# Part-Guided Attention Learning for Vehicle Instance Retrieval

Xinyu Zhang*, Rufeng Zhang*, Jiewei Cao, Dong Gong, Minyu You#,  Chunhua Shen

*Abstract*—Vehicle instance retrieval (IR) often requires one to recognize the fine-grained visual differences between vehicles. Besides the holistic appearance of vehicles which is easily affected by the viewpoint variation and distortion, vehicle parts also provide crucial cues to differentiate near-identical vehicles. Motivated by these observations, we introduce a *Part-Guided Attention Network* (PGAN) to pinpoint the prominent part regions and effectively combine the global and local information for discriminative feature learning. PGAN first detects the locations of different part components and salient regions regardless of the vehicle identity, which serves as the *bottom-up attention* to narrow down the possible searching regions. To estimate the importance of detected parts, we propose a *Part Attention Module* (PAM) to adaptively locate the most discriminative regions with high-attention weights and suppress the distraction of irrelevant parts with relatively low weights. The PAM is guided by the identification loss and therefore provides *top-down attention* that enables attention to be calculated at the level of car parts and other salient regions. Finally, we aggregate the global appearance and local features together to improve the feature performance further. The PGAN combines part-guided bottom-up and top-down attention, global and local visual features in an end-to-end framework. Extensive experiments demonstrate that the proposed method achieves new state-of-the-art vehicle IR performance on four large-scale benchmark datasets.[1]

*Index Terms*—Vehicle instance retrieval, bottom-up attention, top-down attention.



**FIG. 1:** Illustration of the part-guided attention. (a) The rear and front views of two different vehicles with the same car model. (b) The detected candidate part regions from the part extraction module. (c) The heatmaps of part features from the part attention module. The prominent part regions like annual signs are highlighted, while the wrong candidates and insignificant parts like background and back mirror are suppressed.

## I. INTRODUCTION

VEHICLE instance retrieval (IR) aims to verify whether or not two vehicle images captured by different cameras belong to the same identity. Vehicle IR is also known as vehicle re-identification. With the growth of road traffic, it plays an increasingly important role in urban systems and intelligent transportation [1], [2], [3], [4], [5], [6], [7], [8], [9], [10].

Different levels of granularity of visual attention are required under various IR scenarios. In the case of comparing vehicles of different car models, we can easily distinguish their identities by examining the overall appearances, such as car
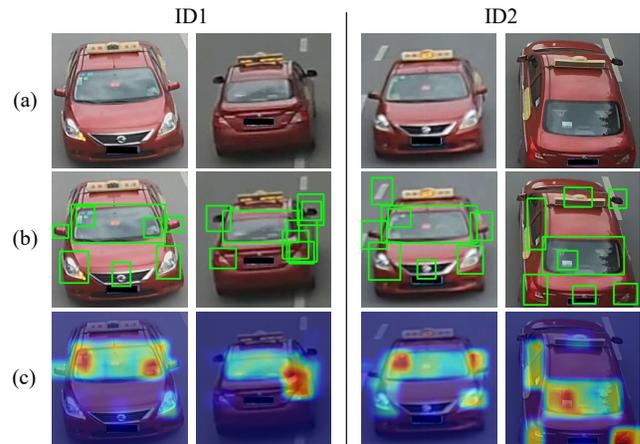
types and headlights [3]. However, most production vehicles can exhibit near-identical appearances since they may be mass-produced by the same manufacturer. When two vehicles with the same car model are presented, more fine-grained details (*e.g.*, annual service signs, customize paintings, and personal decorations) are required for comparison, as shown in Figure 1 (a) that ID1 looks similar like ID2 since they are from the same car mode. Therefore, the key challenge of vehicle IR lies in how to recognize the subtle differences between vehicles and locate the prominent parts that characterize their identities.

Most existing works focus on learning global appearance features with various vehicle attributes, including model type [6], [11], [12], license plate [11], spatial-temporal information [13], [14], orientation [5], [15], [16], [17], *etc*. The main disadvantage of global features is the lack of capability to capture more fine-grained visual differences, which is crucial in vehicle IR. Despite the help of auxiliary attributes, the supervision is still weak. For instance, the license plates are usually not available for privacy protection, while the two extremely similar vehicles from the same model type can not be distinguished (as shown in Figure 1). Also, they are easily degraded by the viewpoint variation, distortion, occlusion, motion blur and illumination, especially in the unconstrained real-world environment. Therefore, it is important to explore more robust and environment-invariant information to represent specific vehicles. Recent many works tend to explore subtle variances from car parts [18], [19], [20] to learn the

X. Zhang, R. Zhang and M. You are with Department of Control Science and Engineering, Tongji University, Shanghai 201804, China (e-mail: zhangxinyu@tongji.edu.cn; cxrfzhang@tongji.edu.cn; myyou@tongji.edu.cn). M. You is also with Shanghai Institute of Intelligent Science & Technology, Tongji University, Shanghai 201804, China.

J. Cao, D. Gong and C. Shen are with The University of Adelaide, Adelaide, SA 5005, Australia (e-mail: jonbakerfish@gmail.com; edgong01@gmail.com; chunhua.shen@adelaide.edu.au). JC, DG, CS and their employer received no financial support for the research, authorship, and/or publication of this article.

*Part of this work was done when X. Zhang was visiting The University of Adelaide. First two authors contributed to this work equally.

#Correspondence should be addressed to M. You.

local information. However, these methods mainly focus on the localization of the spatial part regions without considering how these regions are subject to attention with different degree.

To address above problems, we propose a novel *part-guided attention network* (PGAN) to improve the performance effectively by focusing on the most prominent part regions, which is implemented by integrating the bottom-up attention and the top-down attention systematically. Specially, we first utilize a bottom-up attention module to extract the related vehicle part regions, which is called the *part extraction module* in our work. With the established object detectors [3], [20] that are pre-trained on the vehicle attributes, we consider the extracted part regions from the part extraction module as candidates, which is beneficial for narrowing down the searching area for network learning. Importantly, this bottom-up attention can effectively take advantage of the context correlation among pixels in the same part via assigning same values for all pixels in a specific part region, which is superior to grid attention that gives no consideration on the pixel relationships. Besides, we call these candidate regions as *coarse part regions* since the quality of the detection may be not accurate with the pre-trained part extraction module and some part regions with less information may be included in these candidates.

To extract more effective local information, we apply a top-down attention process to select the most prominent part regions as well as assign appropriate importance scores to them after obtaining the above candidate part regions. Here, we introduce a *part attention module* (PAM), which is guided by the identification loss to allocate the importance for each coarse part region. PAM adaptively locates the discriminative regions with high-attention weights and suppresses the distraction of irrelevant parts with relatively low weights, as shown in Figure 1 (c). It is also beneficial for filling out the wrongly detected part regions by giving a weight near to zero (as shown in the rear view of ID2 in Figure 1). In detail, PAM can assign a special attention weight for each corresponding part region, *i.e.*, all pixels in this region share the same weight, reflecting the importance of the selected part regions by considering all pixels in a part region as a whole. Therefore, PAM is more efficient than grid attention or evenly decomposed part attention [21], [22], [19], [23], since PAM is able to provide more fine-grained attention which is conducted only on the selected part regions by taking the context information among pixels into consideration instead of all spatial pixels. We call these selected-weighted part regions as *fine part regions*.

With the combination of bottom-up and top-down attention, our attention mechanism can provide more prominent part regions for improving the feature representation. Finally, we aggregate the vehicle's holistic appearance and part characteristics with a *feature aggregation module* to improve the performance further. Figure 2 shows the whole training process. To summarize, our main contributions are as follows:

- We design a novel Part-Guided Attention Network (PGAN), which effectively combines part-guided bottom-up and top-down attention together to capture both local and global information.
- We propose to extract Top-$D$ part regions without part alignment to maintain more prominent yet less available parts
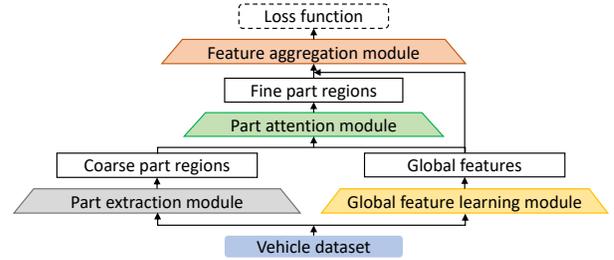


**FIG. 2:** The flow of the overall training process of our framework. The trapezoid represents the modules involved in our framework, while the solid rectangle denotes the outputs of the modules.

effectively in the part extraction module.

- We propose a part attention module (PAM) to evaluate the relative importance of the selected Top-$D$ part regions, which further focuses more on prominent parts and reducing the distraction of wrongly detected or irrelevant parts.
- Extensive experiments on four challenging benchmark datasets demonstrate that our proposed method achieves new state-of-the-art vehicle IR performance.

## II. RELATED WORK

### A. Global Feature-based Methods

**Feature Representation** Vehicle IR aims at learning discriminative feature representation to deal with significant appearance changes for different vehicles. Public large-scale datasets [3], [4], [6], [11], [8], [25], [26] are widely collected with annotated labels and abundant attributes under unrestricted conditions. These datasets face huge challenges on occlusion, illumination, low resolution and various views. One way to deal with these datasets uses deep features [4], [5], [11], [27], [26] instead of hand-crafted features to describe vehicle images. To learn more robust features, some methods [6], [11], [12], [13], [14], [28] try to explore details of vehicles using additional attributes, such as model type, color, spatial-temporal information, *etc*. Moreover, works of [15], [17] propose to use synthetic multi-view vehicle images from a generative adversarial network (GAN) [29] to alleviate cross-view influences among vehicles. In [5], [16] authors also implement view-invariant inferences effectively by learning a viewpoint-aware representation. Although great progress has been obtained by these methods, there is a huge drop when encountering invisible variances of different vehicles as well as large diversities in the same vehicle identity.

