

# Attribute-guided Feature Learning Network for Vehicle Re-identification

Huibing Wang <sup>†</sup>, Jinjia Peng <sup>†</sup>, Dongyan Chen <sup>†</sup>, Guangqi Jiang <sup>†</sup>, Tongtong Zhao <sup>‡</sup> and Xianping Fu <sup>†</sup>

**Abstract**—Vehicle re-identification (reID) plays an important role in the automatic analysis of the increasing urban surveillance videos, which has become a hot topic in recent years. However, it poses the critical but challenging problem that is caused by various viewpoints of vehicles, diversified illuminations and complicated environments. Till now, most existing vehicle reID approaches focus on learning metrics or ensemble to derive better representation, which are only take identity labels of vehicle into consideration. However, the attributes of vehicle that contain detailed descriptions are beneficial for training reID model. Hence, this paper proposes a novel Attribute-Guided Network (AGNet), which could learn global representation with the abundant attribute features in an end-to-end manner. Specially, an attribute-guided module is proposed in AGNet to generate the attribute mask which could inversely guide to select discriminative features for category classification. Besides that, in our proposed AGNet, an attribute-based label smoothing (ALS) loss is presented to better train the reID model, which can strength the distinct ability of vehicle reID model to regularize AGNet model according to the attributes. Comprehensive experimental results clearly demonstrate that our method achieves excellent performance on both VehicleID dataset and VeRi-776 dataset.

**Index Terms**—Attribute-guided Model, Attribute-based Label Smoothing Loss, Vehicle Re-identification.

## I. INTRODUCTION

VEHICLE-related researches are of vital significance for intelligent transport, which have attracted more and more attention widely. Some progresses have been made in computer vision community, such as vehicle detection [1], [2], tracking [3], [4] and classification [5], [6]. Specially, vehicle reID is different from those tasks above and aiming to search a certain vehicle across large images captured from multiple non-overlapping cameras. Meanwhile, it could automatically carry out with less time consuming and manual labor by vehicle reID, which plays an important role in modern smart surveillance systems.

Despite recent progress in vehicle reID, in particular deep learning models have made some progress [7], [8], [9], [10], it still suffers from lots of difficulties caused by various viewpoints of vehicles, complicated environments and diversified illuminations, which makes a great difference in the visual appearance of vehicles. Different from other vision tasks [11], [12], such as person ReID [13], [14], [15], [16] and fine-grained [17], [18], [19], [20], that can extract rich features from images with various poses and colors, vehicles usually have a few attributes that could be utilized to help

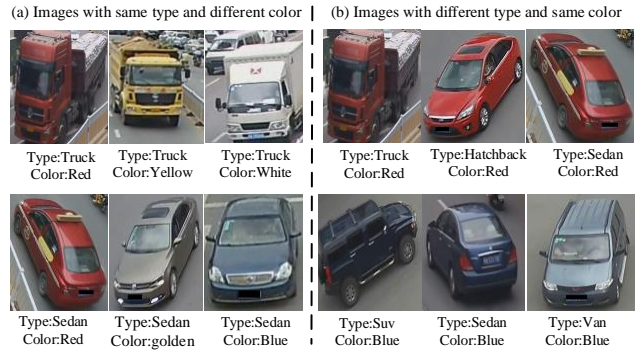


Fig. 1: Vehicle images with different attributes. (a) are the vehicle images with same type and different colors for each row. (b) are the vehicle images with same color and different types for each row.

extract distinctive features for similar vehicles. In particular, distinguishing vehicles that belong to the same or similar models can be more challenging.

To solve these problems, some existing methods focus on learning the spatio-temporal relationship between similar vehicles, such as [21], [22]. However, the spatio-temporal information are not annotated in all existing datasets, which sets a limit to us for exploring it. Besides that, attributes that could provide rich information to learn the correlation among vehicles are important auxiliary signals in vehicle reID task. Hence, several approaches explore powerful features using the attribute of vehicles, such as [23], [24]. Attributes usually describe the high-level properties, which are discriminative for the vehicles. Taking images in Fig.1 as an example, which are selected from VeRi-776 dataset. Each row in the first part shows different vehicles, which have same model and different colors, while the images in the second column have the same color and different models. It is obvious that through the attribute of vehicles, it's easy to distinguish some vehicles. Hence, the attribute is one of the important cues for vehicle reID task, which could help make good use of some local details corresponding to the attribute.

Most existing methods train the attribute features by adding special single branch for different attributes, which neglects the relationship between the attributes and identity. Hence, in our paper, different from these methods, an Attribute-Guided Network (AGNet) is proposed to select details in category feature maps that are most relevant to intrinsic attributes, which improves the discriminative ability of vehicle reID model.

<sup>†</sup>College of Information and Science Technology, Dalian Maritime University, Dalian, Liaoning, 116021, China, e-mail: {huibing.wang,jinjiapeng,chendongyan,guangqi-j,fxp}@dlmu.edu.cn  
<sup>‡</sup>zhaotongtong94@163.com

