

An End-to-End Foreground-Aware Network for Person Re-Identification

Yiheng Liu, Wengang Zhou, Jianzhuang Liu, Guojun Qi, Qi Tian, *Fellow, IEEE*,
Houqiang Li, *Senior Member, IEEE*,

Abstract—Person re-identification is a crucial task of identifying pedestrians of interest across multiple surveillance camera views. In person re-identification, a pedestrian is usually represented with features extracted from a rectangular image region that inevitably contains the scene background, which incurs ambiguity to distinguish different pedestrians and degrades the accuracy. To this end, we propose an end-to-end foreground-aware network to discriminate foreground from background by learning a soft mask for person re-identification. In our method, in addition to the pedestrian ID as supervision for foreground, we introduce the camera ID of each pedestrian image for background modeling. The foreground branch and the background branch are optimized collaboratively. By presenting a target attention loss, the pedestrian features extracted from the foreground branch become more insensitive to the backgrounds, which greatly reduces the negative impacts of changing backgrounds on matching an identical across different camera views. Notably, in contrast to existing methods, our approach does not require any additional dataset to train a human landmark detector or a segmentation model for locating the background regions. The experimental results conducted on three challenging datasets, *i.e.*, Market-1501, DukeMTMC-reID, and MSMT17, demonstrate the effectiveness of our approach.

Index Terms—Person re-identification, background, end-to-end, attention.

I. INTRODUCTION

PERSON re-identification aims to match persons across non-overlapping surveillance camera views. The growing demand for video surveillance and public security has drawn increasing attention to this task. Although great advance has been witnessed in recent years, there are still many challenging issues towards its application. In this work, we are dedicated to learning robust and discriminative pedestrian features insensitive to backgrounds for effective person re-identification.

In person re-identification, a pedestrian is usually represented by a rectangular image region, which inevitably contains some background regions due to the irregular shape of pedestrians. Without precisely localizing the foreground and neglecting the background, the diverse background clutters incurs noise to the model learning and degrades the accuracy.

Yiheng Liu, Wengang Zhou and Houqiang Li are with the CAS Key Laboratory of Technology in Geo-spatial Information Processing and Application System, Department of Electronic Engineering and Information Science, University of Science and Technology of China, Hefei, 230027, China.
E-mail: lyh156@mail.ustc.edu.cn, {zhwg, lihq}@ustc.edu.cn.

Jianzhuang Liu and Qi Tian are with Huawei Noah's Ark Laboratory.
E-mail: {liu.jianzhuang, tian.qi1}@huawei.com

Guojun Qi is with Huawei Cloud EI Product Department.
E-mail: guojun.qi@huawei.com

Corresponding authors: Wengang Zhou and Houqiang Li.

To alleviate the adverse influence from the backgrounds, numerous methods [1], [2], [3], [4], [5], [6], [7], [8], [9] have been proposed. In [1], [2], [3], human landmark detectors are used to extract human keypoints and generate human part bounding boxes. In [5], [7], [6], segmentation models on pedestrian are applied to generate the whole body masks or multiple semantic regions. Thanks to the pre-trained landmark detection models and segmentation models, the body regions can be well separated from the background areas. Compared with the global features extracted from the whole images, the features from the body regions are more discriminative for person re-identification tasks without background noise. Some methods [4], [8], [9] design different attention models to help the networks focus on discriminative human body regions. The performance gain achieved by these methods demonstrates that removing the influence from the backgrounds is beneficial for person re-identification.

Although the existing methods achieve promising results on mitigating the effects of the backgrounds, they still suffer one of more of the following limitations. 1) The human landmark detection model and segmentation model need to be pre-trained with additional labeled human pose and segmentation respectively, which requires extra overhead for model training and data collection. 2) The data-bias between the source datasets and the target person re-identification datasets could deteriorate the estimation of the keypoints and body masks. In particular, the existing large person re-identification datasets are usually composed of low-resolution pedestrian images, which brings remarkable challenges to adapt the pre-trained models. 3) Limited by the data, the human landmark detection model and segmentation model are difficult to be trained together with the person re-identification model in an end-to-end manner to mutually promote each other. 4) During the inference stage, it is time-consuming to generate the keypoints and body masks for individual images by these pre-trained models. 5) Existing attention-based methods do not require additional training data, but the lack of strong supervision for training the networks makes them vulnerable in focusing the model attentions on body regions.

To address the above issues, we propose to use the camera identity information contained in the existing person re-identification datasets to help the models separate the foreground human bodies from the background regions. Some existing methods [10], [11], [12], [13] explore the camera network topology and spatiotemporal constraints between cameras to refine the person similarities. In contrast, we exploit the camera identity information to directly train the network

to learn background feature representations. By introducing the camera identity information, we can learn more discriminative person representations with the foreground branch and alleviate the negative effects of the backgrounds.

Based on the above discussion, we design an end-to-end foreground-aware network FA-Net, which aims to learn foreground-aware features effectively and efficiently. Our method contains two branches. One is the foreground feature extraction branch trained by pedestrian identity information. The other is the background feature extraction branch trained by camera identity information, which is used to constrain the target enhancement module to better distinguish foreground and background. In the inference, only the foreground branch is needed to extract pedestrian features, which is very efficient. To suppress the responses in the nontarget regions, we further propose a target attention loss, which provides strong supervision for training the target enhancement module to focus on the target regions. Unlike existing works [1], [2], [3], [5], [7], [6], our method does not require additional human pose or segmentation but still well discriminates the foregrounds from the backgrounds with promising recognition accuracy.

In the rest of this paper, we first make a survey on related works in Section II. Then, we elaborate our proposed framework in Section III. After that, we evaluate our method with extensive experiments in Section IV. Finally, we conclude this work in Section V.

II. RELATED WORKS

In this section, we first briefly introduce the progress of person re-identification. Then, we review the most related works from two aspects, *i.e.*, one with similar purpose to refine the foreground features and the other with camera clue considered.

A. Brief Overview of Person Re-identification

Most existing person re-identification works focus on two key issues: discriminative feature representation [14], [15] and effective distance measurement [16], [17], [18]. The background clutter, occlusion, and the dramatic variations in viewpoints and pedestrian postures make it critical to extract more discriminative and robust features for person re-identification. On the other hand, given the discriminative feature representation, an effective distance metric is expected to well measure the similarities between pedestrians.

In recent years, the rapid development of Convolutional Neural Networks (CNNs) has greatly promoted the advance of person re-identification. Based on CNNs, many discriminative feature learning methods [19], [20], [21], [8], [22], [23] and distance measurement methods [24], [25], [26] have been proposed. Suh *et al.* [20] design a network to learn a part-aligned representation for person re-identification. A two-stream network is adopted to extract appearance representations and part representations, which are further aggregated to generate the part-aligned features. In [21], a network named Part-based Convolutional Baseline (PCB) is proposed to extract part-level features. Shen *et al.* [24] propose a Kronecker Product Matching module to measure the similarities of the feature maps of different persons.

B. Methods towards Refining Foreground Features

To obtain robust representations, a key challenge is how to alleviate the influence of the backgrounds and make the network focus more on discriminative human bodies. In order to solve this problem, many effective methods [1], [2], [3], [4], [5], [6], [7], [8], [9] have been proposed. These methods fall mainly into the following three categories.

Human landmark detection. Zhao *et al.* [1] train a model to estimate body joint locations and obtain several body subregions. In [2], Wei *et al.* use a model pre-trained on the MPII human pose dataset [27] to estimate keypoints and crop three local body regions. In [3], a pose-driven deep convolutional (PDC) model is proposed to learn improved feature extraction and matching models, in which a human pose estimation algorithm pre-trained on human pose datasets is used to generate human keypoints.

Segmentation-based methods. Kalayeh *et al.* [5] design a SPReID model to integrate human semantic parsing in person re-identification. A human semantic parsing model is trained to segment a human body into multiple semantic regions, which are used to exploit local cues for person re-identification. Song *et al.* [7] use a mask-guided contrastive attention model, which extracts features separately from the body and background regions. A pre-trained human segmentation model is adopted to generate a binary segmentation mask corresponding to the body and background regions. Tian *et al.* [6] learn more discriminative person-part features based on human parsing maps generated by a person parsing network pre-trained on labeled human parsing datasets.