**Metric Learning** To alleviate the above limitation, deep metric learning methods [30], [31], [32], [33] use powerful distance metric expression to pull vehicle images in the same identity closer while pushing dissimilar vehicle images further away. The core idea of these methods is to utilize the matching relationship between image pairs or triplets as much as possible, which are widely used in IR works [34], [35], [24]. Whereas, sampling strategies in deep metric learning lead to suboptimal results and also lack of abilities to recognize more meaningful unobtrusive details. It is thus limited by the complex differences of the vehicle appearances.
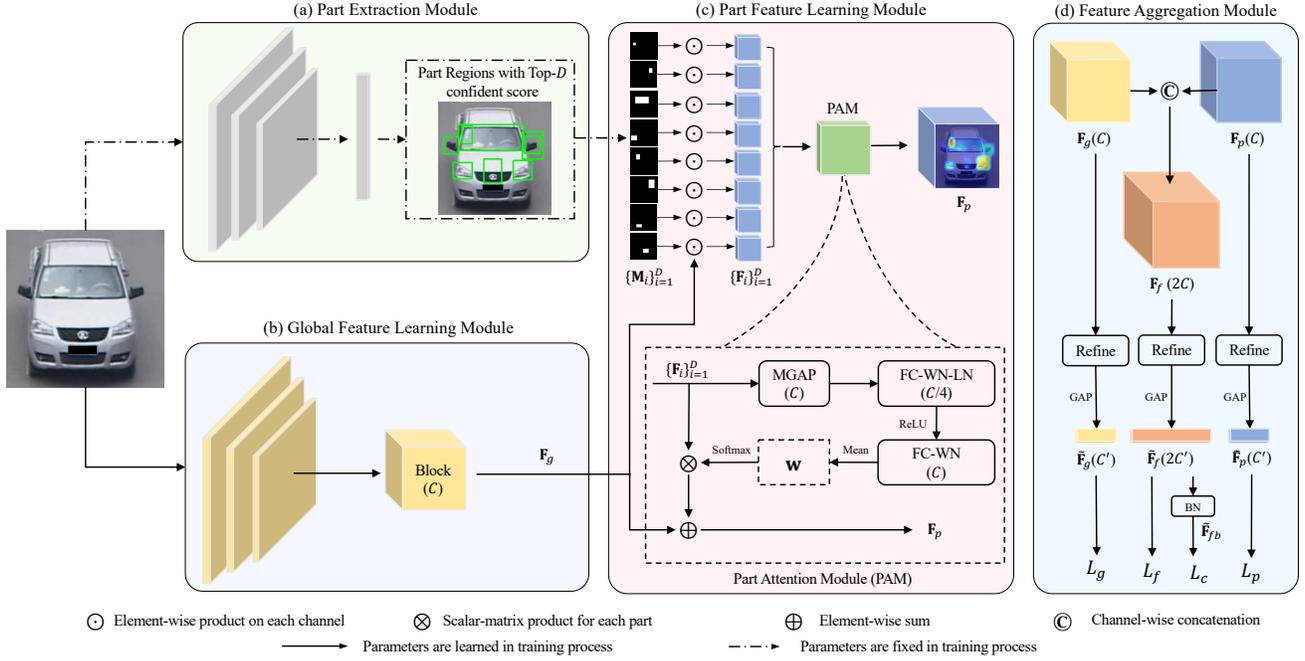
**FIG. 3:** Part-Guided Attention Network (PGAN) pipeline. The model consists of four modules: Part Extraction Module, Global Feature Learning Module, Part Feature Learning Module and Feature Aggregation Module. The input vehicle image is first processed to obtain the global feature $\mathbf{F}_g$ and the part masks $\{\mathbf{M}_i\}_{i=1}^{D}$ of Top-$D$ candidate parts. The part mask features $\{\mathbf{F}_i\}_{i=1}^{D}$ is then obtained via Eq. (2), after which $\{\mathbf{F}_i\}_{i=1}^{D}$ is fed into a Part Attention Module (PAM) to obtain the part-guided feature $\mathbf{F}_p$. PAM is a compact network, learning a soft attention weight $\mathbf{w} \in \mathbb{R}^{D}$, which is composed of a mask-guided average pooling (MGAP) layer and some linear and non-liner layers. Subsequently, the fusion feature $\mathbf{F}_f$ is obtained by concatenating $\mathbf{F}_g$ and $\mathbf{F}_p$. After the refinement and global average pooling (GAP) operation, $\widetilde{\mathbf{F}}_g$, $\widetilde{\mathbf{F}}_p$ and $\widetilde{\mathbf{F}}_f$ are all used for the optimization of triplet loss functions $L_f$, $L_g$ and $L_p$, respectively. Besides, as [24], $\widetilde{\mathbf{F}}_f$ is followed by a BN layer and the normalized feature $\widetilde{\mathbf{F}}_{fb}$ is used for optimizing the softmax cross-entropy loss $L_c$. Here, FC, WN, LN and BN represent fully-connected layer, weight normalization, layer normalization and batch normalization respectively. Mean denotes a channel-wise mean operation. $C$ and $C'$ are channel dimension before and after refine operation.

## B. Part Feature-based Methods

Similar as [36], [37], [38], [39], [40] focusing on object patches in other IR works, a series of part-based learning methods explicitly exploit the discriminative information from multi-part locations of vehicles. [20] provides an attribute detector while [41] provides a marker detector. [19], [21], [22], [23], [42] take great efforts on separating feature maps into multiple even partitions to extract specific features of respective regions. However, it is difficult for vehicles to directly apply this naive partitions since the vehicle appearances change a lot. In other words, almost all pedestrian images have relative regular appearances from top to bottom (representing head to feet), while vehicle appearances suffer from multiple views without unique commons. For example, the bottom partition of the front vehicle is wheels, while that is the vehicle back of the rear one in ID2 in Figure 1. Although [19] utilizes discriminative features from quadruple directions for each vehicle, it still suffers from the misalignment problem due to inaccurate grid partitions.

Another line of part-based methods [5], [43], [44], [45] bring informative key-points to put more attention on effective localized features. In particular, although [43] attempts to detect and use keypoints, it defines a heuristic rule to choose keypoints for every input image. Actually, [43] extracts a subgroup of keypoint features based on the vehicle orientation, in which the choice of the keypoint groups is manually pre-defined. [5] applies an aggregation module on the local

features based on orientation. However, keypoints in each orientation are treated equally and the detail information is easy to be ignored. In contrast, we adaptively learn a soft-attention for each detected local part feature, conditional on the input image. The soft-attention coefficients, measuring the importance of a local feature for the target task, are learnt by using the sole target identification loss. In other word, we do not rely on extra information while [43] uses orientation as extra supervision.

Besides, [18], [20] denote to design part-fused networks using ROI features of each part on vehicles from a pre-trained detection model to extract discriminative features. However, there is no importance selection on the candidate part regions in [18], which considers all part regions equally. Instead, our PGAN can select the most prominent part regions, *e.g.*, annual service signs and hungs, which are subtle yet important to distinguish different vehicles. In addition, we apply a single tailor-designed supervision for the soft-weighted part features together. Compared with [18] applying separate supervision to each part feature, our PGAN can provide more accurate supervision with the aggregated feature. Although [46], [47], [48] also utilize attention in the feature maps, the attention mechanism is applied on each pixel in the feature map. Our part attention module focuses on the pixel sets, *i.e.*, detected part regions on the feature maps. Thus, the context correlation in a same part can be integrally considered. In this way, we can not only consider all part features together as a whole but

also pay more attention to the prominent part regions as well as alleviate the influence of irrelative ones.

## III. METHODOLOGY

We firstly define each vehicle image as $x$ and the unique corresponding identity label as $y$. Given a training set $X^t = \{(x_n^t, y_n^t)\}_{n=1}^{N^t}$, the main goal of the vehicle IR is to learn a feature embedding function $\phi(x^t; \theta)$ for measuring the vehicle similarity under certain metrics, where $\theta$ denotes the parameters of $\phi(\cdot)$. It is important to learn a $\phi$ with good generalization on unseen testing images since there is no overlap identities in training and testing dataset. During testing, given a query vehicle image $x^q$, we can find vehicles with the same identity from a gallery set $X^g = \{(x_n^g, y_n^g)\}_{n=1}^{N_g}$ by comparing the similarity between $\phi(x^q; \theta)$ and each $\phi(x_n^g; \theta)$, $\forall x_n^g$.

In this section, we present the proposed Part-Guided Attention Network (PGAN) in detail. The overall framework is illustrated in Figure 3, which consists of four main components: *Part Extraction Module*, *Global Feature Learning Module*, *Part Feature Learning Module* and *Feature Aggregation Module*. We first generate the part masks of vehicles in the part extraction module, which are then applied on the global feature map to obtain the mask-guided part feature. After that, we learn the attention scores of different parts to enhance the part feature via increasing the weights of discriminative parts as well as decreasing that of less informative parts. Subsequently, the three refined features, *i.e.*, global, part, and fusion features are all used for model optimization.

### A. Global Feature Learning Module

For a vehicle image $x$, before obtaining the part features, we first extract a global feature map $\mathbf{F}_g \in \mathbb{R}^{H \times W \times C}$ with a standard convolutional neural network, as shown in Figure 3 (a). Most previous methods [34], [24] directly feed $\mathbf{F}_g$ into a global average pooling (GAP) layer to obtain the embedding feature that mainly considers the global information, which is studied as a *baseline* model in our experiments.

However, due to the lost of the spacial information after GAP, it is difficult to distinguish two near-identical vehicles, as illustrated in ID1 and ID2 in Figure 1. Therefore, it is crucial to maintain the spatial structure of feature maps, which helps describe the subtle visual differences. We thus directly apply $\mathbf{F}_g$ as one of the inputs for the following part learning process and the final optimization, and we explore a novel method to focus on the effective part regions following.