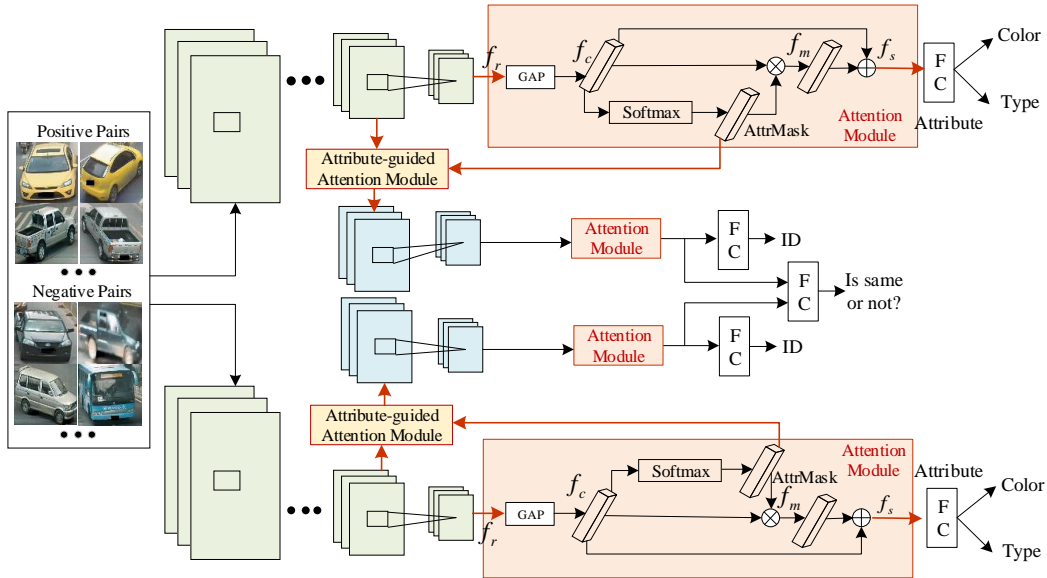


Fig. 2: Illustration of the AGNet network. AGNet is a dual-branch parallel structure. For one branch, after several resnet blocks, it is divided into two sub-branches that one is employed to predict attributes based on the image feature and another is for the category recognition. Specially, an attention module is proposed to generate the attribute mask, which helps the category branch to select better discriminative features for category recognition. Meanwhile, besides these recognition tasks, there is a verification task between two branches to improve the distinct ability of vehicle reID model.

Our intuition is that small attribute cues are usually crucial to distinguish different categories. Hence, besides training branches of attributes and category recognition respectively, an attention module is proposed in attribute branch to generate the mask that is inverse employed in the category branch, as shown in Fig.2. Owing to the generated attribute mask, an attribute-guided attention module is designed in the category branch, which helps to refine category features simultaneously to choose the important attribute features for corresponding input vehicle images. Besides that, due to which the vehicles with same color or type are more similar than others. Hence, it is improper to treat all training sets equally. And the Attribute-based Label Smoothing (ALS) loss is proposed in this paper to better regular the AGNet for training more discriminative features. Our contributions are summarized as follows:

- An Attribute-Guided Network (AGNet) is proposed to simultaneously exploits complementarity between attributes and visual appearance. In AGNet, to better select discriminative features for vehicle reID, an attention module is proposed in attribute recognition branch to generate the mask, which guides the category classification branch to select local attribute features in category feature maps.
- The Attribute-based Label Smoothing (ALS) loss is proposed to assign different weights for the training sets, which takes the attributes of vehicles into consideration. The vehicles with same color or type mean the high possibility of the same vehicles and are assigned greater weights than others.

The rest of this paper is organized as follows. In section 2, we review and discuss related works. Section 3 illustrates the proposed method in detail. Experimental results and compar-

isons on two vehicle reID datasets are discussed in section 4, followed by conclusion in section 5.

## II. RELATED WORK

In this section, existing vehicle reID works are reviewed. With the prosperity of deep learning, vehicle reID has achieved some progress in recent years. Broadly speaking, these approaches could be categorized into three classes, i.e., similarity learning, representation learning and spatio-temporal correlation learning.

A series of metric losses for deep feature embedding to achieve higher performance. In [23], coupled cluster loss was proposed to pull the positive images closer and push those negative ones far away, which minimized intra-class distance and maximized inter-class distance to train the vehicle reID network. GST loss [25] was introduced for CNNs to deal with intra-class variance in learning representation. Besides that, the mean-valued triplet loss was given to alleviate the negative impact of improper triplet sampling during training stage. MGR [26] was presented to enhance the discrimination that not only between different vehicles but also different vehicle models, which further enhanced the discriminative ability of learned features.

Apart from designing losses, some methods attempt to identify vehicles based on the visual appearance. In [27], full-fledged 3D bounding boxes vehicles were detected and then the color histograms and histograms of oriented gradients was used to extract features for vehicle reID. In [28], a ROIs-based vehicle reID method was proposed, which the ROIs' deep features were used as discriminative identifiers, encoding

the structure information of a vehicle for reID task. VAMI [29] transformed single-view feature into a global multi-view feature representation to better optimize the metric learning for training reID model. DHMVI [30] utilized the spatially concatenated convnet and LSTM bi-directional loop to learn transformations across different viewpoints of vehicles, which could infer all viewpoints' information from the only one input view. Additionally, with the popular application of Generative Adversarial Networks (GAN) in person reID, some researchers adopt GAN in vehicle reID task. CV-GAN [31] was first conducted to create the most likely image of other viewpoints to address the viewpoint variation problem. EALN [32] was proposed to automatically generate hard negative samples in the specified embedding space, which improved the capability of the network for discriminating similar vehicles.

Besides, some approaches exploit spatial and temporal information for vehicle images to improve vehicle reID performance. PROVID [33] introduced the information of license plates, visual features and spatial-temporal relations with a progressive strategy to learn similarity scores between vehicle images. Siamese-Cnn+Path-LSTM [22] model was introduced to incorporate complex spatio-temporal information for regularizing the reID results. OIFE [21] employed the log-normal distribution to model the spatio-temporal constrains in camera networks, which refined the retrieval results of vehicles.