Attention-based methods. Zhao *et al.* [4] design an attention model to generate multiple part maps. In [8], a Harmonious Attention CNN (HA-CNN) model is proposed to jointly learn the soft pixel attention and the hard regional attention along with simultaneous optimization of feature representations. Wang *et al.* [9] propose a fully attentional block (FAB) to localize the most discriminative local regions for person re-identification. By applying FAB in different levels of intermediate features, they can acquire different scales of attention responses.

Different from the above works, our method aims to mitigate the influence of backgrounds in a more effective and efficient way. FA-Net does not require additional human pose or segmentation datasets but still has strong supervision information to help locate the body parts and the background parts. Meanwhile, the background feature extraction branch is trained together with the foreground feature extraction branch. This end-to-end training strategy allows the two branches to promote each other to accurately locate the body regions and extract more robust features.

C. Methods Considering Camera Information

In addition to exploiting visual information to match pedestrians, there are some methods [10], [11], [12], [13] using the spatial context of the cameras and the temporal stamp of visual frames to constrain the learning of person similarities. In [11], [13], different approaches are explored to use the spatiotemporal constraint to eliminate the irrelevant gallery

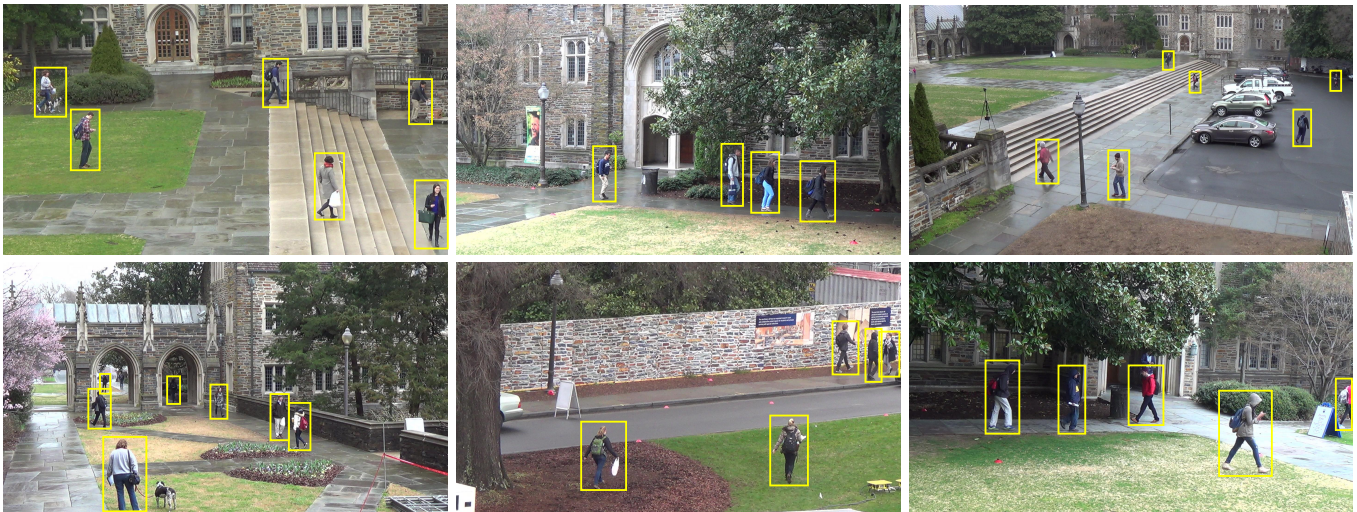


Fig. 1. Illustration of the pedestrian images in the different scenes on DukeMTMC [28]. Although the backgrounds of these pedestrian images captured by the same camera are different, they all belong to the same scene.

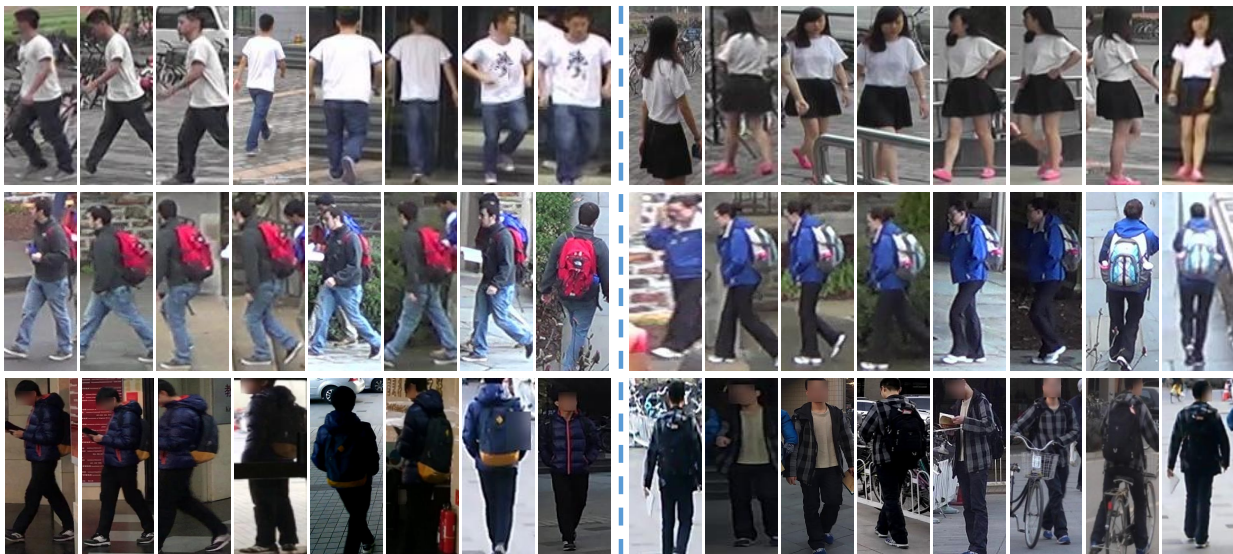


Fig. 2. The pedestrian images on difference datasets. The three lines of images are from Market-1501 [29], DukeMTMC-reID [30] and MSMT17 [31], respectively. Images belonging to the same person have complex and varied backgrounds, which make it difficult to identify persons.

images. Lv *et al.* [12] propose an unsupervised incremental learning algorithm to mine the spatio-temporal patterns using the time interval of pedestrians transferring across different cameras. In [10], a unified framework is designed, which uses the spatiotemporal relations to perform the camera network topology inference.

Different from the above methods, we utilize the camera information from a new perspective. Specifically, we directly use the camera identity information to guide the network to locate the background regions and help the person feature extraction model alleviate the effects from the backgrounds.

III. OUR METHOD

In this section, we first introduce the overall architecture of the proposed end-to-end foreground-aware network for

collaborative learning of foreground feature and background feature in Section III-A. Then, we elaborate our target enhancement module and target attention loss in Section III-B and Section III-C, respectively. Finally, we discuss our overall training objective in Section III-D.

A. Collaborative Learning of Person ID and Camera ID

In the scenario of video surveillance for person re-identification, each person is photographed by a certain camera as illustrated in Fig. 1, and detected in the form of a cropped rectangular image patch, which contains not only the person as foreground but also some portion of scene background, as shown in Fig. 2. As a result, each person image is characterized by two attributes, *i.e.*, the person ID and the camera ID. For cropped images belonging to the same person, they have

