to-fine pose/scale invariant features. Meanwhile, we calculate each detection's occluded attention to differentiate a person's tokens from other people or the background. In this way, we simulate the effect of other objects occluding a person of interest at the token-level. Through comprehensive experiments, we demonstrate the benefits of our method by achieving state-of-the-art performance on two benchmark datasets.*

## 1. Introduction

Person search aims to localize a particular target person from a gallery set of scene images, which is an extremely difficult fine-grained recognition and retrieval problem. A person search system must both *generalize* to separate people from the background, and *specialize* to discriminate identities from each other.

In real-world applications, person search systems must detect people across a wide variety of image sizes and re-identify people despite large changes in resolution and viewpoint. To this end, modern person search methods, either two-step or one-step (*i.e.*, end-to-end), consist of reliable person detection and discriminative feature embedding learning. Two-step methods [5, 10, 13, 18, 30, 38] conduct person re-identification (ReID) on cropped person patches

---

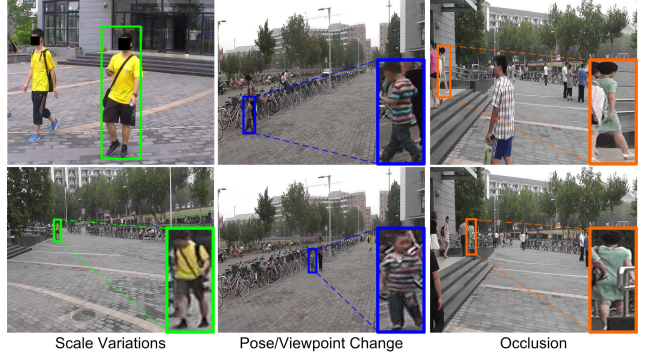[*]Rui Yu's work on this paper was done when he was a summer intern at Kitware.



Figure 1. Main challenges of person search, *e.g.*, scale variations, pose/viewpoint change, and occlusion. The boxes with the same color represent the same ID. For better viewing, we highlight the small-scale individuals at bottom-right corners.

found by a separate object detector. In contrast, end-to-end methods [2, 20, 32–34, 39] jointly solve the detection and ReID sub-problems in a more efficient, multi-task learning framework. However, as shown in Figure 1, they still suffer from three main challenges:

- *There is a conflict in feature learning between person detection and ReID.* Person detection aims to learn features which generalize across people to distinguish people from the background, while ReID aims to learn features which do *not* generalize across people but distinguish people from each other. Previous works follow a "ReID first" [33] or "detection first" [20] principle to give priority to one subtask over the other. However, it is difficult to balance the importance of two subtasks in different situations when relying on either strategy.

- *Significant scale or pose variations increase identity recognition difficulty*; see Figure 1. Feature pyramids or deformable convolutions [14, 18, 33] have been used to solve scale, pose or viewpoint misalignment in feature learning. However, simple feature fusion strategies may introduce additional background noise in feature embeddings, resulting in inferior ReID performance.

- *Occlusions with background objects or other people make appearance representations more ambiguous*, as shown in Figure 1. The majority of previous person search meth-
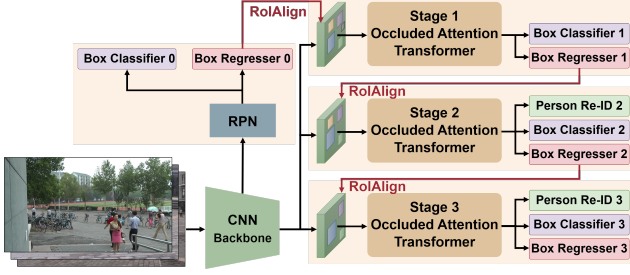
Figure 2. Our proposed cascade framework for person search.

ods focus on holistic appearance modeling of people by anchor-based [20] or anchor-free [33] methods. Despite the improvement of person search accuracy, these are prone to fail with complex occlusions.

To deal with the aforementioned challenges, as shown in Figure 2, we propose a new Cascade Occluded Attention Transformer (COAT) for end-to-end person search. First, inspired by Cascade R-CNN [1], we refine the person detection and ReID quality by a coarse-to-fine strategy in three stages. The first stage focuses on discriminating people from background (detection), but crucially, is not trained to discriminate people from each other (ReID) with a ReID loss. Later stages include both detection and ReID losses. This design improves detection performance (see Section 4.3), as the first stage can generalize across people without having to discriminate between persons. Subsequent stages simultaneously refine the previous stages' bounding box estimates and identity embeddings (see Table 1). Second, we apply multi-scale convolutional transformers at each stage of the cascade. The base feature maps are split into multiple slices corresponding to different scales. The transformer attention encourages the network to learn embeddings on the discriminative parts of each person for each scale, helping overcome the problem of region misalignment. Third, we augment the transformer's learned feature embeddings with an occluded attention mechanism that synthetically mimics occlusions . We randomly mix-up partial tokens of instances in a mini-batch, and learn the cross-attention among the token bank for each instance. This trains the transformer differentiate tokens from other foreground and background detection proposals. Experiments on the challenging CUHK-SYSU [32] and PRW [38] datasets show that the proposed network outperforms state-of-the-art end-to-end methods, especially in terms of the cross-camera setting on the PRW dataset.

**Contributions.** 1) To our knowledge, we propose the first cascaded transformer-based framework for end-to-end person search. The progressive design effectively balances person detection and ReID and the transformers help attend to scale and pose/viewpoint changes. 2) We improve performance with an occluded attention mechanism in the multi-scale transformer that generates discriminative fine-grained

person representations in occluded scenes. 3) Extensive experiments on two datasets show the superiority of our method over existing person search approaches.

## 2. Related Work

**Person Search.** Person search methods can be roughly grouped into two-step and end-to-end approaches. Two-step methods [5, 10, 13, 18, 30] combine a person detector (*e.g.*, Faster R-CNN [27], RetinaNet [22], or FCOS [28]) and a person ReID model sequentially. For example, Wang *et al.* [30] build a person search system including an identity-guided query detector followed by a detection results adapted ReID model. On the other hand, end-to-end methods [6, 20, 32, 33] integrate the two models into a unified framework for better efficiency. Chen *et al.* [6] share detection and ReID features but decompose them in the polar coordinate system in terms of radial norm and angle. Yan *et al.* [33] propose the first anchor-free person search method, which tackles the misalignment issues in different levels (*i.e.*, scale, region, and task). Recently, Li and Miao [20] share the stem representations of person detection and ReID, but solve the two subtasks by two-head networks sequentially. In contrast, inspired by Cascade R-CNN [1], our method follows an end-to-end strategy that balances person detection and ReID progressively via a three-stage cascade framework.