### III. ATTENTION-GUIDED NETWORK

In this paper, we propose an end-to-end Attribute-guided feature learning network with ALS loss for vehicle reID. We firstly introduce the overview of the AGNet in section A. Then the detailed structure of AGNet is illustrated in section B, C and D. Especially, the ALS loss is explained in the section D, which can help the AGNet to select better attribute feature.

#### A. Architecture overview

In this section, a novel Attribute-Guided Network (AGNet) is proposed, which simultaneously exploits complementarity between attributes and visual appearance. The pipeline of the proposed AGNet network is shown in Fig.2. AGNet is a dual-branch parallel structure, which includes identification task and verification task. The images from the generation module are divided into positive and negative samples pairs as inputs for AGNet, which are captured by non-overlapping camera networks together with their corresponding vehicle ID. Images with the same vehicle IDs are positive sample pairs, otherwise, they are defined as negative sample pairs. The objective is to learn a discriminative representation for identifying the same vehicle and distinguishing different vehicles. Taking the upper branch as an example, given a pair of vehicle images, the AGNet extracts the feature representation  $f_g$  by several resnet blocks [34]. Subsequently, the branch is divided into two sub-branches, which one is employed to predict attributes based on the  $f_g$  and another is for the category recognition. Specially, in the attribute sub-branch, an attention module is proposed to generate an attribute mask. The mask is then as additional cues to help the category branch select better discriminative features for category recognition. Meanwhile, in AGNet, besides the

recognition task, it also contains the verification task between two parallel branches, which compares the two input features to judge that the vehicles are the same or not, and improve the distinct ability of vehicle reID model.

#### B. Attribute recognition

Attribute is a type of important auxiliary information for vehicle reID task. The same vehicles always have the same color and model. If vehicle images have different colors or models, they can't be the same vehicle. Hence, in our paper, an attribute recognition branch is designed in AGNet. As shown in Fig.2, for one branch, the input image is fed into resnet blocks to output the feature map  $f_r$  with the size of  $2048 \times 7 \times 7$ . To focus on the meaningful parts of vehicle images and neglect the background when training the feature learning model, an attention module is proposed to generate distinct features. In the attention module, after a global average pooling layer, we employ the Softmax layer to re-weight the feature maps and generate the mask, which could be computed as:

$$M = \text{Softmax}(\text{Conv}(\text{GAP}(f_r))) \quad (1)$$

where the  $\text{Conv}$  operator is  $1 \times 1$  convolution. The  $M$  is the weight matrix, which contains cues of local attribute information. Hence, we call  $M$  attribute mask (AttrMask). After obtaining the attention map  $M$ , the attended feature map could be calculated by  $f_m = f_a \otimes M$ . The operator  $\otimes$  is performed in an element-wise product. Then the attended feature map  $f_m$  is fed into the subsequent structure. At last, the prediction attribute classification is given by the fully connected layer with the cross-entropy loss that could be described as:

$$\ell_{attr} = \sum_{i=1}^k -p_i \log(q_i) \quad (2)$$

Where  $t$  is the target class and  $q_i$  is the predicted probability.  $k$  is the number of labels for attribute in the attribute recognition network. Specifically,  $q_i$  could be calculated by Softmax.  $p_i$  is the target probability.  $p_i = 0$  for all  $i \neq y$  except  $p_y = 1$ . There are vehicle color and vehicle type two attributes employed in this paper. Hence, there are two full connectional layers for different attributes.

#### C. Category recognition

To select regions in category feature maps which are most relevant to intrinsic attributes, we design an attribute-guided category recognition network to better generate discriminative features for category recognition task. Different from the attribute recognition, as shown in Fig. 2, the attribute mask (AttrMask) is generated by the attribute recognition branch, which is then fed into the category recognition branch to produce an attention map for details with intrinsic attributes for the subsequent network. As shown in Fig.3, the attribute-guided attention weights are given in section 3.2 which could

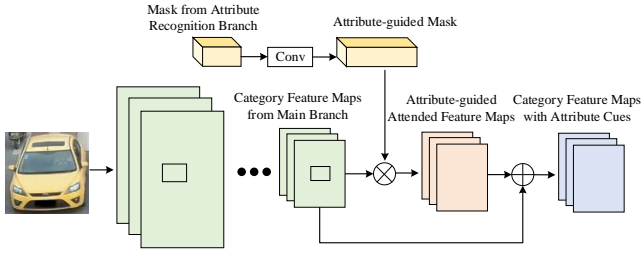


Fig. 3: Attribute-guided attention module. The attribute mask (AttrMask) is generated by the attribute recognition branch, which is then fed into the category recognition branch to produce an attention map for regional features with intrinsic attributes for the subsequent network

be described as  $AttrMask$ . The attribute features are multiplied by the attention weights and summed to produce the features  $f_{attr}$

$$f_{attr} = f_c \otimes Conv(AttrMask) \quad (3)$$

where  $f_c$  is the category feature map. The operator  $\otimes$  is performed in an element-wise product. Then the attended feature map  $f_{attr}$  with attribute information is obtained. In order to fuse the attribute features and category features, a shortcut connection architecture is introduced to embed the input of the attention network directly to its output with an element-wise sum layer, which could be described as  $f_{cs} = f_c + f_{attr}$ . In this way, both the category-involved feature maps and the attribute-involved feature maps are combined to form features  $f_{cs}$  and utilized as the input for the subsequent structure. Similar to the attribute branch, the cross-entropy loss is employed to train the category recognition, which could be described as:

$$l_{category} = \sum_{i=1}^k -p_i \log(q_i) \quad (4)$$

Where  $t$  is the target class and  $q_i$  is the predicted probability.  $k$  is the number of training identities in the dataset in the identification subnetwork.