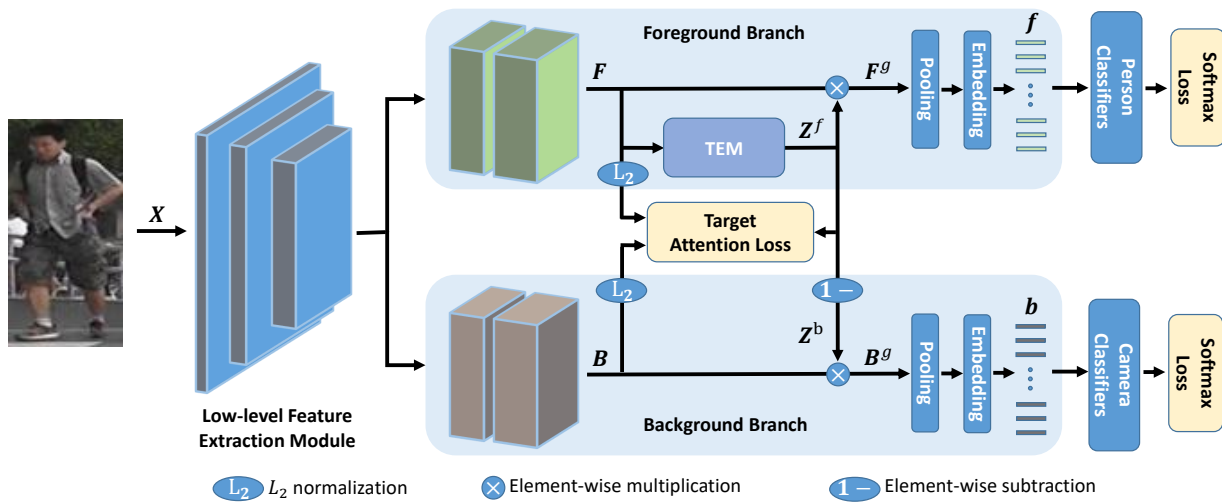


Fig. 3. The overall architecture of the proposed method. The foreground branch and the background branch are independent of each other and do not share their weights. TEM denotes the target enhancement module. The camera classifier is trained to predict which scene the background of one image belongs to. During the inference stage for person re-ID, the background branch is no longer needed.

similar foreground but different backgrounds, which indicates that, to identify the person ID, it is necessary to focus on the foreground and avoid the effect of the background. On the other hand, for pedestrian images captured by the same camera, they are detected from the same scene. The foregrounds, *i.e.*, the pedestrians are usually changing, but the backgrounds are parts of the same scene and share the same camera identity. Therefore, to identify the camera identity of a pedestrian image, we should pay attention to the background and suppress the effects of the foreground. In a nutshell, if a cropped person image can be decomposed into the foreground region and the background region, we can effectively learn the person ID as well as the camera ID separately.

However, in person re-identification task, the foreground mask of a cropped person image is usually unavailable. Since the foreground exactly corresponds to the supplementary region of the background in a cropped person image, the learning of person ID and camera ID can be decoupled by introducing a pseudo mask to indicate the foreground. Based on such observation, we propose a framework with two branches to mutually promote the learning of person ID and camera ID simultaneously. As illustrated in Fig. 3, given an input image \mathbf{X} , we first extract low-level feature maps, which are then fed to two independent branches, *i.e.*, the foreground branch and the background branch. The person representation and the background (*i.e.*, camera) representation learned from the two branches are exploited to predict the person ID and camera ID, respectively. To facilitate the learning, we propose a new target enhancement module as well as a target attention loss, which make two branches interact and promote each other and will be elaborated in the next subsection.

We adopt the ResNet50 [32] as the backbone model, while the global average pooling layer and the fully connected (FC) layer are removed. The layers before the res_conv4 block are adopted as the low-level feature extraction module. The rest blocks of ResNet50 are copied into two independent branches,

i.e., the foreground branch and the background branch. The two branches do not share their weights. Given the low-level features of image \mathbf{X} , the foreground branch first obtains the pedestrian feature map $\mathbf{F} \in \mathbb{R}^{C \times H \times W}$ from the output of the last residual block. Similarly, the background branch extracts the raw background feature map $\mathbf{B} \in \mathbb{R}^{C \times H \times W}$ from the output of the low-level feature extraction module. Based on \mathbf{F} , the foreground target enhancement module (TEM) generates the corresponding spatial attention map. After being enhanced by the spatial attention maps, we obtain the gated foreground feature map $\mathbf{F}^g \in \mathbb{R}^{C \times H \times W}$ and gated background feature map $\mathbf{B}^g \in \mathbb{R}^{C \times H \times W}$.

It has been proven that horizontal pyramid pooling (HPP) [33] successfully enhances the discriminative capabilities of various person parts. Since our network needs to accurately distinguish the foregrounds from backgrounds, the involvement of HPP can further improve the accuracy of spatial attention map prediction in local regions. Therefore, we apply HPP on both of the \mathbf{F}^g and \mathbf{B}^g to obtain features with four horizontal pyramid scales. The four scales have 1, 2, 4 and 8 spatial stripes, respectively. For each scale, feature maps are sliced to the corresponding number of stripes. The features in each stripe are then pooled and embedded into a 256-dim feature vector. Given a foreground feature vector, the corresponding person classifier predicts the person identity and calculates the softmax loss. Each background feature vector is fed to a corresponding camera classifier to predict the camera identity and calculate the softmax loss.

Although the foreground branch and the background branch share a similar architecture, they are trained with different objectives. The foreground branch predicts the person identity as the target by focusing on the foreground human body regions, while the background branch predicts the camera identity as the target by focusing on the background regions. Ideally, there should be no overlap between the focused regions of the two branches. In the following, we will introduce how to make

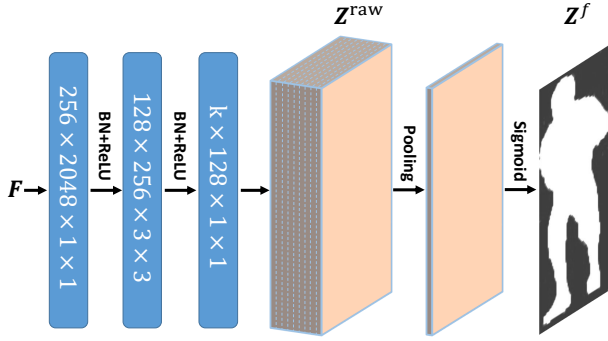


Fig. 4. The architecture of the target enhancement module TEM.

the two branches share their complementary knowledge with a target enhancement module to benefit each other in model training.

B. Target Enhancement Module

The target enhancement module (TEM) aims to generate a pseudo mask to indicate the target (*i.e.*, foreground) region and restrain the responses in the nontarget (*i.e.*, background) region, which is the key to the collaborative learning of the two branches. The architecture of our TEM is shown in Fig. 4. We first feed the raw feature map $\mathbf{F} \in \mathbb{R}^{C \times H \times W}$ into two convolutional blocks. Each consists of three consecutive operations: a convolutional layer, a batch normalization (BN) layer and a rectified linear unit (ReLU). The first convolutional block has 256 filters. The kernel size is set to 1×1 , so as to reduce the feature dimension. For the second convolutional block with 128 filters, the kernel size is set to 3×3 , which increases the receptive field of this module. Then, the output of the two blocks are fed into another convolutional layer with 1×1 kernel size to generate the spatial attention map $\mathbf{Z}^{\text{raw}} \in \mathbb{R}^{k \times H \times W}$. In \mathbf{Z}^{raw} , there are k channels and each channel corresponds to a spatial attention map. We average over these k spatial attention maps into one. Finally, this spatial attention map is normalized into $[0, 1]$ by the sigmoid function, which is formulated as follows,

$$\mathbf{Z}^f = \text{sigmoid}\left(\frac{1}{k} \sum_{c=1}^k \mathbf{Z}^{\text{raw}}(c)\right), \quad (1)$$

where $\mathbf{Z}^{\text{raw}}(c)$ denotes the c^{th} channel of \mathbf{Z}^{raw} . $\mathbf{Z}^f \in \mathbb{R}^{1 \times H \times W}$ is the final foreground spatial attention map, which works as our soft foreground mask.

In some previous works [4], [8], [34], [9], the spatial attention map is directly generated without channel-wise pooling, which suffers an unreliability issue and limits the accuracy of the attention map. In contrast, in the proposed TEM, we average over k spatial attention maps into the final attention map, which is more robust and accurate as justified later in our experiments.

The value of each location in \mathbf{Z}^f denotes the probability that the corresponding spatial location of \mathbf{F} belongs to the foreground target. The higher probability value indicates that the TEM considers the features in this location are more likely

to belong to the body part and should be reserved, while the features in the location with lower probability are more likely to belong to the backgrounds and should be restrained. Therefore, we make use of $1 - \mathbf{Z}^f$ to denote the probability that the corresponding spatial location of \mathbf{F} belongs to the background target. In other words, $\mathbf{Z}^b = 1 - \mathbf{Z}^f$ can be regarded as a soft background mask.