**Visual Transformers in Person ReID.** Based on the original transformer model [29] for natural language processing, Vision Transformer (ViT) [11] is the first pure transformer network to extract features for image recognition. CNNs are widely adopted to extract base features and so reduce the scale of training data required for a pure transformer approach. Luo *et al.* [25] develop a spatial transformer network to sample an affined image from the holistic image to match a partial image. Li *et al.* [19] propose the part-aware transformer to perform occluded person Re-ID through diverse part discovery. Zhang *et al.* [36] introduce a transformer-based feature calibration to integrate large scale features as a global prior. Our paper is the first in the literature to perform person search with multi-scale convolutional transformers . It not only learns discriminative ReID features but also distinguishes people from the background in a cascade pipeline.

**Attention Mechanism in Transformers.** Attention mechanism plays a crucial role in transformers. Recently, many ViT variants [3, 16, 21, 35] have computed discriminative features using a variety of token attention methods. Chen *et al.* [3] propose a dual-branch transformer with a cross-attention based token fusion module to combine two scales of patch features. Lin *et al.* [21] alternate attention in the feature map patches for local representation and attention on the single channel feature map for global representation. Yuan *et al.* [35] introduce the tokens-to-token process to

gradually tokenize images to tokens while preserving structural information. He *et al.* [16] rearrange the transformer layers' patch embeddings via shift and patch shuffle operations. Unlike these methods that rearrange features within an instance, the proposed occluded attention module considers token cross-attention between either positive or negative instances from the mini-batch. Thus our method learns to differentiate tokens from other objects by synthetically mimicking occlusions.

## 3. Cascade Transformers

As discussed in previous works [14, 20, 33], person detection and person ReID have conflicting goals. Hence, it is difficult to jointly learn discriminative unified representations for the two subtasks on the top of the backbone network. Similar to Cascade R-CNN [1], we decompose feature learning into sequential steps in $T$ stages of multi-scale transformers. That is, each head in the transformer refines the detection and ReID accuracy of the predicted objects stage-by-stage. Thus we can progressively learn coarse-to-fine unified embeddings.

Nevertheless, in the case of occlusions by other people, objects or the background, the network may suffer from noisy representations of the target identity. To this end, we develop the occluded attention mechanism in the multi-scale transformer to learn an occlusion-robust representation. As shown in Figure 2, our network is based on the Faster R-CNN object detector backbone with Region Proposal Network (RPN). However, we extend the framework by introducing a cascade of occluded attention transformers (see Figure 3), trained in an end-to-end manner.

### 3.1. Coarse-to-fine Embeddings

After extracting the 1024-dim stem feature maps from the ResNet-50 [15] backbone, we use the RPN to generate region proposals. For each proposal, the RoI-Align operation [27] is applied to pool an $h \times w$ region as the base feature maps $\mathcal{F}$, where $h$ and $w$ denote the height and width of the feature maps respectively, and $c$ is the number of channels.

Afterwards, we employ a multi-stage cascade structure to learn embeddings for person detection and ReID. The output proposals of the RPN are used at the first stage for re-sampling both positive and negative instances. The box outputs of the first stage are then adopted as the inputs of the second stage, and so forth. At each stage $t$, the pooled feature map of each proposal is sent to the convolutional transformers for that stage. To obtain high-quality instances, the cascade structure imposes progressively more strict stage-wise constraints. In practice, we increase the intersection-over-union (IoU) thresholds $u_t$ gradually. The transformers at each stage are followed by three heads, like NAE [6], including a person/background classifier, a box regressor,

and a ReID discriminator. Note that we remove the ReID discriminator at the first stage to focus the network on first detecting all people in the scene before refinement.

### 3.2. Occluded Attention Transformer

In the following, we describe the details of the occluded attention transformers, shown in Figure 3.

**Tokenization.** Given the base feature map $\mathcal{F} \in \mathbb{R}^{h \times w \times c}$, we tokenize it for transformer input at different scales. For multi-scale representation, we first split $\mathcal{F}$ channel-wise into $n$ slices, $\bar{\mathcal{F}} \in \mathbb{R}^{h \times w \times \hat{c}}$, where $\hat{c} = \frac{c}{n}$ to deal with each scale of token. In contrast to ViT [11] with its tokenization of large image patches, our transformer leverages a series of convolutional layers to generate tokens based on the sliced feature maps $\bar{\mathcal{F}}$. Our method benefits from CNNs' inductive biases and learns the CNN's local spatial context. The different scales are realized by different sizes of convolutional kernels.

After converting the sliced feature maps $\bar{\mathcal{F}} \in \mathbb{R}^{h \times w \times \hat{c}}$ to the new token map $\hat{\mathcal{F}} \in \mathbb{R}^{\hat{h} \times \hat{w} \times \hat{c}}$ by one convolutional layer, we flatten it into tokens inputs $\mathbf{x} \in \mathbb{R}^{\hat{h}\hat{w} \times \hat{c}}$ for one instance. The number of tokens calculated as

$$N = \frac{\hat{h}\hat{w}}{d^2} = \frac{\lfloor \frac{h+2p-k}{s} + 1 \rfloor \times \lfloor \frac{w+2p-k}{s} + 1 \rfloor}{d^2}, \quad (1)$$

where we have the kernel size $k$, stride $s$, and padding $p$ for the convolutional layer. $d$ is the patch size of each token.