### B. Part Extraction Module

We first extract the part regions using a pre-trained SSD detector specially trained on vehicle attributes [20]. Here, we only consider 16 of all 21 vehicle attributes as shown in Table I. The reason is that the remaining attributes are vehicle styles, *i.e.*, "car", "trunk", "tricycle", "train" and "bus", representing the whole vehicle image which can be recognized as the global information in our paper. Once detected, we only use the confidence scores to select part regions and ignore the label information of each part. It is reasonable since not all attributes

**TABLE I:** Name and abbreviation of vehicle attributes used in our paper.

| Name | Abbreviation | Name | Abbreviation |
|---|---|---|---|
| annual service signs | anusigns | back mirror | backmirror |
| car light | carlight | carrier | carrier |
| car topwindow | cartopwindow | entry license | entrylicense |
| hanging | hungs | lay ornament | layon |
| light cover | lightcover | logo | logo |
| newer sign | newersign | tissue box | tissuebox |
| plate | plate | safe belt | safebelt |
| wheel | wheel | wind-shield glass | windglass |

are available in each vehicle due to the multi-view variation, so that it is hard to decide a universal rule for reliable part alignments (*i.e.*, selecting same part regions for all vehicles).

Instead of naively selecting relevant part regions by setting a threshold on the confidence scores, we select the most confident top-$D$ proposals as the candidate vehicle parts. The main reasons are twofold: 1) some crucial yet less confident bounding boxes, like annual service signs, play a crucial role in distinguishing different vehicle images; 2) part number is fixed, which is easy to learn the attention model in the following stage. Note that we want to ensure a high recall rate to avoid missing relevant parts. The irrelevant parts are filtered out from the subsequent top-down attention learning.

We use the index $i \in \{1, 2, ..., D\}$ to indicate each of the selected top-$D$ part regions. The spatial area covered by each part is denoted as $A_i$. For each candidate part region $i$, we obtain a binary mask matrix $\mathbf{M}_i \in \{0, 1\}^{H \times W}$ by assigning 1 to the elements inside the part region $A_i$ and 0 to the rest, denoted as:

$$\mathbf{M}_i(\text{pix}) = \begin{cases} 1, & \text{if } \text{pix} \in A_i \\ 0, & \text{if } \text{pix} \notin A_i \end{cases}, \forall i, \qquad (1)$$

where pix indicates a pixel location of $\mathbf{M}_i$. Note that the size of each $\mathbf{M}_i$ is the same as a single channel of $\mathbf{F}_g$. It means that if the parameters of the neural network or the sizes of input images change, the corresponding part locations on $\mathbf{M}_i$ will be changed accordingly and the spacial area $A_i$ is also changed. Although $\mathbf{M}$ can be scaled based on the input of multi-scale images, we resize all images to the same resolution for simplification and thus the size of $\mathbf{M}$ can be regularized to $H \times W$. Besides, during processing, we force all $A_i$ in the range of $H \times W$ to ensure all part regions are located in the range of image areas (*i.e.*, the size of $H \times W$).

After obtaining global feature $\mathbf{F}_g$ and part masks $\{\mathbf{M}_i\}_{i=1}^D$, we project the part masks on the feature map $\mathbf{F}_g$ to generate a set of mask-based part feature representations $\{\mathbf{F}_i\}_{i=1}^D$, which will be taken as the input of the following part feature attention module. For each part region $i$, we can obtain $\mathbf{F}_i$ via the following formula:

$$\mathbf{F}_i = \mathbf{M}_i \odot \mathbf{F}_g, \quad \forall i \in \{1, 2, ..., D\}, \qquad (2)$$

where $\odot$ denotes the element-wise product operation on each channel of $\mathbf{F_g}$. $\mathbf{F}_i$ is the mask-based part feature map of the $i$-th part region. Note that all $\mathbf{F}_i \in \mathbb{R}^{H \times W \times C}$. In each $\mathbf{F}_i$, only the elements in the regions of $i$-th part are activated. The illustration is shown in Figure 3 (c).

We learn an attention module on the part regions in the following section. Unlike the traditional grid attention method
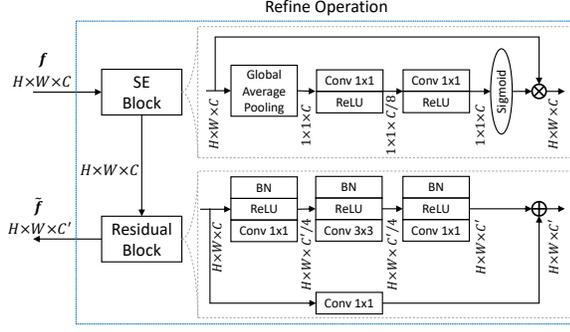
**FIG. 4:** The module structure in the refine operation. $C$ and $C'$ are the size of channel of feature maps before and after the refine operation.

that processes a set of uniform grids, our attention model can focus on the prominent parts by only activating the selected parts. The irrelevant parts can thus be ignored directly. Besides, the context correlation in a same part can be integrally considered, alleviating missing of essential features. Moreover, this part extraction process can be considered as a bottom-up attention mechanism [49] with a set of candidate images regions proposed.

### C. Part Feature Learning Module

Part feature learning module is to produce a weight map across the mask-based part feature maps $\{\mathbf{F}_i\}$. In this way, the network can focus on specific part regions. Recent methods [18], [50] highlight all part regions equally and thus ignores the importance discrepancy among different part regions. Besides, some detected parts might not be informative for some specific cases, such as wrongly detected background or windshield without useful information, which tends to result in degraded results. To tackle the above problems, we propose a *part attention module* (PAM) to adaptively learn the importance of each part so as to take more attention to the most discriminating regions and suppress those with less information. Consequently, PAM can be considered as a part-based top-down attention mechanism, since this attention signal is supervised by the specific identification task to predict an importance distribution over candidate image regions.

**Part attention module (PAM)** Our PAM is designed to obtain a part-guided feature representation $\mathbf{F}_p \in \mathbb{R}^{H \times W \times C}$ relying on a top-down attention mechanism on candidate part regions. From PAM, we can obtain a soft weight vector $\mathbf{w} \in \mathbb{R}^D$ to indicate the importance of each part region, thus the part-guided feature representation $\mathbf{F}_p$ can be obtained as:

$$\mathbf{F}_p = \sum_{i=1}^{D} w_i \mathbf{F}_i + \mathbf{F}_g, \quad (3)$$

where $w_i \in [0, 1]$ denotes the $i$-th element of the soft weight $\mathbf{w}$, which represents a learned weight of $i$-th part feature $\mathbf{F}_i$ obtained via Eq. (4). $\mathbf{w}$ is normalized with sum as 1 so that the relative importance between different parts is obvious. Here, $\mathbf{F}_g$ is added to augment the capability of part regions.

We learn a compact model to predict the attention weights $\mathbf{w}$ for measuring the different importance of each selected part,

as shown in Figure 3 (c). Specifically, we first use a mask-guided global average pooling operation (MGAP) on each $\mathbf{F}_i$ and then learn a mapping function with a softmax layer to obtain $\mathbf{w}$. Each element $w_i$ can be predicted by:

$$w_i = \frac{\exp(\psi(\mathrm{mgap}(\mathbf{F}_i, \mathbf{M}_i), \theta_\psi))}{\sum_{j=1}^{D} \exp(\psi(\mathrm{mgap}(\mathbf{F}_j, \mathbf{M}_j), \theta_\psi))}, \quad (4)$$

where $\psi(\cdot)$ denotes a learnable function that is able to highlight the most important part regions with high values (as shown in Figure 3 (c)). $\theta_\psi$ is the parameter of mapping function $\psi(\cdot)$, and $\mathrm{mgap}(\cdot)$ denotes MGAP operation discussed in the following.

Before feeding $\mathbf{F}_i$ into $\psi$, we average each channel of $\mathbf{F}_i$ as a scalar via the $\mathrm{mgap}(\cdot)$ operator. Note that, in each $\mathbf{F}_i$, only the elements in the part region $i$ are activated and most of the elements in $\mathbf{F}_i$ are zero. Instead of performing the standard global average pooling (GAP), we restrict the average pooling in the areas indicated by the mask $\mathbf{M}_i$ via the MGAP operator. In detail, for each channel of $\mathbf{F}_i$, after summing the nonzero elements, the MGAP operator devides the sum value with the number of elements (*i.e.* $||\mathbf{M}_i||_1 < H \times W$), instead of the number of total elements (*i.e.* $H \times W$) in the GAP.

### D. Feature Aggregation Module

Since global and part-based features provide complementary information, we concatenate the global feature $\mathbf{F}_g$ and part-guided feature $\mathbf{F}_p$ together, which is then denoted as fusion feature $\mathbf{F}_f \in \mathbb{R}^{H \times W \times 2C}$. Furthermore, we adopt a *Refine* operation on $\mathbf{F}_f$ to reduce the dimension of feature representation to speed up the training process. The *Refine* operation is composed of a SE Block [51] and a Residual Block [52], which is illustrated in Figure 4. After a global average pooling (GAP) layer, the refined fusion feature $\widetilde{\mathbf{F}}_f \in \mathbb{R}^{2C'}$, $\widetilde{\mathbf{F}}_g \in \mathbb{R}^{C'}$ and $\widetilde{\mathbf{F}}_p \in \mathbb{R}^{C'}$ are obtained for the whole model optimization. Here, $C'$ is the size of channel of feature maps after the refine operation, while $C$ is that before the refine operation. Note that following [24], an additional batch normalization(BN) layer is adopted on $\widetilde{\mathbf{F}}_f$. It is proved to be beneficial for optimizing the softmax cross-entropy loss [24]. Here, we denote the feature after the BN layer as $\widetilde{\mathbf{F}}_{fb}$.