During the test phrase, the final features are composed of features from category recognition branch and attribute classification branch, which could be described as follows:

$$f = [f_{attr}, f_{category} \times (1 - \alpha)] \quad (5)$$

where  $\alpha$  is the weight for features.  $f_{category}$  is the feature from category recognition branch.  $f_{attr}$  is from attribute recognition branch. The size of features from different branches are all  $1 \times 1 \times 4096$ .

#### D. Verification network

we add the Square Layer into the verification network. The Square Layer is denoted as  $f_v = (f_1 - f_2)^2$ , where  $f_1, f_2$  are the 4096-dim embeddings and  $f_v$  is the output tensor of the Square Layer. Then a convolutional layer and the ALS

output function are added to embed the resulting tensor  $f_v$  to a 2-dim vector. In the verification work, we treat it as a binary classification problem.



Fig. 4: The vehicle examples with different IDs but with the same vehicle color and type. It's difficult to discriminate these vehicles with similar appearance.

Traditionally, the cross-entropy loss is usually employed to train the reID model. However, The vehicles with the same type and color could have high degree of similarity, even though they have different IDs, as shown in Fig.4. Hence, these vehicles may be regard as hard training sets that help the network to effectively distinguish similarity vehicles. In the verification task, the ALS is proposed to regularize AGNet to learn better discriminative features, which could be described as follows:

$$\begin{aligned} l_{verify} &= l_{CE} + \beta l_{\gamma} \\ &= \sum_{i=1}^k -p_i \log(q_i) + \sum_{i=1}^k -\varepsilon \times p_i \log(\alpha + q_i) \end{aligned} \quad (6)$$

where  $\varepsilon$  is defined as:

$$\varepsilon = \begin{cases} \theta, & attr_1 = attr_2, id_1 \neq id_2 \\ 1 - \theta, & id_1 = id_2 \\ 0, & \text{if others,} \end{cases} \quad (7)$$

where  $id$  means the vehicle ID label and  $attr$  represents the attribute of vehicles, such as vehicle type and vehicle color.  $\theta$  is a hyper-parameter. If  $\theta$  is set to 0, Eq.(6) is equivalent to Eq.(4). Besides, the  $\alpha$  is a parameter for adjusting the possibility of the ground-truth. Compared with cross-entropy, the ALS is closer to reality by paying more additional attention to those vehicles with different vehicle IDs while are the same vehicle type and color.

#### E. Training

For the attribute-guided feature learning network, the overall objective function can be formulated as,

$$L_{\theta} = \lambda_1 l_{category} + \lambda_2 (l_{color} + l_{type}) + \lambda_3 l_{verify} \quad (8)$$

where  $\theta$  denotes the parameters in the deep model.  $l_{category}$  is the category recognition loss.  $l_{color}$  and  $l_{type}$  are the attribute classification loss.  $l_{verify}$  represents the loss of verification

task.  $\lambda_1$ ,  $\lambda_2$  and  $\lambda_3$  are the weights for corresponding loss. In our experiments,  $\lambda_1$ ,  $\lambda_2$  and  $\lambda_3$  are set to 0.5, 0.5 and 1, respectively.

#### IV. EXPERIMENTS

##### A. Datasets and evaluation metrics

In this section, some detailed analyses are given to demonstrate the effectiveness of our method. And the proposed AGNet is evaluated utilizing the mean average precision (mAP) [35] and the Cumulative Match Characteristic (CMC) curve, which are widely adopted in vehicle reID. First, we compare our AGNet with state-of-the-art methods. Then, the ablation studies are made to analyze each part in AGNet in detail. Various experiments are conducted on two popular vehicle reID datasets: VeRi-776 [33] and VehicleID [23].

1) *Datasets.*: VeRi-776 [33] is a large urban surveillance vehicle dataset for reID, which contains more than 50,000 images of 776 vehicles with identity annotations, image timestamps, camera geo-locations, vehicle color and type information. In this paper, 37,781 images of 576 vehicles are utilized as the train set and 11,579 images of 200 vehicles are employed as the test set. A subset of 1,678 images in the test set generates the query set.

VehicleID [23] is a surveillance dataset from the real-world scenario, which contains 221,763 images corresponding to 26,267 vehicles in total. From the original testing data, four subsets are extracted, which contain 800, 1,600, 2,400 and 3,200 vehicles, and are searched in different scales. During the phrase of testing, an image is randomly selected from one identity to obtain a gallery set with 800 images, and then the remaining images are all employed as probe images. Three other test sets are processed in the same way.

2) *Evaluation metrics.*: In this paper, we use CMC curve and mAP to evaluate the overall performance for all test images. Each query image in a subset of test images is given for other test images, the average precision for each query  $q$  is calculated by

$$AP(q) = \frac{\sum_{k=1}^n P(k) \times rel(k)}{N_{gt}} \quad (9)$$

Where  $P(k)$  denotes the precision at the  $k_{th}$  position of the results. The  $rel(k)$  is an indicator function equal to 1 if the  $k_{th}$  result is correctly matched or zero otherwise.  $n$  is the number of tests, and  $N_{gt}$  is the ground truths. After experimenting for each query image, the  $mAP$  will be calculated as follows:

$$mAP = \frac{\sum_{q=1}^Q AP(q)}{Q} \quad (10)$$

where  $Q$  is the number of all queries. In this paper, the vehicle images with the same ID and camera number are considered to be junk images in our evaluation of results.