For the raw foreground feature map \mathbf{F} , since we have obtained the soft foreground mask \mathbf{Z}^f , we can use it to enhance the foreground features. The soft foreground mask is applied to each channel of the raw foreground feature map \mathbf{F} , formulated as follows,

$$\mathbf{F}^g = \mathbf{F} \odot \mathbf{Z}^f, \quad (2)$$

where \odot denotes the element-wise multiplication with broadcasting along the channels of \mathbf{F} . The gated person feature map \mathbf{F}^g is fed to the following layers to generate the final person feature representation. Similarly, the raw background feature map \mathbf{B} is gated by the soft background mask \mathbf{Z}^b , which is formulated as follows,

$$\mathbf{B}^g = \mathbf{B} \odot \mathbf{Z}^b = \mathbf{B} \odot (1 - \mathbf{Z}^f), \quad (3)$$

where \mathbf{B}^g is the gated background features. The features in the background regions are enhanced and the features in the foreground regions are restrained.

Under the definition of Eq. 2 and Eq. 3, the soft mask generated by TEM affects both foreground features and background features. This forces TEM to more accurately distinguish the foreground from the background to help both branches focus on the target areas. More importantly, with the help of TEM, the two branches collaboratively promote each other.

C. Target Attention Loss

As discussed in Section III-B, the spatial attention map \mathbf{Z}^f is considered as the soft foreground mask, while the spatial attention map $1 - \mathbf{Z}^f$ is considered as the soft background mask. In principle, in the foreground mask, the values corresponding to the background regions are expected to be close to zero, while the values of the foreground regions are expected to be close to 1. Besides, for the raw foreground features \mathbf{F} , the responses on the nontarget background regions should be small. Similarly, for the raw background features \mathbf{B} , the responses on the nontarget body regions are expected to be small. With such intuition, we design a target attention loss (TAL) as follows,

$$\mathcal{L}_t = \text{avg} \left[\mathbf{F}^{\ell^2} \odot (1 - \mathbf{Z}^f) + \mathbf{B}^{\ell^2} \odot \mathbf{Z}^f \right], \quad (4)$$

where $\text{avg}[\cdot]$ denotes the average operation. \mathbf{F}^{ℓ^2} and \mathbf{B}^{ℓ^2} are the result of performing ℓ^2 normalization over the spatial dimension of \mathbf{F} and \mathbf{B} , which are formulated as

$$\mathbf{F}^{\ell^2}(c) = \frac{\mathbf{F}(c)}{\|\mathbf{F}(c)\|_2}, \mathbf{B}^{\ell^2}(c) = \frac{\mathbf{B}(c)}{\|\mathbf{B}(c)\|_2}, \quad (5)$$

where $\mathbf{F}(c)$ and $\mathbf{F}^{\ell^2}(c)$ correspond to the feature map of the c^{th} channel of \mathbf{F} and \mathbf{F}^{ℓ^2} , respectively. The ℓ^2 normalization applied on the raw feature maps is introduced to avoid the loss simply forcing all values of features to approach zero.

Since $\mathbf{1} - \mathbf{Z}^f$ is the background target mask predicted by TEM, $\mathbf{F}^{\ell^2} \odot (\mathbf{1} - \mathbf{Z}^f)$ denotes the response in the predicted nontarget regions of the foreground features. Then, the minimization of it forces the foreground branch to focus more on the person body and the attention map \mathbf{Z}^f is required to be more accurate. Similarly, the minimization of $\mathbf{B}^{\ell^2} \odot \mathbf{Z}^f$ requires the background branch focus more on the background regions and learn better background feature \mathbf{B} under the guidance of the soft mask achieved from foreground branch. Meanwhile, this also requires TEM to distinguish well between the foreground and the background. Therefore, the minimization of \mathcal{L}_t lets the two branches promote each other, which makes better use of the opposite relationship between the two branches.

D. The Overall Training Objective

With the proposed target attention loss \mathcal{L}_t , the overall training objective of our approach is formulated as follows,

$$\mathcal{L} = \frac{1}{2} (\mathcal{L}_f + \mathcal{L}_b) + \mathcal{L}_t \quad (6)$$

where \mathcal{L}_f and \mathcal{L}_b denote the softmax losses of the foreground branch and background branch for person-ID classification and camera-ID classification, respectively. By minimizing \mathcal{L} , the proposed approach learns the foreground feature representations and the background feature representations simultaneously. Unlike existing works [1], [2], [3], [5], [7], [6], the prediction of the background and the training of person re-identification model are not separate. The addition of the target enhancement module and target attention loss makes the two branches couple and promote each other, which allows our model to obtain a more accurate separation of the foreground and background. Specifically, in the forward propagation, the foreground feature extraction does not depend on the background branch. It is notable that, during the inference stage for person re-ID, the background branch is no longer needed.

IV. EXPERIMENT

In this section, we evaluate the proposed method on three large public image-based person re-identification datasets. We first describe the datasets and implementation details in Section IV-A and Section IV-B, respectively. Then we make an ablation study of our method in Section IV-C. After that, we provide the further analysis and discussion about FA-Net in Section IV-D. Finally, in Section IV-E, we compare our method with the state-of-the-art methods.

A. Datasets and Protocols

To evaluate our proposed methods, we select three large publicly available person re-identification datasets namely Market-1501 [29], DukeMTMC-reID [30] and MSMT17 [31]. Market-1501 contains 32,668 images of 1,501 identities captured by 5 high-resolution cameras and one low-resolution camera. Images are detected by deformable part model (DPM) [35]. The dataset is split into the training set and testing set. 12,936 images of 751 identities are selected as the training set.

The rest 750 identities are used to create the gallery and query sets, which contain 19,734 and 3,368 images, respectively.

DukeMTMC-reID [30] contains the person images extracted from the DukeMTMC [28] tracking dataset. These hand-annotated images are captured from 8 high-resolution cameras. In the standard evaluation protocol, the training set consists of 16,522 images of 702 identities. The remaining 702 identities are used as the testing set with 2,228 query images and 17,661 gallery images. This dataset is very challenging due to the large variations within the same identity and high similarity across persons.

MSMT17 [31] is a newly released large-scale person re-identification dataset, which consists of 126,441 images of 4,101 identities. The images are captured by 12 outdoor cameras and 3 indoor cameras. 4 days with different weather conditions in a month are selected for video collection. Meanwhile, the videos of 3 hours in each day are taken in the morning, noon and afternoon, respectively. The bounding boxes are detected by Faster RCNN [36]. In the standard evaluation protocol, 30,248 images of 1,041 identities are sampled as the training set. The rest images of the 1,041 identities are used as the validation set. The 3060 identities that do not appear in the training set are selected as the testing set with 11,659 query images and 82,161 gallery images.

Compared with Market-1501 and DukeMTMC-reID, MSMT17 contains more identities and images. The more camera views, both indoor and outdoor scenes and the lighting changes at different times of one day make the backgrounds more complex and challenging than previous datasets.

Following most of the previous works, we adopt the Cumulated Matching Characteristics (CMC) table and the mean Average Precision (mAP) to evaluate the performance of each method. All experiments are conducted with the single query setting.

B. Implementation Details

The backbone model ResNet50 is pre-trained on the ImageNet dataset. In order to increase the spatial resolution, following [21], [33], the last spatial down-sampling operation in the backbone network is removed. The input images of the proposed model are resized to 384×128 . Random horizontal flipping is adopted for data augmentation. In each iteration, we select images of 16 pedestrians each with 8 images as the inputs of the network in a mini-batch. The images of each pedestrian are taken from different cameras as much as possible.

The network is updated for 100 epochs by the stochastic gradient descent algorithm with a weight decay of 5×10^{-4} . Following [37], the warmup learning rate adjustment strategy is applied to bootstrap the network for better performance. The learning rate linearly increases from 0.06 to 0.6 in the first 10 epochs. Then, the learning rate is decayed to 6×10^{-2} and 6×10^{-3} at 40th and 80th epoch respectively. The learning rate of the pre-trained layers is set to $0.1 \times$ of the base learning rate. During the evaluation, the averaged feature of the original image and the horizontally flipped version is extracted for each pedestrian image. We use the cosine distance to measure the similarity of two images.