### E. Model Training

In the training process, we adopt softmax cross-entropy loss and triplet loss [34] as a joint optimization. In specific, we apply triplet loss on $\widetilde{\mathbf{F}}_f$ and softmax cross-entropy loss on $\widetilde{\mathbf{F}}_{fb}$, denoted as $L_f$ and $L_c$. In order to make full use of the global and part information separately, we also optimize the refined global feature $\widetilde{\mathbf{F}}_g$ and part-guided feature $\widetilde{\mathbf{F}}_p$ with triplet loss, which are denoted as $L_g$ and $L_p$, respectively. Overall, the total loss function can be formulated as:

$$L = \lambda L_c + L_{tri} = \lambda L_c + L_f + L_g + L_p, \quad (5)$$

where $\lambda$ is the loss weight to trade off the influence of two types of loss functions, *i.e.*, softmax cross-entropy loss $L_c$ and triplet loss $L_{tri}$. Experiments show that joint optimization could improve the ability of feature representation.

For evaluation, we use the normalized fusion feature $\widetilde{\mathbf{F}}_{fb}$ as the final feature representation in our work.

**TABLE II:** Comparison on the different optimization methods of PGAN on VeRi-776. $\widetilde{\mathbf{F}}_{fb}$ is used as the feature representation. For fairness, the feature dimension is fixed to 512 for all methods including the baseline model. We omit the loss weight in Eq. 5 for clarity and set $\lambda$ to 2 here.

| Method | Optimization | mAP | Top-1 | Top-5 |
|---|---|---|---|---|
| Baseline | $L_c + L_f$ | 75.7 | 95.2 | 98.2 |
| Ours | $L_c + L_f$ | 77.7 | 95.9 | **98.5** |
| | $L_c + L_f + L_g$ | 78.0 | 95.1 | 97.7 |
| | $L_c + L_f + L_p$ | 78.5 | 95.8 | 98.3 |
| | $L_c + L_f + L_g + L_p$ (PGAN) | **79.3** | **96.5** | 98.3 |

**TABLE III:** Performance comparison on different attention methods, *i.e.*, grid attention, PGAN without Part Attention Module (PAM) and our PGAN on VeRi-776.

| Method | Dimension | mAP | Top-1 | Top-5 |
|---|---|---|---|---|
| Baseline | 256 | 75.3 | 95.3 | 98.2 |
| Grid Attention | | 76.1 | 95.3 | 97.7 |
| PGAN w/o PAM | | 77.9 | 95.6 | **98.4** |
| PGAN | | 78.6 | 95.4 | 98.0 |
| Baseline | 512 | 75.7 | 95.2 | 98.2 |
| Grid Attention | | 77.0 | 95.8 | 98.0 |
| PGAN w/o PAM | | 78.0 | 95.5 | 98.2 |
| PGAN | | **79.3** | **96.5** | 98.3 |

## IV. Experiments

### A. Datasets and Evaluation Metrics

We evaluate our PGAN method on four public large-scale Vehicle IR (*a.k.a.*, re-identification) benchmark datasets.

*VeRi-776* [11] is a challenging benchmark in vehicle IR task that contains about $50,000$ images of 776 vehicle identities across 20 cameras. Each vehicle is from 2-18 cameras with various viewpoints, illuminations and occlusions. All datasets are split into a training set with $37,778$ images of 576 vehicles and a testing set with $11,579$ images with 200 vehicles.

*VehicleID* [25] is a widely-used vehicle IR dataset which contains vehicle images captured in the daytime by multiple cameras. There are total of $221,763$ images with $26,267$ vehicles, where each vehicle has either front or rear view. The training set contains $110,178$ images of $13,134$ vehicles while the testing set comprises $111,585$ images of $13,133$ vehicles. The evaluation protocol of the large test subset VehicleID is randomly selecting one image from each vehicle to generate a gallery set (2400 images) while the remaining images are used as query set. The random selection process was repeated for 10 times and the mean result is used as the final performance.

*VRIC* [26] is a realistic vehicle IR benchmark with unconstrained variations of images in resolution, motion blur, illumination, occlusion, and multiple viewpoints. It contains $60,430$ images of $5,622$ vehicle identities captured from 60 different traffic cameras during both daytime and nighttime. The training set has $54,808$ images of $2,811$ vehicles, while the rest is used for testing with $5,622$ images of another $2,811$ vehicle IDs.

*VERI-Wild* [8] is recently released with $416,314$ vehicle images of $40,671$ IDs captured by 174 cameras. The training set consists of $30,671$ IDs with $277,797$ images. The small test subset consists of $3,000$ IDs with $41,816$ images while the medium and large subset consist of $5,000$ and $10,000$ IDs with $69,389$ and $138,517$ images respectively.

*Evaluation metrics.* To measure the performance for vehicle IR, we utilize the Cumulated Matching Characteristics (CMC) and the mean Average Precision (mAP) as evaluation criterions. The CMC calculates the cumulative percentage of correct matches appearing before the top-$K$ candidates. We report Top-1 and Top-5 scores to represent the CMC criterion. Given a query image, Average Precision (AP) is the area under the Precision-Recall curve while mAP is the mean value of AP across all query images. The mAP criterion reflects both precision and recall, which provides a more convincing evaluation on IR task.

### B. Implementation Details

*Part extraction.* we directly conduct the inference process to extract part regions using the pretrained detector [20]. There is no re-train or finetune process in our method since the attribute annotations in the four datasets, *i.e.*, VeRi-776, VehicleID, VRIC and VERI-Wild, are not available. In the training process of the detector in [20], which is based on the SSD model [53], the VOC21_S dataset [20] is used as the training data. The VOC21_S dataset is captured during both daytime and nightime by multiple real-world cameras in several cities, so that this dataset shares similar scenarios with the four datasets we used. Since these datasets are collected by different cameras in not exactly the same environment, there is a domain gap issue to some extent. During the inference, the NMS threshold is set to 0.45 in all experiments. For each image, we extract Top-$D$ part regions according to confident scores, where $D = 8$ without specification.

*Vehicle IR model.* We adopt ResNet50 [52] without the last classification layer as the backbone model in the global feature learning module, which is pre-trained on ImageNet [54] initially. The model modification follows [24], *i.e.*, removing the last downsample operation and adding a BN layer before softmax cross-entropy loss.

All images are resized to 224×224. The data augmentations, *i.e.*, random horizontal flipping and random erasing [55] with a probability of 0.5, are used as in [24]. We use Adam optimizer [56] with a momentum of 0.9 and a weight decay $5 \times 10^{-4}$. For all experiments without other specification, we set the batch size to 64 with 16 IDs randomly selected. The learning rate starts from $1.75 \times 10^{-4}$ and is multiplied by 0.5 every 20 epochs. The total number of epochs is 130.

### C. Ablation Study

*1) Effectiveness of joint optimization:* We first design an ablation experiment analyzing the effectiveness of joint optimization with different features and loss functions. For our method, we use the normalized feature $\widetilde{\mathbf{F}}_{fb}$ as the feature representation and fix the feature dimension of $\widetilde{\mathbf{F}}_{fb}$ to 512, *i.e.*, $C' = 256$. For a fair comparison, we also set the feature dimension to 512 in a baseline model. As reported in Table II, we can observe that only using optimization on the fusion feature, *i.e.*, $L_c + L_f$, can improve the performance by 2% on mAP comparing with baseline model, which confirms that PAM can provide important part information that is better for model optimization. After adding $L_g$ and $L_p$ separately, mAP can improve by about 1%. It shows that combining
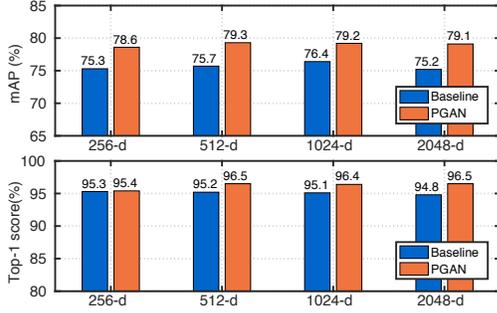
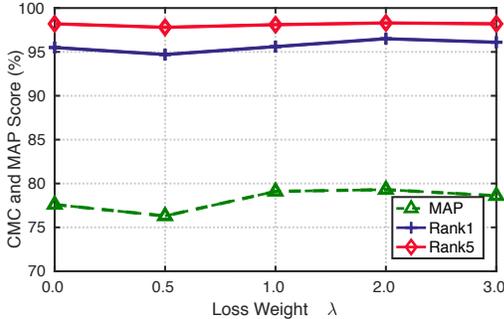**FIG. 5:** Parameter analysis of the feature dimension on VeRi-776.



**FIG. 6:** Parameter analysis of the loss weight $\lambda$ on VeRi-776.

**TABLE IV:** Performance comparison on different part number $D$ of PGAN on VeRi-776. $\widetilde{\mathbf{F}}_{fb}$ is used as the feature representation.

| Part Number $D$ | Dimension 256 | | | Dimension 512 | | |
|---|---|---|---|---|---|---|
| | mAP | Top-1 | Top-5 | mAP | Top-1 | Top-5 |
| Baseline | 75.3 | 95.3 | 98.2 | 75.7 | 95.2 | 98.2 |
| 4 | 76.9 | 95.4 | 97.8 | 76.8 | 94.8 | 98.0 |
| 6 | 78.5 | **95.8** | 98.0 | 78.7 | 96.2 | 98.0 |
| 8 | **78.6** | 95.4 | 98.0 | **79.3** | **96.5** | **98.3** |
| 10 | 77.6 | 95.5 | **98.3** | 79.1 | 95.9 | 98.2 |
| 12 | 77.9 | 94.7 | 97.9 | 77.5 | 95.6 | 98.1 |

**TABLE V:** Performance comparison on different part number $D$ of PGAN on VRIC and VERI-Wild (large subset). $\widetilde{\mathbf{F}}_{fb}$ is used as the feature representation. The feature dimension is fixed to 512.

| Part Number $D$ | VRIC | | | VERI-Wild (large) | | |
|---|---|---|---|---|---|---|
| | mAP | Top-1 | Top-5 | mAP | Top-1 | Top-5 |
| Baseline | 83.5 | 76.1 | 93.0 | 69.4 | 88.1 | 95.4 |
| 4 | 84.1 | 76.8 | 93.0 | **70.8** | 89.5 | **95.9** |
| 6 | 84.6 | 77.6 | 93.6 | 70.7 | 89.4 | 95.8 |
| 8 | **84.8** | **78.0** | 93.2 | 70.6 | 89.2 | 95.7 |
| 10 | 84.1 | 76.8 | 93.2 | 70.7 | 89.6 | 95.8 |
| 12 | 84.3 | 77.1 | **93.8** | 70.6 | **89.7** | 95.8 |

with the additional optimizations on the global and part feature can provide more useful information for the model training. Furthermore, with the joint optimization with all these loss functions, the result improves to 79.3% mAP, which outperforms the baseline model by 3.6%.