**Occluded attention.** To handle occlusions, we introduce a new token-level occluded attention mechanism into the transformers to mimic occlusions found in real applications. Specifically, we first collect the tokens from all the detection proposals in a mini-batch, denoted as *token bank* $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_P\}$, where $P$ is the number of detection proposals in the batch at each stage. Since the proposals the from RPN contain positive and negative examples, the token bank is composed of both foreground person parts and background objects. We exchange tokens among the token bank, based on the same exchange index set $\mathcal{M}$ for all the instances. As shown in Figure 3, the exchanged tokens correspond to a semantically consistent but randomly selected sub-regions in the token maps. Each exchanged token is denoted as

$$\mathbf{x}_i = \{\mathbf{x}_i(\bar{\mathcal{M}}), \mathbf{x}_j(\mathcal{M})\}, \quad i = 1, 2, \cdots, P, i \neq j, \quad (2)$$

where $\mathbf{x}_j$ denotes another sample randomly selected from the token bank. $\bar{\mathcal{M}}$ indicates the complementary set of $\mathcal{M}$, *i.e.*, $\mathbf{x}_i = \mathbf{x}_i(\bar{\mathcal{M}}) \bigcup \mathbf{x}_i(\mathcal{M})$. Given the exchanged token bank $\mathbf{X}$, we compute the multi-scale self-attention among them, as shown in Figure 3. In terms of each scale of tokens, we run two sub-layers of the transformers (*i.e.*, Multi-head Self-Attention (MSA) and a Feed Forward Network (FFN) as in [29]). Specifically, the mixed tokens $\mathbf{x}$ are transformed
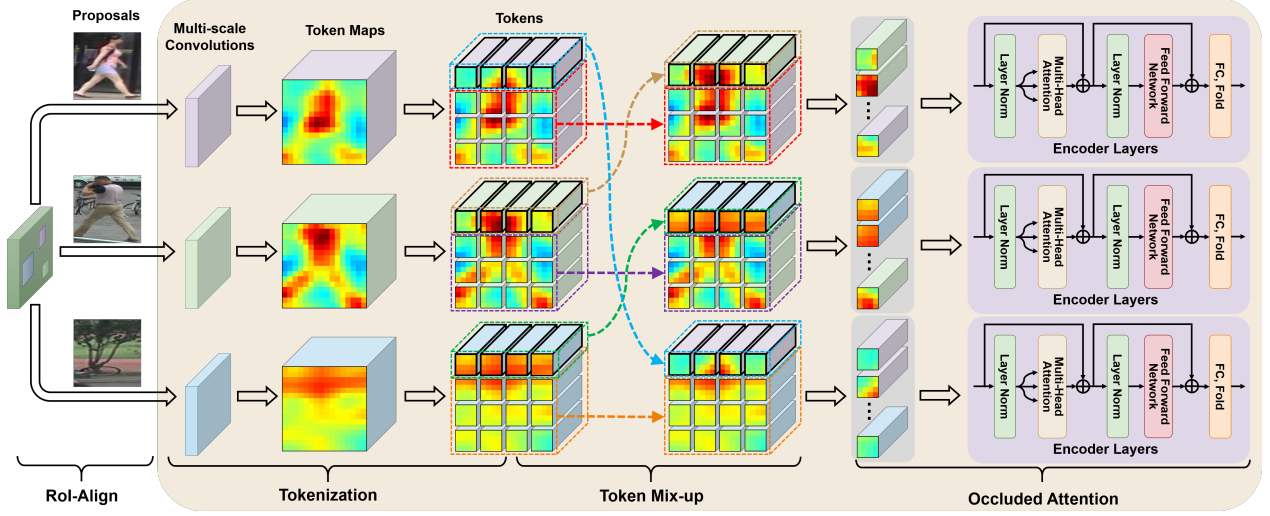
Figure 3. Architecture of occluded attention transformer. The randomly selected regions for token exchange are the same within one mini-batch. For clarity, we only show three instances in a mini-batch and occluded attention for one scale. Best view in color.

into *query* matrices $\mathbf{Q} \in \mathbb{R}^{\hat{h}\hat{w} \times \hat{c}}$, *key* matrices $\mathbf{K} \in \mathbb{R}^{\hat{h}\hat{w} \times \hat{c}}$, and *value* matrices $\mathbf{V} \in \mathbb{R}^{\hat{h}\hat{w} \times \hat{c}}$ by three individual fully connected (FC) layers. We can further compute multi-head attention and the weighted sum over all values as

$$\mathrm{MSA}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \mathrm{softmax}(\frac{\mathbf{Q}\mathbf{K}^{\mathrm{T}}}{\sqrt{\hat{c}/m}})\mathbf{V}, \quad (3)$$

where we split queries, keys, and values into $m$ heads for more diversity, *i.e.*, from tensor with the size of $\hat{h}\hat{w} \times \hat{c}$ to $m$ pieces with the size of $\hat{h}\hat{w} \times \frac{\hat{c}}{m}$. The independent attention outputs are then concatenated and linearly transformed into the expected dimension. Following the MSA module, the FFN module nonlinearly transforms each token to enhance its representation ability. The enhanced feature is then projected to the size of $\hat{h} \times \hat{w} \times \hat{c}$ as the transformer's output.

Finally, we concatenate the outputs of the $n$ scales of transformers to original spatial size $\hat{h} \times \hat{w} \times c$. Note that there is a residual connection outside each transformer. After the global average pooling (GAP) layer, the extracted features are fed into subsequent heads for box regression, person/background classification, and person re-identification.

**Relations to concurrent works.** There are two concurrent ViT based works [3, 16] in different fields. Chen *et al.* [3] develop a multi-scale transformer including two separate branches with small-patch and large-patch tokens. The two-scale representation is learned based on a cross-attention token fusion module, where a single token for each branch is treated as a query to exchange information with other branches. Instead, we leverage a series of convolutional layers with different kernels to generate multi-scale tokens. Finally, we concatenate the enhanced feature maps corresponding to each scale in specific slice of the transformers.

To deal with occlusion and misalignment in person ReID, He *et al.* [16] shuffle person part patch embeddings and re-group them, each group of which contains several random patch embeddings of an individual instance. In contrast, our method first exchanges partial tokens of instances in a mini-batch, and then calculate the occluded attention based on mixed tokens. Thus the final embeddings partially cover the target person with extracted features from a different person or a background object, yielding more occlusion-robust representations.

### 3.3. Training and Inference

In the training phase, the proposed network is trained end-to-end for person detection and person ReID. The person detection loss $\mathcal{L}_{\mathrm{det}}$ consists of regression and classification loss terms. The former is a Smooth-L1 loss of regression vectors between ground-truth and foreground boxes, while the latter computes the cross-entropy loss of predicted classification probabilities of the estimated boxes.