##### B. Implementation details

We implement the proposed vehicle reID model in the Matconvnet [36] framework. The stochastic gradient descent is utilized with a momentum of  $\mu = 0.0005$  during the training procedure on both VeRi-776 and VehicleID. Due to the limit

of the memory of GPU, the batch size is set to 32 on VeRi-776 and VehicleID. The learning rate of the first 50 epochs is set to 0.1, and the last 25 to 0.01. The mini-batch stochastic gradient descent (SGD) is adopted to update the parameters of the network. During the phrase of training, all images with full annotations are employed. However, there are only 10086 vehicles are annotated by attribute labels, which could be used in our proposed AGNet. Hence, there are 78956 images are employed in our experiments.

TABLE I: Experimental results on VeRi-776. The mAP (%) and cumulative matching scores (%) at rank 1, 5 are listed.

Method	mAP	Rank1	Rank5
LOMO [37]	9.64	25.33	46.48
DGD [38]	17.92	50.70	67.52
GoogLeNet [39]	17.81	52.12	66.79
FACT+Plate-SNN+STR [40]	27.77	61.64	78.78
NuFACT+Plate-REC [33]	48.55	76.88	91.42
PROVID [33]	53.42	81.56	95.11
Siamese-Visual [22]	29.48	41.12	60.31
Siamese-Visual+STR [22]	40.26	54.23	74.97
Siamese-CNN+Path-LSTM [22]	58.27	83.49	90.04
OIFE+ST [21]	51.42	68.30	89.70
VAMI [29]	50.13	77.03	90.82
VAMI+ST [29]	61.32	85.92	91.84
AGNet-ASL	66.32	90.90	96.20
AGNet-ASL+STR	71.59	95.61	96.56

##### C. Comparison with the state-of-the-art methods

1) *Comparison on VeRi-776*: The results of the proposed method is compared with state-of-the-art methods on VeRi-776 dataset in Tables I II, which includes: (1) LOMO [37]; (2) DGD [38]; (3) GoogLeNet [39] (4) FACT+Plate-SNN+STR [40]; (5) NuFACT+Plate-REC [33]; (6) PROVID [33]; (7) Siamese-Visual [22]; (8) Siamese-Visual+STR [22]; (9) Siamese-CNN+Path-LSTM [22]; (10) OIFE+ST [21]; (11) VAMI [29]; (12) VAMI+ST [29]. From the Tables I II, it should be noted that the proposed method achieves the best performance among the compared with methods with rank-1 = 90.90%, mAP = 66.32% on VeRi-776, which acquires the highest mAP and rank-1 among all methods under comparisons. More details are analyzed as follows.

Firstly, the proposed AGNet obtains much better performance than those hand-crafted feature representation methods, such as LOMO [37] and DGD [38], which achieves 56.68 and 48.40 points in mAP improvements, respectively. This verifies that the features obtained from deep model are more robust than the hand-crafted feature that are severely affected by the complicated environment.

Secondly, spatio-temporal information is one of the most important cues for vehicle reID. Compared with deep learning methods that exploit vehicle reID task with spatio-temporal information, such as FACT+Plate-SNN+STR [40], PROVID [33], Siamese-Visual+STR [22], Siamese-CNN+Path-LSTM [22], OIFE+ST [21] and VAMI+ST [29], the proposed AGNet

TABLE II: Experimental results on VehicleID. The mAP (%) and cumulative matching scores (%) at Rank 1, 5 are listed.

Method	Test size = 800			Test size = 1600			Test size = 2400		
	mAP	Rank1	Rank5	mAP	Rank1	Rank5	mAP	Rank1	Rank5
BOW-SIFT [40]	-	2.81	4.23	-	3.11	5.22	-	2.11	3.76
LOMO [37]	-	19.76	32.14	-	18.95	29.46	-	15.26	25.63
DGD [38]	-	44.80	66.28	-	40.25	65.31	-	37.33	57.82
VGG+T [23]	-	40.4	61.7	-	35.4	54.6	-	31.9	50.3
VGG+CCL [23]	-	43.6	64.2	-	42.8	66.8	-	32.9	53.3
Mixed DC [23]	-	49.0	73.5	-	42.8	66.8	-	38.2	61.6
FACT [33]	-	49.53	67.96	-	44.63	64.19	-	39.91	60.49
NuFACT [33]	-	48.90	69.51	-	43.64	65.34	-	38.63	60.72
OIFE [21]	-	-	-	-	-	-	-	67.0	82.9
VAMI [29]	-	63.12	83.25	-	52.87	75.12	-	47.34	70.29
TAMR [26]	67.64	66.02	79.71	63.69	62.90	76.80	60.97	59.69	73.87
AGNet-ASL	74.05	71.15	83.78	72.08	69.23	81.41	69.66	65.74	78.28

has higher mAP, rank-1 and rank-5 than them, which demonstrates that our AGNet could extract more discriminative features without other information besides the vehicle images.

Thirdly, compared with those single modal deep learning based methods, which trains the reID model without the spatio-temporal information, the proposed AGNet shows a larger accuracy improvement. Specifically, the best single modal deep learning method VAMI [29] is also lower than the proposed AGNet in mAP, rank-1 and rank-5, which significantly shows the effectiveness of the proposed AGNet.

2) *Comparison on VehicleID*: There are 9 methods are compared with our proposed method, which are (1) LOMO [37]; (2) DGD [38]; (3) VGG+T [23]; (4) VGG+CCL [23]; (5) Mixed DC [23]; (6) FACT [33]; (6) NuFACT [33]; (7) OIFE [21]; (8) VAMI [29]; (9) TAMR [26]. Firstly, it can be observed that deep learning based methods obviously outperform traditional methods. And compared with traditional methods LOMO [37] and DGD [38], the proposed method AGNet has 68.34% and 51.39% gains on the test size with 800 vehicles, respectively. Secondly, Different VeRi-776, there is no spatio-temporal labels in VehicleID. Hence, there are no methods that consider the spatio-temporal information. All compared methods utilize the information only from vehicle images. The proposed AGNet outperforms all deep learning based methods under comparison on the test sets with different sizes on VehicleID, which obtains 71.15%, 69.23%, 65.74% in rank-1, respectively. And this also shows that our proposed AGNet could generate more distinct features for different vehicle reID datasets.