*2) Analysis of different attention method:* We first implement traditional grid attention by removing part extraction module, *i.e.*, PAM is directly used on each grid of $\mathbf{F}_g \in \mathbb{R}^{H \times W \times C}$. As shown in Table III, grid attention can only achieve 77.0% mAP and 95.8% Top-1 accuracy when the feature dimension is 512, showing that part guidance is crucial for filtering invalid information like background. Moreover, we also use the identical weight for each part region by removing PAM. It can be seen as a bottom-up attention with the part guidance from a detection model. From Table III, we can find 0.7% and 1.3% mAP decrease when feature dimension is 512 and 256 without PAM. It proves that PAM is beneficial for focusing on prominent parts as well as suppressing the impact of some wrongly detected or useless regions. We exactly note that our PGAN w/o PAM is still better than grid attention by 1.0% mAP, which also proves the important role of the part-guided bottom-up attention.

*3) Parameter analysis of the feature dimension:* We first analyze the effectiveness of different feature dimension. The dimension $2C'$ of fusion feature $\widetilde{\mathbf{F}}_{fb}$ on VeRi-776 is used as the variable. As shown in Figure 5, our PAM module has consistent improvement compared with the baseline model whatever the dimension is. In particular, when $2C' = 2048$, our PGAN outperforms baseline model by 3.9% and 1.7% in mAP and Top-1 respectively. Besides, it is worth noting that our PGAN with low dimension still performs better than baseline with high dimension. For example, our PGAN

with 256 dimension surpasses the baseline model with 512 dimension by a large margin (78.6% *vs.* 75.7% mAP), which highly proves the effectiveness of our PGAN.

*4) Parameter analysis of the loss weight $\lambda$:* In Figure 6, we conduct experiments to compare different values of the loss weight $\lambda$ in Eq. (5), which evaluates the trade off between softmax cross-entropy loss and triplet loss. When $\lambda = 0$, we only use triplet loss on $\widetilde{\mathbf{F}}_f$ as the optimization. It is clear that when adding softmax cross-entropy loss on $\widetilde{\mathbf{F}}_{fb}$ into the model, our approach can obtain further improvement when the range of $\lambda$ is from 1.0 to 2.0. However, too small $\lambda$ and too large $\lambda$ both lead to the bad influence for the model training. We believe that there is a trade off between these two-type loss functions. Too large $\lambda$ means the relation restriction among samples from the triplet loss does less effort for the model optimization, while too small $\lambda$ means less effectiveness of the global structure from the softmax cross-entropy loss. Form Figure 6 we can see that the best result is obtained when $\lambda$ is set to 2. Without specification, we use $\lambda = 2$ as the default loss weight in our paper.

*5) Parameter analysis of the number of part regions $D$:* In addition, we analyse how the number of part regions $D$ in the part extraction module affects the IR results. We test the performance with $D = \{4, 6, 8, 10, 12\}$ of our PGAN on VeRI-776 in Table IV and on VRIC and VERI-Wild in Table V. The evaluated feature dimension is set to 256 and 512.

As shown in Table IV and Table V, there is a consistent improvement when utilizing the part guidance in our PGAN compared with the baseline model, which clearly verifies the effectiveness of our PGAN method. When the part number $D$ is not large, our PGAN can gradually improve the IR performance with the number of part regions increasing. It shows that the fixed number of part regions is able to narrow down the possible searching regions, which is helpful for focusing on the valid part components. Besides, our PAM can further improve the effectiveness of the part guidance by applying more concentration on the prominent part regions. Especially, when $D$ is changed from 6 to 8 on VeRi-776

**TABLE VI:** Comparisons with state-of-the-art IR methods on VeRi-776, VehicleID, VRIC and VERI-Wild. In each column, the first and second highest results are highlighted by red and blue respectively. The results of Siamese-CNN+Path-LSTM [14] and OIFE [5] on VRIC is reported by MSVR [26]. * denotes that VANet uses ResNet50 for VehicleID dataset while uses GoogLeNet for VeRi.

| Method | Backbone | VeRi-776 | | VehicleID | | | VRIC | | | VERI-Wild | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | Small | | Medium | | Large | |
| | | mAP | Top-1 | mAP | Top-1 | Top-5 | mAP | Top-1 | Top-5 | mAP | Top-1 | mAP | Top-1 | mAP | Top-1 |
| FACT+Plate-SNN+STR [11] | GoogleNet | 27.8 | 61.4 | - | - | - | - | - | - | - | - | - | - | - | - |
| Siamese+Path-LSTM [14] | ResNet50 | 58.3 | 83.5 | - | - | - | - | 30.6 | 57.3 | - | - | - | - | - | - |
| OIFE [5] | GoogleNet | 51.4 | 92.4 | - | 67.0 | 82.9 | - | 24.6 | 51.0 | - | - | - | - | - | - |
| PROVID [13] | GoogleNet | 53.4 | 81.6 | - | - | - | - | - | - | - | - | - | - | - | - |
| VAMI [16] | Self-design | 50.1 | 77.0 | - | 47.3 | 70.3 | - | - | - | - | - | - | - | - | - |
| MSVR [26] | MobileNet | 49.3 | 88.6 | - | 63.0 | 73.1 | - | 46.6 | 65.6 | - | - | - | - | - | - |
| RNN-HA [12] | ResNet50 | 56.8 | 74.8 | - | 81.1 | 87.4 | - | - | - | - | - | - | - | - | - |
| SCAN [46] | VGG16 | 49.9 | 82.2 | - | 65.4 | 78.5 | - | - | - | - | - | - | - | - | - |
| RAM [23] | VGGM | 61.5 | 88.6 | - | 67.7 | 84.5 | - | - | - | - | - | - | - | - | - |
| AAVER [43] | ResNet50 | 66.4 | 90.2 | - | 63.5 | 85.6 | - | - | - | - | - | - | - | - | - |
| Part-Regular [18] | ResNet50 | 74.3 | 94.3 | - | 74.2 | 86.4 | - | - | - | - | - | - | - | - | - |
| FDA-Net [8] | Self-design | 55.5 | 84.3 | 61.8 | 55.5 | 74.7 | - | - | - | 35.1 | 64.0 | 29.8 | 57.8 | 22.8 | 49.4 |
| QD-DLF [19] | Self-design | 61.8 | 88.5 | 68.4 | 64.1 | 83.4 | - | - | - | - | - | - | - | - | - |
| VANet [57]* | ResNet50 | 66.3 | 89.8 | - | 80.4 | 93.0 | - | - | - | - | - | - | - | - | - |
| TAMR [48] | ResNet18 | - | - | 61.0 | 59.7 | 73.9 | - | - | - | - | - | - | - | - | - |
| GRF+GGL [42] | VGGM | 61.7 | 89.4 | - | 70.0 | 87.1 | - | - | - | - | - | - | - | - | - |
| MVAN [47] | ResNet50 | 72.5 | 92.6 | 76.8 | 72.6 | 83.1 | - | - | - | - | - | - | - | - | - |
| Baseline [24] | ResNet50 | 75.7 | 95.2 | 83.5 | 77.5 | 91.0 | 83.5 | 76.1 | 93.0 | 82.6 | 94.0 | 77.2 | 91.7 | 69.4 | 88.1 |
| PGAN | ResNet50 | 79.3 | 96.5 | 83.9 | 77.8 | 92.1 | 84.8 | 78.0 | 93.2 | 83.6 | 95.1 | 78.3 | 92.8 | 70.6 | 89.2 |

dataset, the performance can be improved by 3% to 3.6% in mAP when the feature dimension is 512 and 3.2% to 3.3% in mAP when the feature dimension is 256 compared with the baseline model. When $D = 8$, we can obtain the relatively best result. However, the performance decreases when the part number continually increases. The reasons are twofold: 1) many detected part regions are covered with each other, which provide no further part information for the model learning; 2) more wrongly detected parts are extracted that results in the distraction of the model learning via providing large invalid information. We believe that if we use a better detector, the performance will be further improved.

The similar trend is observed on VRIC dataset, as shown in Table V. However, for VERI-Wild dataset, our PGAN is relatively robust to the part number. The reason is that images in VERI-Wild are high resolution and the part regions are detected more accurately. A few part regions are satisfactory to distinguish different vehicles. Although there exists an optimal $D$ for a specific dataset, we use 8 as the default setting for simplification.