To supervise person ReID, we use the classic non-parametric Online Instance Matching (OIM) loss [32] $\mathcal{L}_{\mathrm{OIM}}$, which maintains a lookup table (LUT) and a circular queue (CQ) to store the features of all the labeled and unlabeled identities from recent mini-batches, respectively. We can efficiently compute the cosine similarities between the samples in the mini-batch and LUT/CQ for embedding learning. Moreover, inspired by [24], we add another cross-entropy loss function $\mathcal{L}_{\mathrm{ID}}$ to predict the identities of people for an additional ID-wise supervision. In summary, we train the proposed COAT by using the following multi-stage loss:

$$\mathcal{L} = \sum_{t=1}^{T} \mathcal{L}_{\mathrm{det}}^{t} + \mathbb{I}(t > 1)(\lambda_{\mathrm{OIM}}\mathcal{L}_{\mathrm{OIM}}^{t} + \lambda_{\mathrm{ID}}\mathcal{L}_{\mathrm{ID}}^{t}), \quad (4)$$

4

where $t \in \{1, 2, \ldots, T\}$ denotes the index of the stage and $T$ is the number of cascade stages. The coefficients $\lambda_{\text{OIM}}$ and $\lambda_{\text{ID}}$ are used to balance the OIM and ID loss terms. $\mathbb{I}(t > 1)$ is the indicator function to indicate that we do not consider person ReID loss at the first stage.

In the inference phase, we replace the occluded attention mechanism with the classic self-attention module in the transformers by removing the token mix-up step in Figure 3. We output the detection bounding boxes with corresponding embeddings at the last stage and use NMS operations to remove redundant boxes.

## 4. Experiments

All experiments are conducted in PyTorch with one NVIDIA A100 GPU. For a fair comparison with prior works, we use the first four residual blocks (`conv1`∼`conv4`) of ResNet-50 [15] as the backbone and resize the images to $900 \times 1500$ as the input.

### 4.1. Datasets

We evaluate our method on two publicly available datasets. The **CUHK-SYSU** dataset [32] annotates $8,432$ identities and $96,143$ bounding boxes in $18,184$ images. The default gallery size is set as $100$ for the $2,900$ testing identities in $6,978$ images. The **PRW** dataset [38] collects data from 6 cameras, including $932$ identities and $43,110$ pedestrian boxes in $11,816$ frames. PRW is divided into a training set with $5,704$ frames and $482$ identities and a testing set with $2,057$ query persons in $6,112$ frames.

We follow the standard evaluation metrics for person search [32, 38]. A box is matched if the overlap ratio between the predicted and ground-truth boxes with the same identity is more than $0.5$ IoU. For person detection, we use Recall and Average Precision (AP). For person ReID, we use the mean Average Precision (mAP) and cumulative matching characteristics (top-1) scores.

### 4.2. Implementation Details

Similar to Cascade R-CNN [1], we use $T = 3$ stages in the cascade framework, where 128 detection proposals are extracted per image for each stage. Following [6, 20, 32], the scale of the base feature map is set as $h = w = 14$. The index of exchanging tokens in Eq. (2) is set as the random horizontal or vertical strip in the token map. The number of heads in Eq. (3) is set as $m = 8$. The IoU thresholds $u_t$ for detection are set as $0.5, 0.6, 0.7$ for the three sequential stages. The kernel sizes of the convolutional layers to compute the tokens are set as $k = \{1 \times 1, 3 \times 3\}$ for the three stages, with corresponding strides $s = \{1, 1\}$ and paddings $p = \{0, 1\}$ to guarantee the same size of output feature maps. Due to the small feature size, we set $d = 1$ in Eq. (2), *i.e.*, conducting pixel-wise tokenization. The CQ

| Stage1 | Stage2 | Stage3 | mAP | top-1 |
|--------|--------|--------|-----|-------|
| *(a) w/o Transformers*: | | | | |
| ✗ | | | 43.5 | 81.2 |
| †✗ | ✗ | | 47.7 | 84.6 |
| †✗ | †✗ | ✗ | 48.4 | 85.2 |
| †✗ | ✗ | ✗ | 49.5 | 85.5 |
| ✗ | ✗ | ✗ | 47.2 | 84.9 |
| *(b) w/ Transformers*: | | | | |
| ✓ | | | 43.3 | 78.7 |
| †✓ | ✓ | | 50.8 | 84.9 |
| †✓ | †✓ | ✓ | 51.3 | 85.5 |
| †✓ | ✓ | ✓ | **53.3** | **87.4** |
| ✓ | ✓ | ✓ | 50.3 | 84.0 |
| *(c) IoU Thresholds*: | | | | |
| 0.5 | 0.5 | 0.5 | 52.5 | 86.0 |
| 0.6 | 0.6 | 0.6 | 52.6 | 86.2 |
| 0.7 | 0.7 | 0.7 | 51.0 | 85.5 |
| 0.5 | 0.6 | 0.6 | 52.6 | 86.3 |
| 0.5 | 0.6 | 0.7 | **53.3** | **87.4** |

Table 1. Comparison with different cascade variants of COAT on PRW [38]. "✗" means using the same ResNet block (`conv5`) as [6, 20, 32], while "✓" means using the proposed transformers at each stage. "†" means the heads without the ReID loss. Gray highlighting indicates the parameters selected for our final system.

size of the OIM loss is set as $5,000$ and $500$ for CUHK-SYSU and PRW respectively. The loss weights in Eq. (4) are set as $\lambda_{\text{OIM}} = \lambda_{\text{ID}} = 0.5$.

We use the SGD optimizer with momentum $0.9$ to train our model for 15 epochs, with an initial learning rate warming up to $0.003$ during the first epoch, being reduced by a factor of 10 at the 10-th epoch. At the inference phase, we use NMS with $0.4/0.4/0.5$ threshold to remove redundant boxes detected by the first/second/third stage.

### 4.3. Ablation Studies

We conduct a series of ablation studies on the PRW dataset [38] to analyze our design decisions.
**Contribution of cascade structure.** To show the cascade structure's contribution, we evaluate coarse-to-fine constraints in terms of the number of cascade stages and IoU thresholds.

First, we replace the occluded attention transformer with the same ResNet block (`conv5`) as [6, 20, 32] at each stage. As shown in Table 1(a), the cascade structure significantly improves person search accuracy when adding more stages, *i.e.*, from $43.5\%$ to $49.5\%$ in mAP and $81.2\%$ to $85.5\%$ in top-1 accuracy. As we introduce the proposed occluded attention transformer, the performance is further improved (see Table 1(b)), which demonstrates our occluded attention transformer's effectiveness .