#### D. Evaluation of proposed method

To verify the effectiveness of the proposed method, some ablation experiments are conducted. The comparison results on VeRi-776 and VehicleID are presented in Table III and Table IV. Fig. 5 and Fig. 6 show the CMC curves of the compared methods.

Firstly, our descriptor is learned by multiple branches in the proposed network. Every branch is trained with sharing parameters in part of convolutional layers. We thus compare different features to test the effectiveness of our descriptor.

TABLE III: Performance of features fusion on VeRi-776. The mAP (%) and cumulative matching scores (%) at Rank 1, 5 are listed.

Descriptor	mAP	Rank1	Rank5
Only-ID	57.82	87.24	93.32
AGNet-CE-Attr	44.09	81.34	89.74
AGNet-CE-ID	64.09	89.78	94.79
AGNet-CE-All	64.78	89.86	95.17
AGNet-ASL-Attr	43.68	80.15	90.22
AGNet-ASL-ID	62.61	89.69	94.75
AGNet-ASL-All	66.32	90.90	96.20

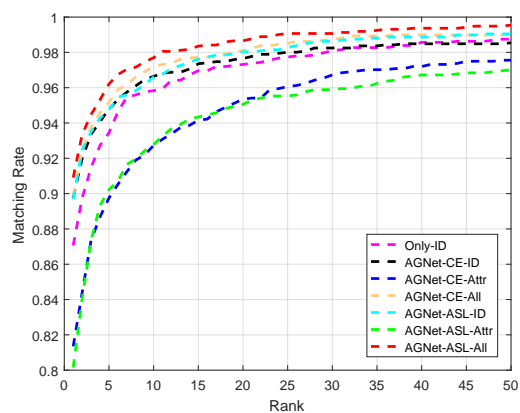


Fig. 5: The CMC results of different methods on VeRi-776.

“AGNet-ASL-All” is our proposed method that combines all features for reID task. “AGNet-ASL-ID” and “AGNet-ASL-Attr” denote the features are extracted by category branch and attribute branch, respectively. Different from “AGNet-CE-All”, “AGNet-CE-All” means that the training network has the same structure with “AGNet-ASL-All” except the training loss. “AGNet-ASL-All” is trained with the proposed ASL while “AGNet-CE-All” is CE loss. “AGNet-CE-Attr” and “AGNet-CE-ID” are similar with “AGNet-ASL-ID” and “AGNet-ASL-Attr”. From Table III and Table IV, it can be observed that “AGNet-ASL-ID” achieves higher mAP and

TABLE IV: Performance of features fusion on VehicleID. The mAP (%) and cumulative matching scores (%) at Rank 1, 5 are listed.

Descriptor	Test size = 800			Test size = 1600			Test size = 2400			Test size = 3200		
	mAP	Rank1	Rank5	mAP	Rank1	Rank5	mAP	Rank1	Rank5	mAP	Rank1	Rank5
Only-ID	70.59	67.56	80.61	68.97	65.87	79.36	65.06	61.88	75.26	63.49	60.62	72.60
AGNet-CE-Attr	67.58	63.60	81.01	62.47	58.39	76.30	59.17	54.95	73.37	56.60	52.71	69.23
AGNet-CE-ID	71.54	68.73	80.90	69.22	66.41	78.20	65.67	62.74	75.15	64.34	61.69	72.71
AGNet-CE-All	72.20	69.19	82.48	69.52	66.52	79.25	67.50	63.99	76.86	65.24	62.51	73.97
AGNet-ASL-Attr	66.77	62.69	81.29	63.22	59.20	76.78	59.59	55.53	73.06	57.09	53.34	69.29
AGNet-ASL-ID	72.39	69.66	81.55	70.68	67.99	79.40	66.95	64.09	76.46	65.63	63.11	73.42
AGNet-ASL-All	74.05	71.15	83.78	72.08	69.23	81.41	69.66	65.74	78.28	67.02	64.40	75.21

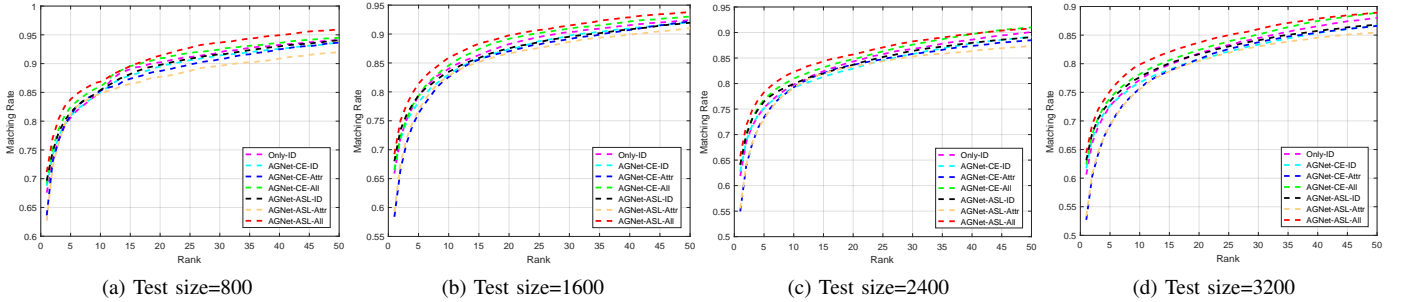


Fig. 6: The CMC curves of different methods on VehicleID.