*6) Effectiveness on different baseline:* In order to fully verify the effectiveness of our PGAN, we apply our method on various of baseline models. As shown in Figure 7, we can see that deeper backbone is beneficial for the performance, *e.g.*, ResNet18 baseline achieves 71.0% mAP while ResNet50 baseline 75.7%. In particular, we discard the last three FC layers in VGGM [58] backbone and the last FC layer in GoogleNet [59] backbone to insert our PGAN into the model. PGAN gains 2.8% and 1.7% mAP increases when applied on the GoogleNet baseline and VGGM baseline respectively. This validates the effectiveness of our proposed PGAN, which can be used as a general module for other tasks to some extent. We can also see that the improvement on VGGM is less than PGAN applied on ResNet50 baseline (3.6% mAP). The reason might be that the size of channels of the output features from
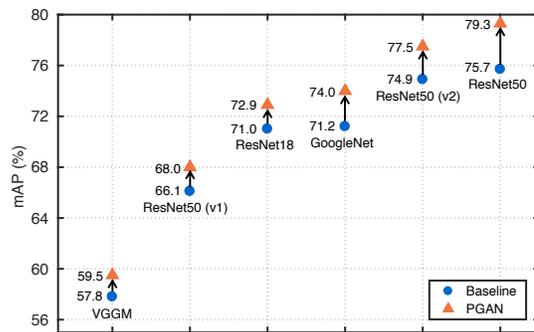


**FIG. 7:** Effectiveness on different baseline. v1 and v2 are the variation versions of our used baseline ResNet50 model [24]. v1 denotes the plain ResNet50 (last stride is 2) in [52] without data augmentation [55]. v2 denotes the v1 with [55]. Here, we remove the downsampling and use [55] in ResNet18, GoogleNet and ResNet50 for fair comparison.

the VGGM and GoogleNet backbone is 512 and 1024, which is smaller than that of ResNet50 (*i.e.*, 2048). Therefore, PAM module has more robust ability of feature representation in ResNet50 than that in VGGM and GoogleNet backbone.

Moreover, we can observe that some training methods in [24] are beneficial for the performance increase. In particular, when we use ResNet50 (v1), the original model with last stride as 2 and without data augmentation [55], our PGAN obtains about 3% improvement. Our PGAN is also useful for the version of ResNet50 (v2), *i.e.*, adding [24] on ResNet50 (v1). We can get 2.6% mAP improvement when comparing with the baseline. Overall, our PGAN can provide consistent improvement whatever the baseline is.

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

Finally, we compare our PGAN against other state-of-the-art vehicle IR methods, shown in Table VI. All reported results of our method are based on 512-dimension $\widetilde{\mathbf{F}}_{fb}$.
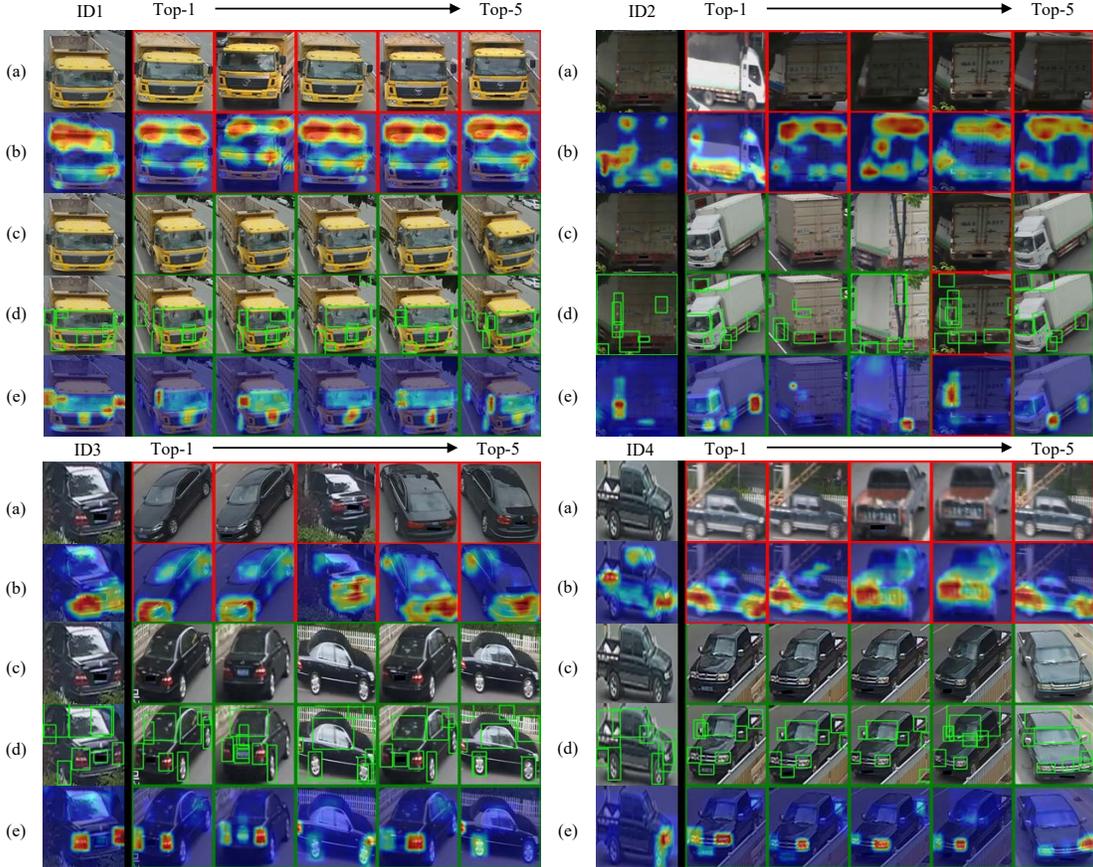
**FIG. 8:** Illustration of visualized comparison between traditional grid attention and our PGAN on VeRi-776 dataset. For a query image, we draw: (a) Top-5 retrieval results and (b) the corresponding heatmaps of $\mathbf{F}_p$ from PAM in grid attention; (c) Top-5 retrieval results, (d) the detected candidate part regions and (e) the corresponding heatmaps of $\mathbf{F}_p$ from PAM in PGAN. The correct and false matched vehicle images are enclosed in green and red rectangles respectively. It shows that our PGAN can put more attention on the most prominent part regions, such as back mirrors, windshield stickers and car brands. However, the grid attention mainly focuses on some insignificant regions like the car roof, resulting in the attention distracting. (Best viewed in color)

For VeRi-776, we strictly follow the cross-camera-search evaluation protocol as [11]. From Table VI, it is clear that our PGAN outperforms all the existing method for a large margin. For instance, the performance of PGAN is better than the state-of-the-art method, *i.e.*, Part-Regular [18], for 5% mAP and 2.3% Top-1 respectively. RAM [23] concatenates all the global and local features together as the final representation, which achieves 61.5% mAP with VGGM backbone. Similar as RAM, our PGAN can achieve 63.6% mAP with VGGM when combining $\widetilde{\mathbf{F}}_{fb}$, $\widetilde{\mathbf{F}}_g$ and $\widetilde{\mathbf{F}}_p$ together.

For VehicleID, we only report the result of the large test subset on Top-1 and Top-5. Our method surpasses almost all the methods except RNN-HA [12] at Top-1 and VANet [57]. Notice that RNN-HA uses the additional supervision of the vehicle model and the size of input image is $672 \times 672$ (9 times bigger than ours). However, as reported in [12], the performance of RNN-HA is extremely dropped by a large margin on VeRi-776 when the image size is set to $224 \times 224$, which is lower than our PGAN for about 22% in Top-1. In addition, VANet uses a specific viewpoint-based loss function, in which viewpoint labels are generated from a viewpoint model that is trained on manually annotated training samples. It is specially good for the front and rear view that appears in all vehicles in VehicleID. For VeRi-776, containing multi-

view vehicles, the mAP of VANet is lower than our PGAN by 13%. Since VANet only reports the result on GoogleNet backbone on VeRi-776, we also use the same backbone in our PGAN. From Figure 7, we can see that our PGAN obtains 71.2% mAP using GoogleNet as the backbone model, which is largely higher than VANet (66.3% mAP). It means that our PGAN is more beneficial for improving the performance in the multi-view scenario.

For VRIC, one of the largest dataset in vehicle IR, our proposed PGAN achieves satisfactory performance with 78.0% mAP and 93.2% Top-1. Note that we only use single resolution in both training and inference stages, achieving higher performance than MSVR [26] using multi-scale feature representations. We can also observe that although our baseline model has achieved satisfactory results, our PGAN can still improve the performance. It proves that although suffering from extreme motion blur, low resolution and various complex environment in VRIC dataset, our method can still extract useful and valid information. Since the detector [20] is applied without finetune, we believe that with a more accurate detector, our PGAN is able to perform better.

VERI-Wild is a newly released large vehicle dataset with more unconstrained variations in resolutions, illuminations, occlusion, and viewpoints, *etc*. There are only a few meth-
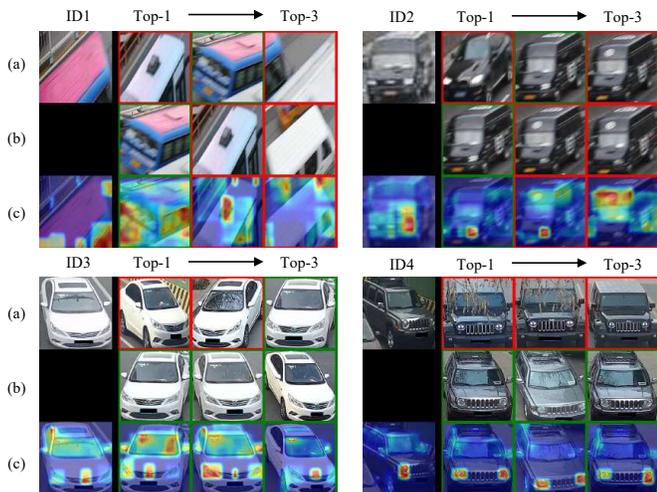
**FIG. 9:** Visualization of Top-3 retrieval images of baseline and our PGAN. ID1 and ID2 are from VRIC dataset, while ID3 and ID4 are from VERI-Wild dataset. For a query image, we draw: (a) Top-3 retrieval results from baseline; (b) Top-3 retrieval results and (c) the corresponding heatmaps of $\mathbf{F}_p$ from PAM in PGAN. The correct and false matched vehicle images are enclosed in green and red rectangles respectively.