TABLE I

THE ABLATION STUDY OF THE PROPOSED METHOD ON THE MARKET-1501, DUKEMTMC-reID, AND MSMT17 DATASETS. THE CMC RESULTS AND mAP ACCURACY ARE REPORTED. FOR THE BASELINE NETWORK, A GLOBAL AVERAGE POOLING IS DIRECTLY APPLIED TO THE OUTPUTS OF THE MODIFIED RESNET50 BACKBONE MODEL TO GENERATE THE FINAL FEATURE REPRESENTATIONS. *B/L* DENOTES THE BASELINE NETWORK. *TEM* IS THE TARGET ENHANCEMENT MODULE. *BG* MEANS THAT THE BACKGROUND BRANCH IS ADDED, WHILE THE SOFT MASK GENERATED BY *TEM* IS NOT APPLIED TO BACKGROUND FEATURES. *IA* DENOTES THAT THE IMPORTANT INTERACTION BETWEEN THE TWO BRANCHES IS ADOPTED, *i.e.*, THE BACKGROUND FEATURES ARE GATED BY THE SOFT MASK GENERATED BY THE FOREGROUND BRANCH. *TAL* MEANS THAT THE TARGET ATTENTION LOSS IS ADOPTED. *FA-Net* IS THE FINAL ARCHITECTURE OF OUR METHOD, WHERE THE HORIZONTAL PYRAMID POOLING HPP [33] IS APPLIED.

Method	Market-1501				DukeMTMC-reID				MSMT17			
	R1	R5	R10	mAP	R1	R5	R10	mAP	R1	R5	R10	mAP
Baseline	89.7	96.4	97.8	72.7	77.9	88.9	91.6	60.1	64.4	77.4	81.9	31.4
B/L+TEM	90.9	96.5	97.7	73.5	79.5	88.7	91.7	60.7	65.4	78.3	82.8	33.0
B/L+TEM+BG	92.5	97.1	97.9	79.3	83.4	91.8	93.5	67.6	70.8	82.9	86.8	41.1
B/L+TEM+BG+IA	92.9	97.2	98.2	79.6	84.3	92.0	94.6	67.7	71.6	83.4	86.8	41.2
B/L+TEM+BG+IA+TAL	93.3	97.4	98.3	80.1	85.2	91.9	94.0	67.9	72.3	83.5	87.2	42.1
FA-Net (B/L+TEM+BG+IA+TAL+HPP)	95.0	97.9	98.6	84.6	88.7	93.8	95.5	77.0	76.8	86.8	89.8	51.0

TABLE II

THE IMPACT OF TARGET ATTENTION LOSS TAL WITH DIFFERENT SETTING ON THE PERFORMANCE OF THE PROPOSED METHOD. THE DEFINITIONS OF TAL^{v1} AND TAL^{v2} IS GIVEN IN EQ. 7. IN THE LAST EXPERIMENTS, TAL IS DEFINED AS EQ. 4.

Method	Market-1501					DukeMTMC-reID					MSMT17				
	R1	R5	R10	R20	mAP	R1	R5	R10	R20	mAP	R1	R5	R10	R20	mAP
B/L+TEM+BG+IA	92.9	97.2	98.2	99.1	79.6	84.3	92.0	94.6	96.1	67.7	71.6	83.4	86.8	89.8	41.2
B/L+TEM+BG+IA+TAL ^{v1}	92.6	97.3	98.5	99.2	79.7	82.8	91.8	93.5	95.4	68.5	68.9	81.1	85.1	88.5	39.1
B/L+TEM+BG+IA+TAL ^{v2}	93.1	97.7	98.6	99.2	79.8	83.3	91.5	93.7	95.2	67.1	70.0	81.9	85.7	89.0	39.4
B/L+TEM+BG+IA+TAL	93.3	97.4	98.3	99.1	80.1	85.2	91.9	94.0	95.6	67.9	72.3	83.5	87.2	90.1	42.1

C. Ablation Study

Impact of each component. As shown in Table I, we evaluate the effect of each component of our network. The baseline network directly applies the global average pooling on the outputs of the modified ResNet50 backbone model to generate the final feature representations. After TEM is added, the network achieves 1.2%, 1.6% and 1.0% improvement in the rank-1 accuracy and 0.8%, 0.6% and 1.6% improvement in the mAP accuracy on Market-1501, DukeMTMC-reID and MSMT17, respectively. This indicates that TEM effectively helps the network focus more on discriminative regions.

In *B/L+TEM+BG*, the background branch is added, but the background features are not gated by soft mask. The joint training of the two branches brings significant improvements in the rank-1 accuracies and mAP accuracies on all these three datasets. Because of the addition of the background branch, the low-level feature extraction module is shared by the two branches. To predict the camera identities, the background branch requires the low-level feature extraction module to learn additional texture and color patterns, which provides richer patterns for the extraction of foreground features.

When the main interaction between the two branches is applied, the gated features are obtained according to Eq. 2 and Eq. 3. The network achieves 0.4%, 0.9% and 0.8% improvement in the rank-1 accuracy on Market-1501, DukeMTMC-reID and MSMT17, respectively. This is because the prediction of the soft mask generated by TEM simultaneously affects the features of both branches. To identify the camera identities of images, the network requires the soft mask accurately distinguish between foreground and background. The addition of new supervision information better guides the training of

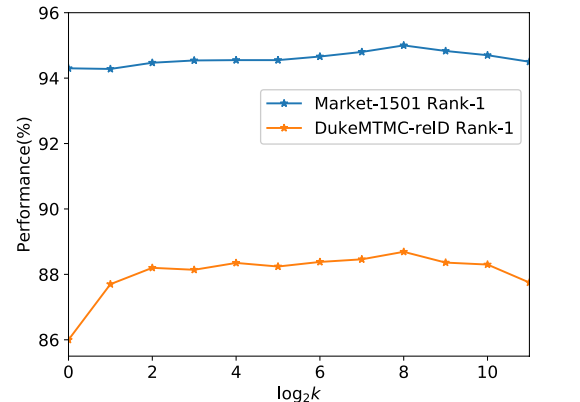


Fig. 5. The rank-1 accuracies of FA-Net with different k values on the Market-1501 and DukeMTMC-reID datasets.

TEM.

After adding TAL, the performances are improved by 0.4%, 0.9% and 0.7% in rank-1 accuracy and 0.5%, 0.2% and 0.9% in mAP accuracy on Market-1501, DukeMTMC-reID and MSMT17, respectively. This shows that TAL helps the two branches interact better, which makes each branch focus more on its target regions and makes TEMs learn more accurate attention maps. When HPP is adopted, another performance gain is obtained. This is because the addition of HPP makes better use of TEM. HPP forces the network to pay more attention to the local regions and helps TEM improve the accuracy of the prediction.

Analysis on the target enhance module. In TEM, we adopt the averaged results of k spatial attention maps as the final

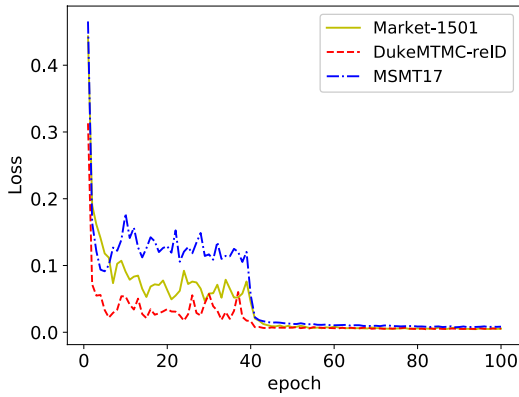


Fig. 6. Visualization of changes in the loss values of camera identity prediction on three datasets.

attention map. Fig. 5 shows the performances achieved by our method with different k . We find that FA-Net achieves the best performance when $k = 256$. When $k = 1$, the rank-1 accuracies drop on both of the two datasets compared to $k = 256$. This is because directly generating the final attention map is unreliable. The spatial noise could corrupt the accuracy of the attention map. However, when the final attention map is the averaged result of several attention maps, it is more robust to noise, which is important to accurately enhance the target regions.