Moreover, the increasing IoU thresholds $u_t$ in the cascade design improve person search performance. As reported in Table 1(c), equal IoU thresholds at each stage pro-
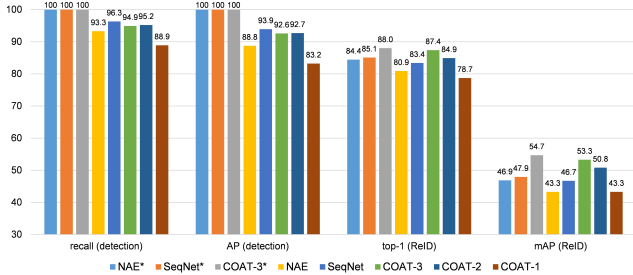
Figure 4. Detection and person search results for COAT and two compared methods on PRW, both with (person ReID only) and without (person search) ground-truth detection boxes being provided. * denotes the oracle results using the ground-truth boxes.

| Method | Tokens | Feats | mAP | top-1 |
|---|:---:|:---:|:---:|:---:|
| Vanilla Attention | | | 52.9 | 86.4 |
| CrossViT [3] | ✓ | | 49.9 | 86.1 |
| Jigsaw [16] | ✓ | | 51.9 | 86.0 |
| Batch DropBlock [7] | | ✓ | 52.7 | 86.7 |
| Cutout [8] | | ✓ | 53.2 | 86.6 |
| Mixup [37] | | ✓ | 52.8 | 86.6 |
| Occluded Attention | ✓ | | **53.3** | **87.4** |

Table 2. Comparison of our attention mechanisms and other related modules. "Tokens" and "Feats" denote token-level enhanced attention and feature-level augmentation respectively.

duce lower accuracy than our method. For example, more false positives or false negatives are introduced if $u_t = 0.5$ or $u_t = 0.7$. In contrast, our method can select detection proposals with increasing quality for better performance, *i.e.*, generating more candidate detections in the first stage and only highly-overlapping detections by the third stage.

**Relations between person detection and ReID.** As discussed in the introduction, there is a conflict between person detection and ReID. In Figure 4, we explore the relationship between the two subtasks. We compare our COAT with state-of-the-art NAE [6] and SeqNet [20], which share the same Faster R-CNN detector. We also construct three COAT variants with different stages, *i.e.*, COAT-$t$, where $t = 1, 2, 3$ denotes the number of stages. When looking solely at person ReID rather than person search, *i.e.*, when ground-truth detection boxes are given, COAT outperforms the two competitors with an over 3% gain in top-1 and over 6% gain in mAP. Meanwhile, our is slightly worse in person detection accuracy than SeqNet [20]. These results indicate that our improved ReID performance comes from coarse-to-fine person embeddings rather than more precise detections.

We also observe that the person detection performance is improved from $t = 1$ to $t = 2$ but then slightly reduced with $t = 3$. We speculate that this is because, when trading-off person detection and ReID, our method focuses more on learning discriminative embeddings for person ReID, while slightly sacrificing detection performance.

In addition, from Table 1(a)(b), note that the COAT variant with ReID loss in the first stage performs worse than our method (50.3 vs. 53.3 for mAP). Simultaneously learning a discriminative representation for person detection and ReID is extremely difficult. Therefore, we remove the ReID discriminator head at Stage 1 in the COAT method (*c.f.* Figure 2). If we continue removing the ReID discriminator at the second stage, the ReID performance is reduced by $\sim 2\%$ in mAP. This shows the ReID embeddings do benefit from multi-stage refinement.

**Comparison with other attention mechanisms.** To verify the effectiveness of our occluded attention mechanism in

the transformer, we apply the recently proposed Jigsaw [16] and CrossViT [3] in our method. As discussed in Section 3.2, Jigsaw Patch [16] is used to generate robust ReID features by shift and patch shuffle operations. CrossViT [3] is a dual-branch transformer to learn multi-scale features. It is also noteworthy that they leverage large image patches as the input for pure vision transformers. We also evaluate the COAT variant a vanilla self-attention mechanism, denoted as vanilla attention.

In Table 2, CrossViT [3] focuses on exchanging information between two scales of tokens, achieving inferior mAP. The results show that Jigsaw [16] also hurts mAP. We speculate that either exchanging query information in CrossViT [3] or the shift and shuffle feature operations in Jigsaw [16] are ambiguous in such small $14 \times 14$ base feature maps, limiting the power of them for person search. In contrast, our occluded attention is designed for small feature maps and obtains better performance, *i.e.*, both $0.4\%$ gain in mAP and $1.0\%$ gain in top-1 score. Instead of sharing class tokens in different branches or shuffling channels of feature maps based on an individual instance, we effectively learn context information across different instances in a mini-batch, and differentiate the person from other people or the background to synthetically mimic occlusion.

**Comparison with feature augmentation.** Our method is related to previous augmentation strategies for person ReID, such as Batch DropBlock Network [7], Cutout [8] and Mixup [37]. As presented in Table 2, person search accuracy is not improved by using feature augmentation, simply augmenting feature patches with zeros.

**Influence of occluded attention mechanism.** As discussed in Section 3.2, we use occluded attention to calculate discriminative person embeddings. We evaluate the use of occluded attention (token mixup) and different scales in Table 3. Note, the top-1 score is improved from $86.4$ to $87.4$ with occluded attention and that multiple convolutional kernels for tokenization improve performance. Note that multiple convolutions do not increase the model size, since the feature maps $\mathcal{F}$ are channel-wise sliced for each scale.