Fig. 7: Illustration of the features obtained from our proposed network on a test set of VeRi-776 with t-SNE. Best viewed in color. In addition to vehicle images with the same ID, the images with similar colors or types are more easily concentrated together.

rank-1 than “AGNet-ASL-Attr”. This is because that “AGNet-ASL-ID” obtained from category branch which is trained with identity information has more details than the features from attribute branch. Besides that, it is worth noting that “AGNet-ASL-All” has 3.71% improvements in mAP than the “AGNet-

ASL-ID” on VeRi-776. And it also has 22.64% increases than the ‘AGNet-ASL-Attr’. The similar improvements could be observed from the comparison of “AGNet-CE-All”, “AGNet-ASL-ID” and “AGNet-ASL-Attr”. These could demonstrate that the attribute branch learns some distinctive features which are helpful for vehicle reID task.

Secondly, to demonstrate the effectiveness of the proposed attribute-guided attention model, the “Only-ID” is compared, which trains the reID model only with the ID labels with the attention structure. It is observed that, compared with “Only-ID”, our proposed “AGNet-ASL-All” achieves 8.5% improvements in mAP on VeRi-776. For VehicleID, it has 3.94%, 3.36%, 3.86%, 3.78% gains in rank-1 on test sets with different size.

Third, to demonstrate the effectiveness of the proposed ASL loss, ASL loss is compared with the “AGNet-CE-All” that means trains the reID model with cross-entropy loss and the same structure as “AGNet-ASL-All”. Compared with “AGNet-CE-All”. The significant improvements have achieved on both VeRi-776 and VehicleID, which verifies that the proposed ASL loss is more adaptive for the vehicle reID task. This is because that the ASL considers the vehicles with the same model and color should have high degree of similarity, even though they have different IDs. So it gives different weights for different conditions which could better regularize the network to train the vehicle reID model.

#### E. Qualitative analysis

To better illustrate that vehicle attribute is effective for the vehicle reID, we visualize the features of selected vehicle images in the VeRi-776 test set and project the features to



(a) VeRi-776



(b) VehicleID

Fig. 8: The retrieval results on the VehicleID and VeRi-776. (a) The results on VeRi-776. (b) The results on VehicleID. The left column shows query images while the images of right-hand side are retrieval results obtained by proposed method

2-dimensional space. Then the t-SNE [41] is employed for dimension reduction and visualization. We show all the vehicle test images in a non-occlusion form. As show in Fig.7, it can be seen that not only the same identity can be clustered together, the vehicle images which have the same color or type are also clustered together, which demonstrates that the color and type are major cues for the vehicle reID.

To further illustrate the effectiveness of the proposed framework in this paper, some results are visualized. Examples of vehicle reID results on VeRi-776 and VehicleID by our approach are shown in Fig.8. In Fig.8, for VeRi-776, the left column shows query images, while images on the right-hand are the top-11 results obtained by the algorithm. Vehicle images with red border are error results while other images are right results. For VehicleID, the left column shows query images, while images on the right-hand side are the top-5 results obtained by the algorithm. Vehicle images with green

border are right results while other images are wrong results. The number on the left-top means Vehicle ID/Camera ID. The same Vehicle ID represents the same vehicle. The Camera ID is the camera number that images are captured. From Fig.8, it could be observed that we could observe that our proposed method has high accuracy and good robustness to different viewpoints and illumination.

## V. CONCLUSION

In this paper, we propose AGNet with attribute-guided attention module which could learn global representation with the abundant attribute features in an end-to-end manner. Besides that, to better train the reID model, the ALS loss is presented, which can strength the distinct ability of vehicle reID model according to the attributes to regularize AGNet model. It can be observed from the results that compared with other existing vehicle reID methods, AGNet could achieve



competitive results. However, it is difficult to distinct some vehicles that are occluded by other vehicles. Hence, in our future studies, we would focus on exploiting the local regions with distinct features to improves the performance of reID model.