Analysis on the target attention loss. In Table II, we show the impact of TAL with different settings on the performance of the proposed method. The other two versions of TAL are analyzed, which is formulated as

$$\begin{aligned} \mathcal{L}_t^{v1} &= \text{avg} [\mathbf{F} + \mathbf{B}], \\ \mathcal{L}_t^{v2} &= \text{avg} [\mathbf{F} \odot (\mathbf{1} - \mathbf{Z}^f) + \mathbf{B} \odot \mathbf{Z}^f], \end{aligned} \quad (7)$$

where \mathcal{L}_t^{v1} and \mathcal{L}_t^{v2} correspond to TAL^{v1} and TAL^{v2} , respectively. In TAL^{v1} , the loss just regularizes the foreground and background features, which caused a slight degradation in the rank-1 accuracy of the model. This denotes that simply regularizing the features can not boost the performance. Compared to TAL^{v1} , TAL^{v2} achieves slight performance improvements on the rank-1 accuracy, while still damages the performance of model on DukeMTMC-reID and MSMT17. For TAL defined in Eq. 4, the ℓ^2 normalization on \mathbf{F} and \mathbf{B} is applied, which boost both of the rank-1 accuracy and mAP accuracy of the model. This is because ℓ^2 normalization avoids the loss simply minimizing the values of all locations. The suppression of the responses of the nontarget areas makes the model focus on the target regions and learn better soft mask.

D. Further Analysis and Discussion

Is it feasible to use the camera identity information to guide the network to learn background features? There are two problems with the training data that may challenge the proposed approach. 1) For the same scene captured by the same camera, the pedestrian images are cropped from different locations, whose backgrounds may be very different. In fact,

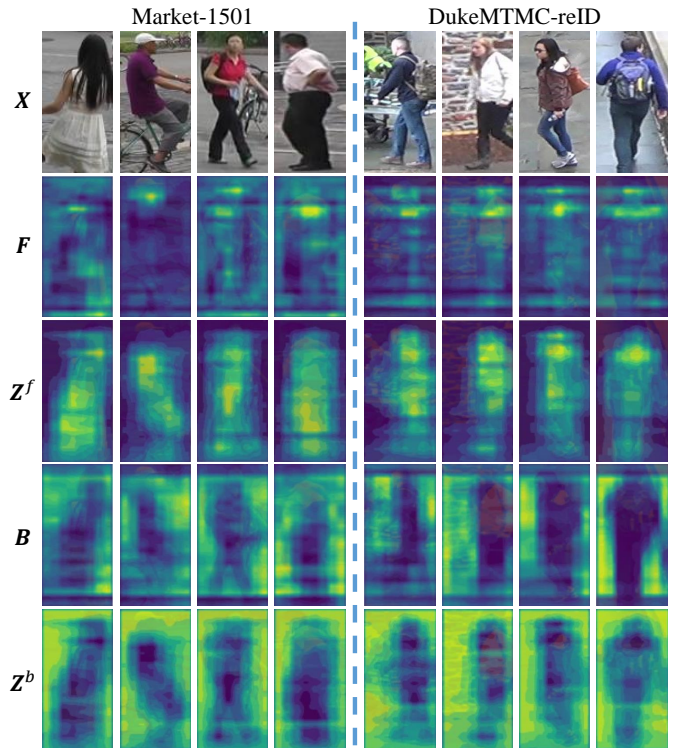


Fig. 7. Visualization of the features and attention maps of FA-Net on the testing set. The first row is the image input \mathbf{X} . The following rows are the corresponding foreground features \mathbf{F} , foreground attention map \mathbf{Z}^f , background features \mathbf{B} and background attention map \mathbf{Z}^b of each image. The features and attention maps are displayed above the original images. The strip-shaped responses are due to the addition of horizontal pyramid pooling (HPP).

this situation does not bother the network. The camera classifiers predict camera identities of images based on whether the backgrounds of images belong to the corresponding scene, rather than focusing on a certain background. 2) Some local scene regions captured by different cameras may be similar and cannot be distinguished by humans. Nevertheless, the difference in the viewpoints and local fine textures can help the network distinguish between them. For example, some paths of different scenes in Fig. 1 are similar, but the angles of the bricks are different because of the variations in viewpoints. Notably, there may be some local regions of different scenes in the training data that are too similar to the model to distinguish. However, this case only occupies a small part of the data, which is considered as noise data. Otherwise, this person re-identification problem with similar backgrounds is so easy that a simple network can solve it. Therefore, it is feasible to use the camera identity information to help network learn background representations.

To verify our motivation, we show the changes in the loss value of camera identity prediction as the training epoch increases in Fig. 6. The final loss values tend to approximate zero, which indicates that the camera classifiers well predict the camera identities of training images. We show some examples of the features and attention maps of testing images generated by FA-Net in Fig. 7. The responses of the background features are mainly in the background regions.

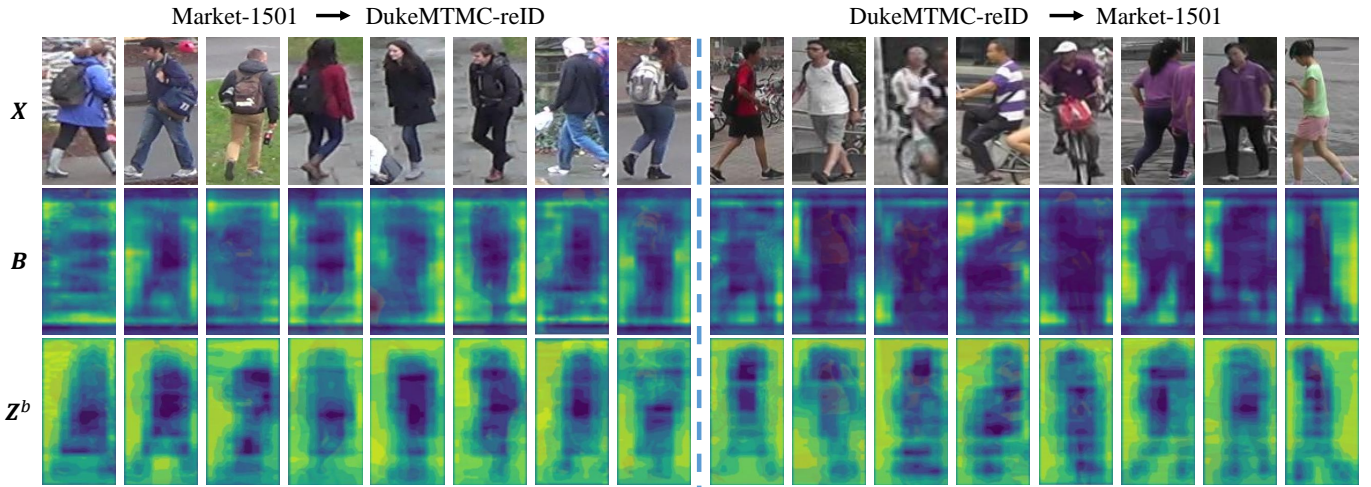


Fig. 8. Visualization of the background features and background attention maps on the unseen scenes. The first row is the image input \mathbf{X} . The following rows are the corresponding background features \mathbf{B} and background attention map \mathbf{Z}^b of each image.

TABLE III

THE PERFORMANCES OF OUR METHOD ON UNSEEN SCENES. THE NETWORK IS TRAINED ON ONE DATASET AND DIRECTLY TESTED ON ANOTHER DATASET. THEREFORE, THERE IS NO SCENE OVERLAP BETWEEN THE TRAINING SET AND THE TESTING SET.

Method	Market-1501 \rightarrow DukeMTMC-reID					DukeMTMC-reID \rightarrow Market-1501				
	R1	R5	R10	R20	mAP	R1	R5	R10	R20	mAP
Baseline	28.5	43.5	50.3	56.8	13.9	50.4	67.4	74.2	80.8	21.5
B/L+TEM	31.8	46.5	53.1	59.5	16.1	52.6	69.7	76.0	81.9	23.2
B/L+TEM+BG	36.9	52.0	58.3	65.1	20.3	54.3	71.0	77.2	82.7	24.7
B/L+TEM+BG+IA	38.0	53.7	60.2	66.1	20.7	55.1	72.0	78.3	83.8	25.3
B/L+TEM+BG+IA+TAL	39.4	54.9	60.6	66.0	21.5	56.2	73.0	78.5	83.0	25.5
FA-Net (B/L+TEM+BG+IA+TAL+HPP)	49.3	63.5	69.0	73.7	30.7	65.1	79.3	84.4	89.0	34.2

The convergence of the background branch on the training set and the responses in the background regions of the testing images verify our motivation that the camera identity information can guide the background branch to learn background representations. By observing the soft masks of the foreground and background, we see that TEM well distinguishes foregrounds and the backgrounds. This is benefited from the guidance from both branches and the addition of target attention loss. With the soft masks, TEM forces the two branches to focus on the target regions and learn better foreground and background representations.