6

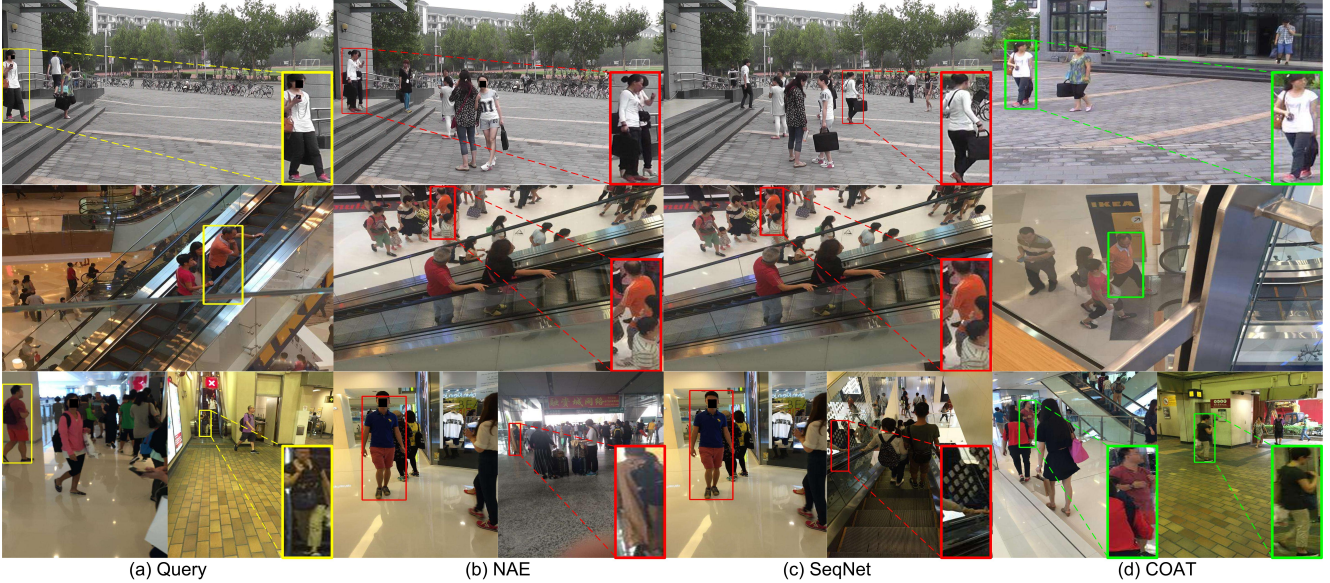(a) Query      (b) NAE      (c) SeqNet      (d) COAT

Figure 5. Qualitative examples of top-1 person search results of NAE [6], SeqNet [20] and COAT on PRW (1st row) and CUHK-SYSU (2nd and 3rd rows) datasets, where small query, failure and correct cases are highlighted in yellow, red and green boxes respectively.

| Method | Token Mixup | Scales | mAP | top-1 |
|--------|:-----------:|:------:|:---:|:-----:|
| Vanilla Attention | | $\{1 \times 1\}$ | 52.1 | 85.3 |
| Vanilla Attention | | $\{3 \times 3\}$ | 53.1 | 86.0 |
| Vanilla Attention | | $\{1 \times 1, 3 \times 3\}$ | 52.9 | 86.4 |
| Occluded Attention | ✓ | $\{1 \times 1\}$ | 52.2 | 86.5 |
| Occluded Attention | ✓ | $\{3 \times 3\}$ | 52.5 | 86.4 |
| Occluded Attention | ✓ | $\{1 \times 1, 3 \times 3\}$ | **53.3** | **87.4** |

Table 3. Comparison of our attention mechanisms and other related modules. "Scales" denotes the used convolutional kernels.

### 4.4. Comparison with State-of-the-art

As presented in Table 4, we compare our COAT with state-of-the-art algorithms, including both two-step methods [5, 10, 13, 18, 30, 38] and end-to-end methods [2, 4, 6, 9, 12, 17, 20, 23, 26, 31–34, 39], on two datasets.

**Results on CUHK-SYSU.** With the gallery size of 100, our method achieves the best 94.2% mAP and comparable 94.7% top-1 scores compared to the best two-step method TCTS [30] with explicitly trained bounding box and ReID feature refinement modules. Among end-to-end methods, our method performs better than state-of-the-art AlignPS+ [33] with a multi-scale anchor-free representation [28], SeqNet [20] with two-stage refinement and AGWF [12] with part classification based sub-networks. The results indicate the effectiveness of our cascaded multi-scale representation. Using the post-processing operation Context Bipartite Graph Matching (CBGM) [20], both mAP and top-1 scores of our method can be further improved slightly. For a comprehensive evaluation, as shown in Figure 6, we compare mAP scores of competitive methods as we increase gallery size. Since it is challenging to consider more distracting people in the gallery set, the performance



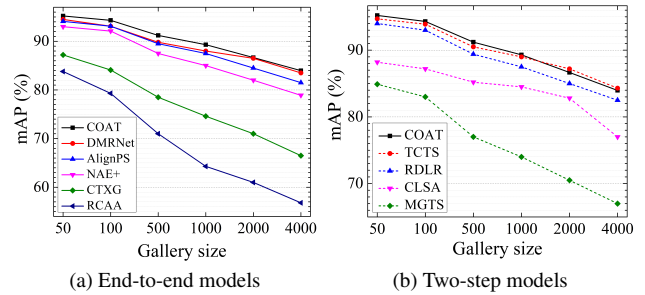(a) End-to-end models      (b) Two-step models

Figure 6. Comparison with (a) end-to-end models and (b) two-step models on CUHK-SYSU with different gallery sizes.

of all compared methods is reduced as the gallery size increases. However, our method consistently outperforms all the end-to-end methods and the majority of two-step methods. When the gallery size is larger than $1,000$, our method performs slightly worse than the two-step TCTS [30].

**Results on PRW.** Although the PRW dataset [38] is more challenging, with less training data but larger gallery size, than the CUHK-SYSU dataset [32], the results show a similar trend. Our method achieves comparable performance as AGWF [12] and a significant gain of 6.7% mAP and 4.0% top-1 scores than SeqNet [20]. DMRNet [14] and AlignPS [33] leverage stronger object detectors, such as RetinaNet [22] and FCOS [28], than the Faster R-CNN [27] in our method, but still achieve inferior performance. Further, we compare performance on PRW's multi-view gallery (see the group marked by † in Table 4). Our method outperforms existing methods in terms of both mAP and Top-1 scores with a clear margin. We attribute this to our cascaded transformer structure which generates more discriminative ReID features, especially in the cross-camera setting with significant pose/viewpoint changes.