## REFERENCES

- [1] X. Hu, X. Xu, Y. Xiao, H. Chen, S. He, J. Qin, and P.-A. Heng, "Sinet: A scale-insensitive convolutional neural network for fast vehicle detection," *IEEE Transactions on Intelligent Transportation Systems*, vol. 20, no. 3, pp. 1010–1019, 2018.
- [2] W. Chu, Y. Liu, C. Shen, D. Cai, and X.-S. Hua, "Multi-task vehicle detection with region-of-interest voting," *IEEE Transactions on Image Processing*, vol. 27, no. 1, pp. 432–441, 2018.
- [3] 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8797–8806.
- [4] Y. Fang, C. Wang, W. Yao, X. Zhao, H. Zhao, and H. Zha, "On-road vehicle tracking using part-based particle filter," *IEEE Transactions on Intelligent Transportation Systems*, 2019.
- [5] Y. Wang, L. Wu, X. Lin, and J. Gao, "Multiview spectral clustering via structured low-rank matrix factorization," *IEEE transactions on neural networks and learning systems*, vol. 29, no. 10, pp. 4833–4843, 2018.
- [6] Z. Ma, D. Chang, J. Xie, Y. Ding, S. Wen, X. Li, Z. Si, and J. Guo, "Fine-grained vehicle classification with channel max pooling modified cnns," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 4, pp. 3224–3233, 2019.
- [7] 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.
- [8] B. He, J. Li, Y. Zhao, and Y. Tian, "Part-regularized near-duplicate vehicle re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3997–4005.
- [9] 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 Transactions on Intelligent Transportation Systems*, 2019.
- [10] H. Guo, C. Zhao, Z. Liu, J. Wang, and H. Lu, "Learning coarse-to-fine structured feature embedding for vehicle re-identification," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [11] Y. Wang, X. Huang, and L. Wu, "Clustering via geometric median shift over riemannian manifolds," *Information Sciences*, vol. 220, pp. 292–305, 2013.
- [12] L. Wu, Y. Wang, and L. Shao, "Cycle-consistent deep generative hashing for cross-modal retrieval," *IEEE Transactions on Image Processing*, vol. 28, no. 4, pp. 1602–1612, 2018.
- [13] L. Wu, Y. Wang, X. Li, and J. Gao, "What-and-where to match: deep spatially multiplicative integration networks for person re-identification," *Pattern Recognition*, vol. 76, pp. 727–738, 2018.
- [14] L. Wu, Y. Wang, L. Shao, and M. Wang, "3-d personvlad: Learning deep global representations for video-based person reidentification," *IEEE transactions on neural networks and learning systems*, 2019.
- [15] L. Wu, R. Hong, Y. Wang, and M. Wang, "Cross-entropy adversarial view adaptation for person re-identification," *IEEE Transactions on Circuits and Systems for Video Technology*, 2019.
- [16] D. Li, X. Chen, Z. Zhang, and K. Huang, "Learning deep context-aware features over body and latent parts for person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 384–393.
- [17] Z. Yang, T. Luo, D. Wang, Z. Hu, J. Gao, and L. Wang, "Learning to navigate for fine-grained classification," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 420–435.
- [18] L. Wu, Y. Wang, X. Li, and J. Gao, "Deep attention-based spatially recursive networks for fine-grained visual recognition," *IEEE transactions on cybernetics*, vol. 49, no. 5, pp. 1791–1802, 2018.
- [19] W. Ge, X. Lin, and Y. Yu, "Weakly supervised complementary parts models for fine-grained image classification from the bottom up," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3034–3043.
- [20] Y. Wang, W. Zhang, L. Wu, X. Lin, and X. Zhao, "Unsupervised metric fusion over multiview data by graph random walk-based cross-view diffusion," *IEEE transactions on neural networks and learning systems*, vol. 28, no. 1, pp. 57–70, 2015.
- [21] 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 *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 379–387.
- [22] 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 *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 1900–1909.
- [23] H. Liu, Y. Tian, Y. Yang, L. Pang, and T. Huang, "Deep relative distance learning: Tell the difference between similar vehicles," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2167–2175.
- [24] X. Liu, S. Zhang, Q. Huang, and W. Gao, "Ram: a region-aware deep model for vehicle re-identification," in *2018 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2018, pp. 1–6.
- [25] Y. Bai, Y. Lou, F. Gao, S. Wang, Y. Wu, and L.-Y. Duan, "Group-sensitive triplet embedding for vehicle reidentification," *IEEE Transactions on Multimedia*, vol. 20, no. 9, pp. 2385–2399, 2018.
- [26] H. Guo, K. Zhu, M. Tang, and J. Wang, "Two-level attention network with multi-grain ranking loss for vehicle re-identification," *IEEE Transactions on Image Processing*, 2019.
- [27] D. Zapletal and A. Herout, "Vehicle re-identification for automatic video traffic surveillance," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2016, pp. 25–31.
- [28] Y. Zhao, C. Shen, H. Wang, and S. Chen, "Structural analysis of attributes for vehicle re-identification and retrieval," *IEEE Transactions on Intelligent Transportation Systems*, 2019.
- [29] Y. Zhou and L. Shao, "Aware attentive multi-view inference for vehicle re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 6489–6498.
- [30] Y. Zhou, L. Liu, and L. Shao, "Vehicle re-identification by deep hidden multi-view inference," *IEEE Transactions on Image Processing*, vol. 27, no. 7, pp. 3275–3287, 2018.
- [31] Y. Zhou and L. Shao, "Cross-view gan based vehicle generation for re-identification," in *BMVC*, vol. 1, 2017, pp. 1–12.
- [32] Y. Lou, Y. Bai, J. Liu, S. Wang, and L.-Y. Duan, "Embedding adversarial learning for vehicle re-identification," *IEEE Transactions on Image Processing*, 2019.
- [33] X. Liu, W. Liu, T. Mei, and H. Ma, "Provid: Progressive and multi-modal vehicle reidentification for large-scale urban surveillance," *IEEE Transactions on Multimedia*, vol. 20, no. 3, pp. 645–658, 2017.
- [34] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *arXiv preprint arXiv:1512.03385*, 2015.
- [35] Y. Lin, L. Zheng, Z. Zheng, Y. Wu, Z. Hu, C. Yan, and Y. Yang, "Improving person re-identification by attribute and identity learning," *Pattern Recognition*, 2019.
- [36] A. Vedaldi and K. Lenc, "Matconvnet: Convolutional neural networks for matlab," in *Proceedings of the 23rd ACM international conference on Multimedia*. ACM, 2015, pp. 689–692.
- [37] 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*, 2015, pp. 2197–2206.
- [38] T. Xiao, H. Li, W. Ouyang, and X. Wang, "Learning deep feature representations with domain guided dropout for person re-identification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 1249–1258.
- [39] L. Yang, P. Luo, C. Change Loy, and X. Tang, "A large-scale car dataset for fine-grained categorization and verification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3973–3981.
- [40] X. Liu, W. Liu, T. Mei, and H. Ma, "A deep learning-based approach to progressive vehicle re-identification for urban surveillance," in *European Conference on Computer Vision*. Springer, 2016, pp. 869–884.
- [41] L. Van Der Maaten, "Accelerating t-sne using tree-based algorithms," *The Journal of Machine Learning Research*, vol. 15, no. 1, pp. 3221–3245, 2014.