Is performance improvement due to more parameters?

In the training stage, compared with the baseline method and $B/L+TEM$, $B/L+TEM+BG+IA+TAL$ has more parameters due to the addition of the background branch. However, it has a similar number of parameters to the baseline network and the same number of parameters to $B/L+TEM$ in inference because TEM has very few parameters and the background branch is not used. The results in Table I show that the performance of $B/L+TEM+BG+IA+TAL$ is improved significantly over baseline model. Rank-1 accuracies are improved by **3.6%**, **7.3%** and **7.9%** and mAP accuracies are improved by **7.4%**, **7.8%** and **10.7%** on Market-1501, DukeMTMC-reID and MSMT17, respectively. Compared to $B/L+TEM$, $B/L+TEM+BG+IA+TAL$ boosts rank-1 accuracies by **2.4%**, **5.7%** and **6.9%** and mAP accuracies by **6.6%**, **7.2%** and

9.1% on Market-1501, DukeMTMC-reID and MSMT17, respectively. This means that the significant performance improvements achieved by our method are mainly due to the better feature extraction rather than more parameters. Specifically, the significant performance gains under the same number of parameters indicate that our approach is more efficient and helps reduce the computational overhead of large-scale person re-identification.

Evaluation on unseen scenes. In our method, we use the camera identities to guide the background branch to learn the background features and constrain the learning of the soft spatial mask to help the foreground feature focus on the person body parts. The involvement of the background branch brings one question: whether the proposed method can still improve performance under unknown backgrounds?

To verify the effectiveness of the proposed method on the unseen scenes, we train the network on one dataset and directly test it on another dataset. The collections of these two datasets are in different scenes and use different camera settings. So our model is tested on the new scenes, which have different backgrounds from the training set. As shown in Table III, we observe that even in the unseen scenes, every module of our method improves the performance. This shows that each module of FA-Net is still effective even in a new scene. The background branch is introduced to help the low-level feature extraction module learn richer patterns and regularize the

learning of TEM, and it is abandoned in inference. Therefore, FA-Net still performs well in the unseen scene.

We show some examples of the background features and background soft masks generated by our method on the unseen scenes in Fig. 8. It's observed that even in unseen scenes, the background branch still pays more attention to the backgrounds and the TEM well distinguishes the foregrounds and backgrounds. Specifically, for the first testing image when training on Market-1501, this kind of wall does not appear in the training dataset, but the responses in the features generated by background branch appear at the wall. This indicates that the background model trained using camera identity information is generalized to unseen scenes.

Does using camera identity information require additional data collection overhead? In an intelligent surveillance system, after we retrieve the image of the person of interest, we usually need to further know the person's location. This is available according to the location of the camera that captures this image. It indicates that in practical applications, it is necessary to record which camera each image comes from. Meanwhile, the recording of camera identity information is very easy and does not require manual labeling. Most of the existing person re-identification datasets also record the camera identity of each image. Instead, the methods based on the human landmark detection model and segmentation model require additional manually labeled datasets. This shows that using camera identity information is more economical and does not incur the overhead of additional data collection for many practical applications.

E. Comparison with the State-of-the-Art Methods

As shown in Table IV, we first compare our method with the related works on the Market-1501 and DukeMTMC-reID datasets. Some approaches that try to remove the influence from the backgrounds are included, such as human landmark detection method GLAD [2], segmentation method SPReID [5] and attention-based method HA-CNN [8]. Our approach achieves 95.0% rank-1 accuracy and 84.6% mAP accuracy on the Market-1501 dataset, 88.7% rank-1 accuracy and 77.0% mAP accuracy on the DukeMTMC-reID dataset.

Compared with the segmentation method SPReID [5], our method boosts the rank-1 accuracy by 2.5% and 4.3% and mAP accuracy by 3.3% and 6.0% on Market-1501 and DukeMTMC-reID, respectively. This indicates that our method mitigates the impact of the backgrounds and achieves better performance even without the additional human pose or segmentation datasets. Compared to attention-based method HA-CNN [8], our method improves the rank-1 accuracies by 3.8% and 8.2% and mAP accuracies by 8.9% and 13.2% on Market-1501 and DukeMTMC-reID, respectively. This shows that the additional supervision information (camera identify information and TAL) effectively helps the attention module TEM to predict the target regions. In IANet [38], a spatial interaction-and-aggregation module (SIA) is proposed to deal with large variations in body pose and scale, which makes the network learn more robust foreground features. Compared to IANet [38], our method improves the rank-1 accuracies by

TABLE IV
COMPARISON WITH THE RELATED METHODS ON MARKET-1501 AND DUKEMTMC-REID. THE MAP AND RANK-1 ACCURACIES ARE REPORTED. RK DENOTES THE RE-RANKING OPERATION [39].

Method	Reference	Market-1501		DukeMTMC-reID	
		R1	mAP	R1	mAP
CAN [40]	TIP'17	60.3	35.9	-	-
GLAD [2]	MM'17	89.9	73.9	-	-
AACN [41]	CVPR'18	85.9	66.9	76.8	59.3
HA-CNN [8]	CVPR'18	91.2	75.7	80.5	63.8
SPReID[5]	CVPR'18	92.5	81.3	84.4	71.0
FD-GAN [42]	NIPS'18	90.5	77.7	80.0	64.5
PABR [20]	ECCV'18	91.7	79.6	84.4	69.3
PCB [21]	ECCV'18	92.3	77.4	81.7	66.1
PCB+RPP [21]	ECCV'18	93.8	81.6	83.3	69.2
LITM+GHIS [43]	AAAI'19	93.9	83.9	85.9	74.5
HPM [33]	AAAI'19	94.2	82.7	86.6	74.3
IANet [38]	CVPR'19	94.4	83.1	87.1	73.4
FA-Net	This work	95.0	84.6	88.7	77.0
AACN+RK [41]	CVPR'18	88.7	83.0	-	-
PABR+RK [20]	ECCV'18	93.4	89.9	88.3	83.9
PCB+RPP+RK [21]	ECCV'18	95.1	91.9	-	-
FA-Net+RK	This work	95.8	93.4	91.5	88.9

TABLE V
COMPARISON WITH THE EXISTING METHODS ON MSMT17.

Method	Reference	R1	R5	R10	mAP
GoogLeNet [44]	CVPR'15	47.6	65.0	71.8	23.0
PDC [3]	ICCV'17	58.0	73.6	79.4	29.7
GLAD [2]	MM'17	61.4	76.8	81.6	34.0
PCB + RPP [21]	ECCV'18	68.2	81.2	85.5	40.4
IANet [38]	CVPR'19	75.5	85.5	88.7	46.8
FA-Net	This work	76.8	86.8	89.8	51.0

0.6% and 1.6% and mAP accuracies by 1.5% and 3.6% on Market-1501 and DukeMTMC-reID, respectively. This indicates that it is effective to use the complementary knowledge about foreground and background to learn foreground mask for enhancing foreground features.

In Table V, we compare our method with the existing methods on MSMT17 dataset. Compared with IANet [38], our method boosts the rank-1 accuracy by 1.3% and the mAP accuracy by 4.2%. It is worth noting that the images of MSMT17 have more complex backgrounds due to the 15 camera views with both indoor and outdoor scenes and the lighting changes at different times of one day. The significant performance improvement achieved on such a challenging dataset demonstrates the effectiveness of our method in handling the effects from the backgrounds and extracting more robust and discriminative pedestrian features.

V. CONCLUSION

In this paper, we have proposed an end-to-end foreground-aware network for person re-identification. In order to alleviate the influence from the backgrounds, our method learns a soft foreground mask and locates the background regions using the camera identities available in the existing person re-identification datasets, rather than from additional human pose or segmentation datasets. Benefiting from the target enhancement modules and the target attention loss, the foreground

branch and the background branch simultaneously promote each other and learn more robust and discriminative feature representations. Extensive experiments on three large person re-identification datasets demonstrate the effectiveness of our approach.