**TABLE VII:** The average running time (ms) per frame on VeRi-776.

| Method | Module | | |
|---|---|---|---|
| | Detector | IR | Total |
| Baseline | - | 11.8 | 11.8 |
| PGAN | 39.2 | 12.3 | 51.5 |



**FIG. 10:** Illustrations of some failed samples on the VRIC dataset.

ods that have reported the results. Table VI shows that our proposed PGAN achieves great improvement compared with other methods, *e.g.*, achieving 70.6% mAP at the large test subset. FDA-Net [8] uses grid attention module to strengthen the model ability on local subtle differences, which performs worse than our PGAN. For fairness, we also conduct grid attention instead of PAM in our method and achieve 70.2% mAP, showing that our method is more useful via focusing on subtle differences. Moreover, we apply VGGM as the backbone model for VeRi-776, as shown in Figure 7. Our PGAN gets 59.5% mAP that is better than FDA-Net (55.5%).

We also report the result of the baseline model on all dataset. Note that, we also set the feature dimension to 512 in baseline model for fairness since the feature dimension of our fusion feature $\widetilde{\mathbf{F}}_{fb}$ is set to 512. Experiments show that our method achieves higher results than the baseline model in all datasets.

### E. Visualization

In this section, we visualize some retrieval results of the baseline, grid attention and our part-guided attention method (PGAN), repetively. As shown in Figure 8, we illustrate four different query vehicle images and their corresponding Top-5 most similar images as well as the heatmaps of $\mathbf{F}_p$ from the gallery set on VeRi-776 dataset. Meanwhile, we illustrate the Top-3 retrieval results on VRIC and VERI-Wild datasets in Figure 9 to show the effectiveness of our PGAN. In detail, the main advantages of our PGAN can be summarized as follows:

*1) Insensitive to various situations:* Our PGAN can extract more robust feature representation so as to significantly improve the IR performance. As shown in the ID2 and ID3 in Figure 8, given a rear vehicle image, we can not only find the easy vehicles from the rear views, but also get the side-view vehicle images that are difficult to recognize even by humans. In contrast, the grid attention can only focus on the images from the nearly same views. Moreover, our PGAN is also able

to deal with various situations. As shown in Figure 9, although images in VRIC and VERI-Wild datasets suffer from blur, illumination and occlusion, our PGAN can still find the correct vehicles according to the prominent part regions. It means that our method is more robust to learn discriminative features that is not sensitive to multiple variants from the environment.

*2) The effectiveness of the part extraction module as the bottom-up attention:* The detected part regions play an important role in feature representation. As illustrated in the ID3, it is clear that the wrongly retrieved images from the grid attention method are different from the query image from the car lights. However, a lot of regions representing the body and the bottom of the car are concentrated, which are not the obvious differences between two vehicles. Nevertheless, with the guidance of the detected part regions, our PGAN can only focus on these candidate regions that is beneficial for focusing on useful regions as well as alleviating the bad effect from the other regions. In other words, the part extraction module helps the network learning by narrowing down the searching ranges.

*3) The effectiveness of the part attention module as the top-down attention:* Our PGAN is useful for selecting the most prominent part regions and lighten the influence of invalid and useless regions. As described in the main paper, we propose a part attention module (PAM) that is responsible for learning a soft attention weight for each part. Therefore, the important part regions are underlined by a high-attention value, while the impact of other insignificant parts is relatively suppressed. From the feature maps, we can clearly observe that our PGAN could focus on the most prominent part regions, such as the car lights in ID3, back mirrors in ID1. As shown in ID4, although there are few valid part regions that are extracted, our PGAN can still find the key information to recognize the vehicles, such as the wheel and the cat lights. On the contrary, the grid attention is largely influenced by some invalid regions that are extremely similar in different vehicles, such as the bottom of the vehicle body.

*4) The influence of the overlapped part regions:* When the area of the overlapped region is large, the attention weights of both two regions from our PAM will be tended to be large consistently if this region is prominent, and vice versa, *e.g.*, the side wind-shield glass of ID3 in Figure 8. For another situation that the area of the overlapped region is small, our PAM can provide the overall evaluation for every part region. For example, as the illustrated image in Figure 3, the attention weight of the annul services sign is large due to its unique, while the weight of the wind-shield glass is small because it has relatively less informative information.

*5) The limitation of our PGAN:* From Figure 10, it shows that our PGAN fails to distinguish: i) the extreme similar vehicles that share the same appearance; ii) the public vehicles

without unique features. It is reasonable since our PGAN depends on the discriminative information. If vehicle plates are available, our PGAN can achieve higher performance.

### F. Discussion

As shown in Table I, 16 attributes are included in our part extraction module. Although the attribute information is ignored when selecting the top-$D$ part regions based on the confidence scores, we can still analyze which part regions are prominent. We extract the attention weights from the PAM on the VeRi-776 training dataset. Note that we set $D = 8$ in this section and newer sign attribute is not appeared in VeRi-776 training dataset. From Figure 11(a), we can observe that carlight is the most frequently selected part. It makes sense that carlight appears in almost all vehicles whatever the vehicle view is. Moreover, windglass, backmirror and wheel also appear frequently, while layon (lay ornament), entry license, hungs and tissue appear rarely. Refer to Figure 11(b), it is clear that carlight plays the most important role in distinguishing different vehicles. It is interesting to see that some subtle part regions still have useful information, such as logo, hungs, entrylicense and annusigns, although these attributes appear less than windglass and backmirror that include nearly no information. The analyses show that our PGAN is effective to attend the meaningful and useful information for identifying vehicles in an interpretatable way.

Furthermore, we also report the running time. All experiments are conducted on a GeForce GTX 1080 Ti machine. Table VII shows that the IR module in PGAN can achieve the comparable speed with the baseline despite the additional PAM module and feature aggregation module. Although the SSD detector [20] is time-consuming, our PGAN is still practical in the real world.

## V. CONCLUSION

In this work, we have presented a novel Part-Guided Attention Network (PGAN) for vehicle instance retrieval (IR). First, we extract part regions of each vehicle image from an object detection model. These part regions provide a range of candidate searching area for the network learning, which is regarded as a bottom-up attention process. Then we use the proposed part attention module (PAM) to discover the prominent part regions by learning a soft attention weight for each candidate part, which is a top-down attention process. In this way, the most discriminative parts are highlighted with high-attention weights, while the opposite effects of invalid or useless parts are suppressed with relatively low weights. Furthermore, with the joint optimization of the holistic feature and the part feature, the IR performance can be further improved. Extensive experiments show the effectiveness of our method. The proposed PGAN outperforms other state-of-the-art methods by a large margin. We plan to extend the proposed method to the multi-task learning, i.e., object detection and tracking, for simultaneously improving the performance of these two tasks.
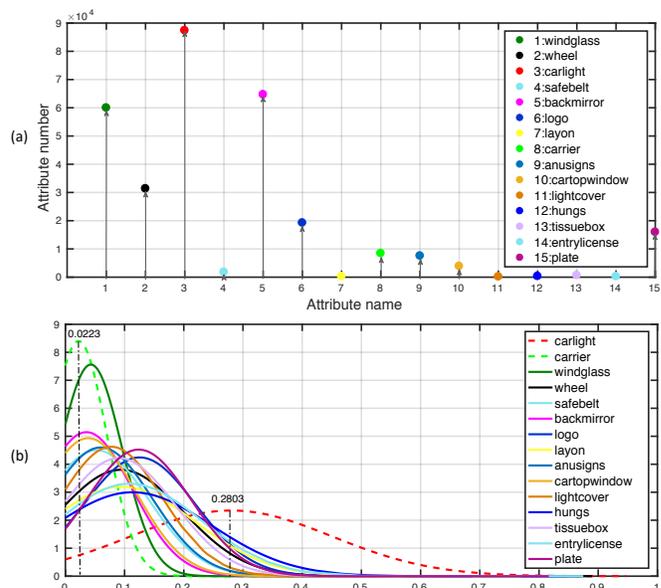


**FIG. 11:** Statistical analysis of the effectiveness of each attribute in Table I on VeRi-776 training dataset. Attribute names are denoted as the abbreviations in Table I. (a) Number statistic for each vehicle attribute. (b) The probability density function of attention weights from PAM for each attribute. The most prominent attribute is shown in red dashed line, while the least informative attribute is shown in green dashed line. Best view in color.

## REFERENCES

[1] C. Arth, C. Leistner, and H. Bischof, "Object reacquisition and tracking in large-scale smart camera networks," in *Int. Conf. Distributed Smart Cameras*, 2007, pp. 156–163.

[2] R. S. Feris, B. Siddiquie, J. Petterson, Y. Zhai, A. Datta, L. M. Brown, and S. Pankanti, "Large-scale vehicle detection, indexing, and search in urban surveillance videos," *IEEE Trans. Multimedia*, pp. 28–42, 2012.

[3] L. Yang, P. Luo, C. Change Loy, and X. Tang, "A large-scale car dataset for fine-grained categorization and verification," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2015, pp. 3973–3981.

[4] X. Liu, W. Liu, H. Ma, and H. Fu, "Large-scale vehicle re-identification in urban surveillance videos," in *Proc. Int. Conf. Multimedia. Expo.*, 2016, pp. 1–6.

[5] Z. Wang, L. Tang, X. Liu, Z. Yao, S. Yi, J. Shao, J. Yan, S. Wang, H. Li, and X. Wang, "Orientation invariant feature embedding and spatial temporal regularization for vehicle re-identification," in *Proc. IEEE Int. Conf. Comp. Vis.*, 2017, pp. 379–387.

[6] H. Guo, C. Zhao, Z. Liu, J. Wang, and H. Lu, "Learning coarse-to-fine structured feature embedding for vehicle re-identification," in *Proc. AAAI Conf. Artificial Intell.*, 2018.

[7] Z. Tang, M. Naphade, M.-Y. Liu, X. Yang, S. Birchfield, S. Wang, R. Kumar, D. Anastasiu, and J.-N. Hwang, "Cityflow: A city-scale benchmark for multi-target multi-camera vehicle tracking and re-identification," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2019, pp. 8797–8806.