| Method | CUHK-SYSU | | PRW | |
|---|---|---|---|---|
| | mAP | top-1 | mAP | top-1 |
| **Two-step** | | | | |
| DPM [38] | - | - | 20.5 | 48.3 |
| MGTS [5] | 83.0 | 83.7 | 32.6 | 72.1 |
| CLSA [18] | 87.2 | 88.5 | 38.7 | 65.0 |
| RDLR [13] | 93.0 | 94.2 | 42.9 | 70.2 |
| IGPN [10] | 90.3 | 91.4 | **47.2** | 87.0 |
| TCTS [30] | **93.9** | **95.1** | 46.8 | 87.5 |
| **End-to-end** | | | | |
| OIM [32] | 75.5 | 78.7 | 21.3 | 49.9 |
| IAN [31] | 76.3 | 80.1 | 23.0 | 61.9 |
| NPSM [23] | 77.9 | 81.2 | 24.2 | 53.1 |
| RCAA [2] | 79.3 | 81.3 | - | - |
| CTXG [34] | 84.1 | 86.5 | 33.4 | 73.6 |
| QEEPS [26] | 88.9 | 89.1 | 37.1 | 76.7 |
| HOIM [4] | 89.7 | 90.8 | 39.8 | 80.4 |
| APNet [39] | 88.9 | 89.3 | 41.9 | 81.4 |
| BINet [9] | 90.0 | 90.7 | 45.3 | 81.7 |
| NAE [6] | 91.5 | 92.4 | 43.3 | 80.9 |
| NAE+ [6] | 92.1 | 92.9 | 44.0 | 81.1 |
| DMRNet [14] | 93.2 | 94.2 | 46.9 | 83.3 |
| PGS [17] | 92.3 | **94.7** | 44.2 | 85.2 |
| AlignPS [33] | 93.1 | 93.4 | 45.9 | 81.9 |
| AlignPS+ [33] | 94.0 | 94.5 | 46.1 | 82.1 |
| SeqNet [20] | 93.8 | 94.6 | 46.7 | 83.4 |
| AGWF [12] | 93.3 | 94.2 | **53.3** | **87.7** |
| COAT | **94.2** | **94.7** | **53.3** | 87.4 |
| AlignPS [33]+CBGM [20] | 93.6 | 94.2 | 46.8 | 85.8 |
| AlignPS+ [33]+CBGM [20] | 94.2 | 94.3 | 46.9 | 85.7 |
| SeqNet+CBGM [20] | **94.8** | **95.7** | 47.6 | 87.6 |
| COAT+CBGM | **94.8** | 95.2 | **54.0** | **89.1** |
| HOIM† [4] | - | - | 36.5 | 65.0 |
| NAE+† [6] | - | - | 40.0 | 67.5 |
| SeqNet† [20] | - | - | 43.6 | 68.5 |
| SeqNet+CBGM† [20] | - | - | 44.3 | 70.6 |
| AGWF† [12] | - | - | 48.0 | 73.2 |
| COAT† | - | - | 50.9 | 75.1 |
| COAT+CBGM† | - | - | **51.7** | **76.1** |

Table 4. Comparison with the state-of-the-art methods. † denotes the performance only evaluated on the multi-view gallery. Bold indicates highest score in the group.

| Method | Params(M) | MACs(G) | FPS | mAP | top-1 |
|---|---|---|---|---|---|
| NAE [6] | 33.43 | 287.35 | 14.48 | 43.3 | 80.9 |
| AlignPS [33] | 42.18 | 189.98 | 16.39 | 45.9 | 81.9 |
| SeqNet [20] | 48.41 | 275.11 | 12.23 | 46.7 | 83.4 |
| COAT | 37.00 | 236.29 | 11.14 | **53.3** | **87.4** |

Table 5. Comparison of person search efficiency.

efficiency.

## 5. Conclusion

We have developed a new Cascade Occluded Attention Transformer (COAT) for end-to-end person search. Notably, COAT learns a discriminative coarse-to-fine representation for both person detection and person ReID via a cascade transformer framework. Meanwhile, the occluded attention mechanism synthetically mimics occlusions from either foreground or background objects. COAT outperforms state-of-the-art methods, which we hope will inspire more research into transformer-based person search methods.

**Ethical considerations.** Like most technologies, person search methods may have societal benefits and negative impacts. How the technology is employed is critical. For example, person search can identify persons of interest to aid law enforcement and counter-terrorism operations. However, the technology should only be used in locations where an expectation of privacy is waived by entering those locations, such as public areas, airports, and private buildings with clear signage. These systems should not be employed without probable cause, or by unjust governments that seek to acquire ubiquitous knowledge of the movements of all of their citizens to enable persecution and repression.

For comparability, this research uses human subjects imagery collected in prior works. CUHK-SYSU [32] was collected from "street snaps" and "movie snapshots", while PRW [38] was collected with video cameras in a public area of a university campus. No mention is made in either paper of review by an ethical board (*e.g.*, an Institutional Review Board), but these papers were published before this new standard was established at CVPR or most major AI conferences. Our preference would be to work with ethically collected person search datasets, and we would welcome a public disclosure from the authors of their ethical compliance. We believe the community should focus resources on developing ethical person search datasets and phase out the use of legacy, unethically collected datasets.

**Qualitative results.** Some example person search results on two datasets are shown in Figure 5. Our method can deal with cases of slight/moderate occlusion and scale/pose variations, while other state-of-the-art methods such as SeqNet [20] and NAE [6] fail in these scenarios.

**Efficiency comparison.** We compare our efficiency with three representative end-to-end networks including NAE [6], AlignPS [33] and SeqNet [20] which have publicly released source code. We evaluate the methods with the same scale test images and on the same GPU.

From Table 5, we compare the number of parameters, the multiply–accumulate operations (MACs), and the running speed in frames per second (FPS). Our method has lower computational complexity and slightly slower speed than other compared methods, but achieved $+6.6\%$ and $+4.0\%$ gains in mAP and top-1 accuracy respectively. In contrast to [11, 16], we employ only one encoder layer in our transformers and use multi-scale convolutions to reduce the number of channels before tokenization, increasing COAT's

# References

[1] Zhaowei Cai and Nuno Vasconcelos. Cascade R-CNN: delving into high quality object detection. In *CVPR*, pages 6154–6162, 2018. 2, 3, 5

[2] 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 *ECCV*, pages 86–102, 2018. 1, 7, 8

[3] Chun-Fu Chen, Quanfu Fan, and Rameswar Panda. Crossvit: Cross-attention multi-scale vision transformer for image classification. In *ICCV*, 2021. 2, 4, 6

[4] Di Chen, Shanshan Zhang, Wanli Ouyang, Jian Yang, and Bernt Schiele. Hierarchical online instance matching for person search. In *AAAI*, pages 10518–10525, 2020. 7, 8

[5] Di Chen, Shanshan Zhang, Wanli Ouyang, Jian Yang, and Ying Tai. Person search via a mask-guided two-stream CNN model. In *ECCV*, pages 764–781, 2018. 1, 2, 7, 8

[6] Di Chen, Shanshan Zhang, Jian Yang, and Bernt Schiele. Norm-aware embedding for efficient person search. In *CVPR*, pages 12612–12621, 2020. 2, 3, 5, 6, 7, 8