REFERENCES

- [1] H. Zhao, M. Tian, S. Sun, J. Shao, J. Yan, S. Yi, X. Wang, and X. Tang, "Spindle net: Person re-identification with human body region guided feature decomposition and fusion," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 1077–1085.
- [2] L. Wei, S. Zhang, H. Yao, W. Gao, and Q. Tian, "Glad: Global-local-alignment descriptor for pedestrian retrieval," in *Proceedings of the ACM International Conference on Multimedia (ACM MM)*. ACM, 2017, pp. 420–428.
- [3] C. Su, J. Li, S. Zhang, J. Xing, W. Gao, and Q. Tian, "Pose-driven deep convolutional model for person re-identification," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 3960–3969.
- [4] L. Zhao, X. Li, Y. Zhuang, and J. Wang, "Deeply-learned part-aligned representations for person re-identification," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 3219–3228.
- [5] M. M. Kalayeh, E. Basaran, M. Gökmen, M. E. Kamasak, and M. Shah, "Human semantic parsing for person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 1062–1071.
- [6] M. Tian, S. Yi, H. Li, S. Li, X. Zhang, J. Shi, J. Yan, and X. Wang, "Eliminating background-bias for robust person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 5794–5803.
- [7] C. Song, Y. Huang, W. Ouyang, and L. Wang, "Mask-guided contrastive attention model for person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 1179–1188.
- [8] W. Li, X. Zhu, and S. Gong, "Harmonious attention network for person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 2285–2294.
- [9] C. Wang, Q. Zhang, C. Huang, W. Liu, and X. Wang, "Manacs: A multi-task attentional network with curriculum sampling for person re-identification," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- [10] Y.-J. Cho, S.-A. Kim, J.-H. Park, K. Lee, and K.-J. Yoon, "Joint person re-identification and camera network topology inference in multiple cameras," *Proceedings of the Computer Vision and Image Understanding (CVIU)*, 2019.
- [11] W. Huang, R. Hu, C. Liang, Y. Yu, Z. Wang, X. Zhong, and C. Zhang, "Camera network based person re-identification by leveraging spatial-temporal constraint and multiple cameras relations," in *Proceedings of the International Conference on Multimedia Modeling (MMM)*. Springer, 2016, pp. 174–186.
- [12] J. Lv, W. Chen, Q. Li, and C. Yang, "Unsupervised cross-dataset person re-identification by transfer learning of spatial-temporal patterns," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [13] G. Wang, J. Lai, P. Huang, and X. Xie, "Spatial-temporal person re-identification," in *AAAI Conference on Artificial Intelligence (AAAI)*, 2019.
- [14] S. Karanam, Y. Li, and R. J. Radke, "Person re-identification with discriminatively trained viewpoint invariant dictionaries," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 4516–4524.
- [15] F. Xiong, M. Gou, O. Camps, and M. Szaier, "Person re-identification using kernel-based metric learning methods," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2014, pp. 1–16.
- [16] R. Zhao, W. Ouyang, and X. Wang, "Learning mid-level filters for person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014, pp. 144–151.
- [17] S. Liao, Y. Hu, X. Zhu, and S. Z. Li, "Person re-identification by local maximal occurrence representation and metric learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 2197–2206.
- [18] X. Yang, M. Wang, and D. Tao, "Person re-identification with metric learning using privileged information," *IEEE Transactions on Image Processing (TIP)*, vol. 27, no. 2, pp. 791–805, 2017.
- [19] S.-Z. Chen, C.-C. Guo, and J.-H. Lai, "Deep ranking for person re-identification via joint representation learning," *IEEE Transactions on Image Processing (TIP)*, vol. 25, no. 5, pp. 2353–2367, 2016.
- [20] Y. Suh, J. Wang, S. Tang, T. Mei, and K. Mu Lee, "Part-aligned bilinear representations for person re-identification," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 402–419.
- [21] Y. Sun, L. Zheng, Y. Yang, Q. Tian, and S. Wang, "Beyond part models: Person retrieval with refined part pooling (and a strong convolutional baseline)," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 480–496.
- [22] F. Zheng, C. Deng, X. Sun, X. Jiang, X. Guo, Z. Yu, F. Huang, and R. Ji, "Pyramidal person re-identification via multi-loss dynamic training," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 8514–8522.
- [23] Y. Sun, Q. Xu, Y. Li, C. Zhang, Y. Li, S. Wang, and J. Sun, "Perceive where to focus: Learning visibility-aware part-level features for partial person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 393–402.
- [24] Y. Shen, T. Xiao, H. Li, S. Yi, and X. Wang, "End-to-end deep kronecker-product matching for person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 6886–6895.
- [25] D. Chen, H. Li, T. Xiao, S. Yi, and X. Wang, "Video person re-identification with competitive snippet-similarity aggregation and co-attentive snippet embedding," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 1169–1178.
- [26] Y. Shen, H. Li, S. Yi, D. Chen, and X. Wang, "Person re-identification with deep similarity-guided graph neural network," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 486–504.
- [27] M. Andriluka, L. Pishchulin, P. Gehler, and B. Schiele, "2d human pose estimation: New benchmark and state of the art analysis," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014, pp. 3686–3693.
- [28] E. Ristani, F. Solera, R. Zou, R. Cucchiara, and C. Tomasi, "Performance measures and a data set for multi-target, multi-camera tracking," in *Proceedings of the European Conference on Computer Vision (ECCV)*. Springer, 2016, pp. 17–35.
- [29] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian, "Scalable person re-identification: A benchmark," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 1116–1124.
- [30] Z. Zheng, L. Zheng, and Y. Yang, "Unlabeled samples generated by gan improve the person re-identification baseline in vitro," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 3754–3762.
- [31] L. Wei, S. Zhang, W. Gao, and Q. Tian, "Person transfer gan to bridge domain gap for person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 79–88.
- [32] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [33] Y. Fu, Y. Wei, Y. Zhou, H. Shi, G. Huang, X. Wang, Z. Yao, and T. Huang, "Horizontal pyramid matching for person re-identification," in *AAAI Conference on Artificial Intelligence (AAAI)*, 2019.
- [34] Y. Liu, Z. Yuan, W. Zhou, and H. Li, "Spatial and temporal mutual promotion for video-based person re-identification," in *AAAI Conference on Artificial Intelligence (AAAI)*, 2019.
- [35] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, "Object detection with discriminatively trained part-based models," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 32, no. 9, pp. 1627–1645, 2010.
- [36] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2015, pp. 91–99.
- [37] H. Luo, Y. Gu, X. Liao, S. Lai, and W. Jiang, "Bag of tricks and a strong baseline for deep person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPR Workshops)*, 2019, pp. 0–0.
- [38] R. Hou, B. Ma, H. Chang, X. Gu, S. Shan, and X. Chen, "Interaction-and-aggregation network for person re-identification," in *Proceedings*

- of the *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 9317–9326.
- [39] Z. Zhong, L. Zheng, D. Cao, and S. Li, “Re-ranking person re-identification with k-reciprocal encoding,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 1318–1327.
- [40] H. Liu, J. Feng, M. Qi, J. Jiang, and S. Yan, “End-to-end comparative attention networks for person re-identification,” *IEEE Transactions on Image Processing (TIP)*, vol. 26, no. 7, pp. 3492–3506, 2017.
- [41] J. Xu, R. Zhao, F. Zhu, H. Wang, and W. Ouyang, “Attention-aware compositional network for person re-identification,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 2119–2128.
- [42] Y. Ge, Z. Li, H. Zhao, G. Yin, S. Yi, X. Wang *et al.*, “Fd-gan: Pose-guided feature distilling gan for robust person re-identification,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2018, pp. 1230–1241.
- [43] Y. Zhang, Q. Zhong, L. Ma, D. Xie, and S. Pu, “Learning incremental triplet margin for person re-identification,” in *AAAI Conference on Artificial Intelligence (AAAI)*, 2019.
- [44] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, “Going deeper with convolutions,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.