[8] Y. Lou, Y. Bai, J. Liu, S. Wang, and L. Duan, "Veri-wild: A large dataset and a new method for vehicle re-identification in the wild," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2019, pp. 3235–3243.

[9] F. Zheng, X. Miao, and H. Huang, "Fast vehicle identification via ranked semantic sampling based embedding," in *Proc. Int. Joint Conf. Artificial Intell.*, 2018, pp. 3697–3703.

[10] Z. Zheng, T. Ruan, Y. Wei, Y. Yang, and T. Mei, "Vehiclenet: Learning robust visual representation for vehicle re-identification," *arXiv preprint arXiv:2004.06305*, 2020.

[11] X. Liu, W. Liu, T. Mei, and H. Ma, "A deep learning-based approach to progressive vehicle re-identification for urban surveillance," in *Proc. Eur. Conf. Comp. Vis.*, 2016, pp. 869–884.

[12] X.-S. Wei, C.-L. Zhang, L. Liu, C. Shen, and J. Wu, "Coarse-to-fine: A rnn-based hierarchical attention model for vehicle re-identification," in *Proc. Asian Conf. Comp. Vis.*, 2018, pp. 575–591.

[13] X. Liu, W. Liu, T. Mei, and H. Ma, "Provid: Progressive and multimodal vehicle reidentification for large-scale urban surveillance," *IEEE Trans. Multimedia*, pp. 645–658, 2017.

[14] Y. Shen, T. Xiao, H. Li, S. Yi, and X. Wang, "Learning deep neural networks for vehicle re-id with visual-spatio-temporal path proposals," in *Proc. IEEE Int. Conf. Comp. Vis.*, 2017, pp. 1900–1909.

[15] Y. Zhou and L. Shao, "Cross-view gan based vehicle generation for re-identification." in *Proc. British Machine Vis. Conf.*, vol. 1, 2017, pp. 1–12.

[16] ——, "Aware attentive multi-view inference for vehicle re-identification," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2018, pp. 6489–6498.

[17] Y. Lou, Y. Bai, J. Liu, S. Wang, and L.-Y. Duan, "Embedding adversarial learning for vehicle re-identification," *IEEE Trans. Image Process.*, 2019.

[18] B. He, J. Li, Y. Zhao, and Y. Tian, "Part-regularized near-duplicate vehicle re-identification," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2019, pp. 3997–4005.

[19] J. Zhu, H. Zeng, J. Huang, S. Liao, Z. Lei, C. Cai, and L. Zheng, "Vehicle re-identification using quadruple directional deep learning features," *IEEE Trans. Intell. Transportation Syst.*, 2019.

[20] Y. Zhao, C. Shen, H. Wang, and S. Chen, "Structural analysis of attributes for vehicle re-identification and retrieval," *IEEE Trans. Intell. Transportation Syst.*, 2019.

[21] H. Chen, B. Lagadec, and F. Bremond, "Partition and reunion: A two-branch neural network for vehicle re-identification," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn. Workshop*, 2019, pp. 184–192.

[22] Y. Chen, L. Jing, E. Vahdani, L. Zhang, M. He, and Y. Tian, "Multi-camera vehicle tracking and re-identification on ai city challenge 2019," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn. Workshop*, 2019, pp. 324–332.

[23] X. Liu, S. Zhang, Q. Huang, and W. Gao, "Ram: a region-aware deep model for vehicle re-identification," in *Proc. Int. Conf. Multimedia. Expo.*, 2018, pp. 1–6.

[24] H. Luo, Y. Gu, X. Liao, S. Lai, and W. Jiang, "Bag of tricks and a strong baseline for deep person re-identification," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn. Workshop*, 2019.

[25] H. Liu, Y. Tian, Y. Yang, L. Pang, and T. Huang, "Deep relative distance learning: Tell the difference between similar vehicles," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2016, pp. 2167–2175.

[26] A. Kanacı, X. Zhu, and S. Gong, "Vehicle re-identification in context," in *Proc. German. Conf. Comp. Vis.*, 2018, pp. 377–390.

[27] Y. Tang, D. Wu, Z. Jin, W. Zou, and X. Li, "Multi-modal metric learning for vehicle re-identification in traffic surveillance environment," in *Proc. IEEE Int. Conf. Image Process.*, 2017, pp. 2254–2258.

[28] X. Liu, W. Liu, H. Ma, and S. Li, "A progressive vehicle search system for video surveillance networks," in *Proc. Int. Conf. Multimedia Big Data*, 2018, pp. 1–7.

[29] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Proc. Advances in Neural Inf. Process. Syst.*, 2014, pp. 2672–2680.

[30] K. Yan, Y. Tian, Y. Wang, W. Zeng, and T. Huang, "Exploiting multi-grain ranking constraints for precisely searching visually-similar vehicles," in *Proc. IEEE Int. Conf. Comp. Vis.*, 2017, pp. 562–570.

[31] A. Sanakoyeu, V. Tschernezki, U. Buchler, and B. Ommer, "Divide and conquer the embedding space for metric learning," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2019, pp. 471–480.

[32] R. Kumar, E. Weill, F. Aghdasi, and P. Sriram, "Vehicle re-identification: an efficient baseline using triplet embedding," *arXiv preprint arXiv:1901.01015*, 2019.

[33] Y. Yuan, K. Yang, and C. Zhang, "Hard-aware deeply cascaded embedding," in *Proc. IEEE Int. Conf. Comp. Vis.*, 2017, pp. 814–823.

[34] A. Hermans, L. Beyer, and B. Leibe, "In defense of the triplet loss for person re-identification," *arXiv preprint arXiv:1703.07737*, 2017.

[35] W. Lin, Y. Li, H. Xiao, J. See, J. Zou, H. Xiong, J. Wang, and T. Mei, "Group reidentification with multigrained matching and integration," *IEEE Trans. Cybernetics.*, 2019.

[36] W. Lin, Y. Shen, J. Yan, M. Xu, J. Wu, J. Wang, and K. Lu, "Learning correspondence structures for person re-identification," *IEEE Trans. Image Process.*, pp. 2438–2453, 2017.

[37] 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 *Proc. Eur. Conf. Comp. Vis.*, 2018, pp. 480–496.

[38] X. Zhang, J. Cao, C. Shen, and M. You, "Self-training with progressive augmentation for unsupervised cross-domain person re-identification," in *Proc. IEEE Int. Conf. Comp. Vis.*, 2019.

[39] G. Wang, Y. Yuan, X. Chen, J. Li, and X. Zhou, "Learning discriminative features with multiple granularities for person re-identification," in *Proc. ACM Int. Conf. Multimedia*, 2018, pp. 274–282.

[40] S. Li, J. Li, W. Lin, and H. Tang, "Amur tiger re-identification in the wild," *arXiv preprint arXiv:1906.05586*, 2019.

[41] C. Cui, N. Sang, C. Gao, and L. Zou, "Vehicle re-identification by fusing multiple deep neural networks," in *Int. Conf. Image Process Theory, Tools and Applications*. IEEE, 2017, pp. 1–6.

[42] X. Liu, S. Zhang, X. Wang, R. Hong, and Q. Tian, "Group-group loss-based global-regional feature learning for vehicle re-identification," *IEEE Trans. Image Process.*, vol. 29, pp. 2638–2652, 2019.

[43] P. Khorramshahi, A. Kumar, N. Peri, S. S. Rambhatla, J.-C. Chen, and R. Chellappa, "A dual path model with adaptive attention for vehicle re-identification," *arXiv preprint arXiv:1905.03397*, 2019.

[44] A. Kanaci, M. Li, S. Gong, and G. Rajamanoharan, "Multi-task mutual learning for vehicle re-identification," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn. Workshop*, 2019, pp. 62–70.

[45] P. Khorramshahi, N. Peri, A. Kumar, A. Shah, and R. Chellappa, "Attention driven vehicle re-identification and unsupervised anomaly detection for traffic understanding," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn. Workshop*, 2019, pp. 239–246.

[46] S. Teng, X. Liu, S. Zhang, and Q. Huang, "Scan: Spatial and channel attention network for vehicle re-identification," in *Pacific Rim Conference on Multimedia*, 2018, pp. 350–361.

[47] S. Teng, S. Zhang, Q. Huang, and N. Sebe, "Multi-view spatial attention embedding for vehicle re-identification," *IEEE Trans. Circuits Syst. Video Technol.*, 2020.

[48] H. Guo, K. Zhu, M. Tang, and J. Wang, "Two-level attention network with multi-grain ranking loss for vehicle re-identification," *IEEE Trans. Image Process.*, vol. 28, no. 9, pp. 4328–4338, 2019.

[49] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang, "Bottom-up and top-down attention for image captioning and visual question answering," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2018, pp. 6077–6086.

[50] P. Wang, B. Jiao, L. Yang, Y. Yang, S. Zhang, W. Wei, and Y. Zhang, "Vehicle re-identification in aerial imagery: Dataset and approach," *arXiv preprint arXiv:1904.01400*, 2019.

[51] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2018, pp. 7132–7141.

[52] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2016, pp. 770–778.

[53] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *Proc. Eur. Conf. Comp. Vis.*, 2016, pp. 21–37.

[54] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2009, pp. 248–255.

[55] Z. Zhong, L. Zheng, G. Kang, S. Li, and Y. Yang, "Random erasing data augmentation," *arXiv preprint arXiv:1708.04896*, 2017.

[56] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[57] R. Chu, Y. Sun, Y. Li, Z. Liu, C. Zhang, and Y. Wei, "Vehicle re-identification with viewpoint-aware metric learning," in *Proc. IEEE Int. Conf. Comp. Vis.*, 2019.

[58] K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman, "Return of the devil in the details: Delving deep into convolutional nets," *arXiv preprint arXiv:1405.3531*, 2014.

[59] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2015, pp. 1–9.