[7] Zuozhuo Dai, Mingqiang Chen, Xiaodong Gu, Siyu Zhu, and Ping Tan. Batch dropblock network for person re-identification and beyond. In *ICCV*, pages 3690–3700, 2019. 6

[8] Terrance Devries and Graham W. Taylor. Improved regularization of convolutional neural networks with cutout. *CoRR*, abs/1708.04552, 2017. 6

[9] Wenkai Dong, Zhaoxiang Zhang, Chunfeng Song, and Tieniu Tan. Bi-directional interaction network for person search. In *CVPR*, pages 2836–2845, 2020. 7, 8

[10] Wenkai Dong, Zhaoxiang Zhang, Chunfeng Song, and Tieniu Tan. Instance guided proposal network for person search. In *CVPR*, pages 2582–2591, 2020. 1, 2, 7, 8

[11] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. 2, 3, 8

[12] Byeong-Ju Han, Kuhyeun Ko, and Jae-Young Sim. End-to-end trainable trident person search network using adaptive gradient propagation. In *ICCV*, pages 925–933, 2021. 7, 8

[13] Chuchu Han, Jiacheng Ye, Yunshan Zhong, Xin Tan, Chi Zhang, Changxin Gao, and Nong Sang. Re-id driven localization refinement for person search. In *ICCV*, pages 9813–9822, 2019. 1, 2, 7, 8

[14] 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 *AAAI*, pages 1505–1512, 2021. 1, 3, 7, 8

[15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 3, 5

[16] Shuting He, Hao Luo, Pichao Wang, Fan Wang, Hao Li, and Wei Jiang. Transreid: Transformer-based object re-identification. In *ICCV*, 2021. 2, 3, 4, 6, 8

[17] Hanjae Kim, Sunghun Joung, Ig-Jae Kim, and Kwanghoon Sohn. Prototype-guided saliency feature learning for person search. In *CVPR*, pages 4865–4874, 2021. 7, 8

[18] Xu Lan, Xiatian Zhu, and Shaogang Gong. Person search by multi-scale matching. In *ECCV*, pages 553–569, 2018. 1, 2, 7, 8

[19] Yulin Li, Jianfeng He, Tianzhu Zhang, Xiang Liu, Yongdong Zhang, and Feng Wu. Diverse part discovery: Occluded person re-identification with part-aware transformer. In *CVPR*, 2021. 2

[20] Zhengjia Li and Duoqian Miao. Sequential end-to-end network for efficient person search. In *AAAI*, pages 2011–2019, 2021. 1, 2, 3, 5, 6, 7, 8

[21] Hezheng Lin, Xing Cheng, Xiangyu Wu, Fan Yang, Dong Shen, Zhongyuan Wang, Qing Song, and Wei Yuan. CAT: cross attention in vision transformer. *CoRR*, abs/2106.05786, 2021. 2

[22] Tsung-Yi Lin, Priya Goyal, Ross B. Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, pages 2999–3007, 2017. 2, 7

[23] Hao Liu, Jiashi Feng, Zequn Jie, Jayashree Karlekar, Bo Zhao, Meibin Qi, Jianguo Jiang, and Shuicheng Yan. Neural person search machines. In *ICCV*, pages 493–501, 2017. 7, 8

[24] Hao Luo, Youzhi Gu, Xingyu Liao, Shenqi Lai, and Wei Jiang. Bag of tricks and a strong baseline for deep person re-identification. In *CVPRW*, pages 1487–1495, 2019. 4

[25] Hao Luo, Wei Jiang, Xing Fan, and Chi Zhang. Stnreid: Deep convolutional networks with pairwise spatial transformer networks for partial person re-identification. *IEEE TMM*, 22(11):2905–2913, 2020. 2

[26] Bharti Munjal, Sikandar Amin, Federico Tombari, and Fabio Galasso. Query-guided end-to-end person search. In *CVPR*, pages 811–820, 2019. 7, 8

[27] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster R-CNN: towards real-time object detection with region proposal networks. *IEEE TPAMI*, 39(6):1137–1149, 2017. 2, 3, 7

[28] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. FCOS: fully convolutional one-stage object detection. In *ICCV*, pages 9626–9635, 2019. 2, 7

[29] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, pages 5998–6008, 2017. 2, 3

[30] Cheng Wang, Bingpeng Ma, Hong Chang, Shiguang Shan, and Xilin Chen. TCTS: A task-consistent two-stage framework for person search. In *CVPR*, pages 11949–11958, 2020. 1, 2, 7, 8

[31] Jimin Xiao, Yanchun Xie, Tammam Tillo, Kaizhu Huang, Yunchao Wei, and Jiashi Feng. IAN: the individual aggregation network for person search. *PR*, 87:332–340, 2019. 7, 8

[32] Tong Xiao, Shuang Li, Bochao Wang, Liang Lin, and Xiaogang Wang. Joint detection and identification feature learning for person search. In *CVPR*, pages 3376–3385, 2017. 1, 2, 4, 5, 7, 8

[33] Yichao Yan, Jingpeng Li, Jie Qin, Song Bai, Shengcai Liao, Li Liu, Fan Zhu, and Ling Shao. Anchor-free person search. In *CVPR*, 2021. 1, 2, 3, 7, 8

[34] Yichao Yan, Qiang Zhang, Bingbing Ni, Wendong Zhang, Minghao Xu, and Xiaokang Yang. Learning context graph for person search. In *CVPR*, pages 2158–2167, 2019. 1, 7, 8

[35] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Francis E. H. Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. In *ICCV*, pages 558–567, 2021. 2

[36] Guowen Zhang, Pingping Zhang, Jinqing Qi, and Huchuan Lu. Hat: Hierarchical aggregation transformers for person re-identification. In *ACMMM*, 2021. 2

[37] Hongyi Zhang, Moustapha Cissé, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *ICLR*, 2018. 6

[38] Liang Zheng, Hengheng Zhang, Shaoyan Sun, Manmohan Chandraker, Yi Yang, and Qi Tian. Person re-identification in the wild. In *CVPR*, pages 3346–3355, 2017. 1, 2, 5, 7, 8

[39] Yingji Zhong, Xiaoyu Wang, and Shiliang Zhang. Robust partial matching for person search in the wild. In *CVPR*, pages 6826–6834, 2020. 1, 7, 8