

Attribute and State Guided Structural Embedding Network for Vehicle Re-identification

Hongchao Li, Chenglong Li, Aihua Zheng, Jin Tang, and Bin Luo

Abstract—Vehicle re-identification (Re-ID) is a crucial task in smart city and intelligent transportation, aiming to match vehicle images across non-overlapping surveillance camera scenarios. However, the images of different vehicles may have small visual discrepancies when they have the same/similar attributes, e.g., the same/similar color, type, and manufacturer. Meanwhile, the images from a vehicle may have large visual discrepancies with different states, e.g., different camera views, vehicle viewpoints, and capture time. In this paper, we propose an attribute and state guided structural embedding network (ASSEN) to achieve discriminative feature learning by attribute-based enhancement and state-based weakening for vehicle Re-ID. First, we propose an attribute-based enhancement and expanding module to enhance the discrimination of vehicle features through identity-related attribute information, and we design an attribute-based expanding loss to increase the feature gap between different vehicles. Second, we design a state-based weakening and shrinking module, which not only weakens the state information that interferes with identification but also reduces the intra-class feature gap by a state-based shrinking loss. Third, we propose a global structural embedding module that exploits the attribute information and state information to explore hierarchical relationships between vehicle features, then we use these relationships for feature embedding to learn more robust vehicle features. Extensive experiments on benchmark datasets VeRi-776, VehicleID, and VERI-Wild demonstrate the superior performance and generalization of the proposed method against state-of-the-art vehicle Re-ID methods. The code is available at https://github.com/ttaalle/fast_assem.

Index Terms—Vehicle re-identification, Attribute-based enhancement, State-based weakening, Global structural embedding.

I. INTRODUCTION

VEHICLE Re-identification (Re-ID) aims to identify vehicle images from the gallery images captured from non-overlapping surveillance cameras that share the same identity as the given probe vehicle. It is an active and challenging task and has drawn much attention due to its wide applications in social security, smart city, and intelligent transportation. The blossom of Deep Convolutional Neural Network (DCNN) has witnessed recent breakthroughs in vehicle Re-ID. However, it still faces two severe challenges. 1) The large intra-class discrepancy among the same vehicle images under different

This research is supported in part by the National Natural Science Foundation of China (61976002, 61976003, 61860206004), and the Natural Science Foundation of Anhui Higher Education Institutions of China (KJ2019A0033).

Hongchao Li, Jin Tang and Bin Luo are with Anhui Provincial Key Laboratory of Multimodal Cognitive Computation, School of Computer Science and Technology, Anhui University, Hefei, 230601, China.

Chenglong Li and Aihua Zheng are with Information Materials and Intelligent Sensing Laboratory of Anhui Province, Anhui Provincial Key Laboratory of Multimodal Cognitive Computation, School of Artificial Intelligence, Anhui University, Hefei, 230601, China.

Corresponding author: Aihua Zheng. Email: ahzheng214@foxmail.com.

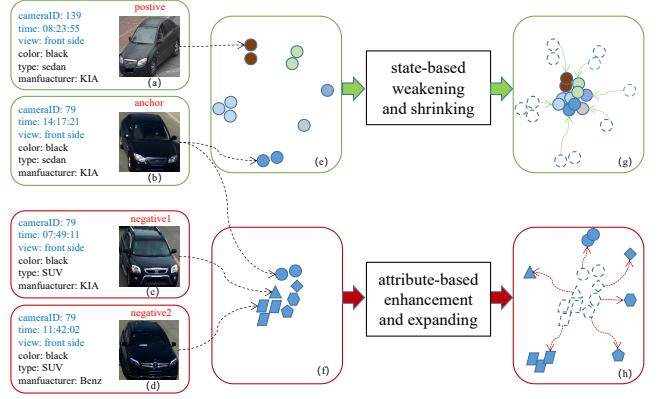


Fig. 1. Illustration of our attribute-based enhancement and state-based weakening framework for vehicle Re-ID. Different colors represent different state information while different markers denote different IDs. In the input image space, vehicle “b, c, d” are more like a same identity vehicle than vehicle “a, b”. Our attribute-based enhancement and expanding module is designed to expand the attribute distribution and re-weight attribute features to the vehicle features to enhance the inter-class difference. For example, the feature distance between manufacture KIA and manufacture Benz is enlarged to force the feature distance between vehicle “b” and vehicle “d” to be greater. In the same way, our state-based weakening and shrinking module is designed to shrink the state distribution and re-weight state features to the vehicle features to weaken the intra-class difference. For example, the feature distance between camera 139 and camera 79 is reduced to force the feature distance between vehicle “b” and vehicle “a” to be smaller. Therefore our attribute-based enhancement and state-based weakening framework can cluster images from the same vehicle compactly and enhance the discrimination between different vehicles.

states, e.g., different camera views, vehicle viewpoints, and capture time as shown in Fig. 1 (a) and (b). 2) The small inter-class discrepancy among different vehicles especially when sharing the same/similar attributes, e.g., the same/similar color, type, and manufacturer as shown in Fig. 1 (b), (c) and (d).

Recent efforts have provided various solutions while handling the above challenges. Representative approaches fall into five categories: 1) Global feature based methods [1-6], which aim to extract the global hand-crafted/deep features of vehicle images by specific metric learning methods. However, global feature based methods are generally hard to capture the intra-class discrepancy and inter-class similarity since only the appearance of vehicle images are considered. 2) Path-based methods [7-9] usually adopt spatial-temporal information to remove unreasonable vehicles for refining the retrieval results in the inference stage. However, the appearance changes of the vehicle due to spatial-temporal changes are ignored in the learning stage of vehicle features. 3) Viewpoint-based methods [10-12], which aim to handle viewpoint changes

and learn multi-view features via metric learning for vehicle Re-ID. Meanwhile, some viewpoint-based methods [13, 14] generate hard negative cross-view and same-view images for more robust training with a Generative Adversarial Network (GAN) [15]. Although these viewpoint-based methods significantly reduce the intra-class difference, they ignore the intrinsic state factors of vehicles (*e.g.*, camera views and capture time) and overlook the challenge of the subtle inter-class discrepancy. 4) Local information enhancement methods [16-21] usually provide some stable discriminative cues to increase the inter-class discrepancy for vehicle Re-ID. However, local region extraction models usually require a large amount of annotated data which are time and labor consuming. Furthermore, the forthcoming Re-ID model may be sensitive to the inaccurate part extraction. 5) Attribute-based methods, which use attribute labels to constrain identity features [22], or directly concatenating [23] or summing weighted [24, 25] identity features and attribute features to boost the Re-ID task. Generally speaking, path-based and viewpoint-based methods devote to reduce the impact of identity-unrelated information on vehicle Re-ID, while local information enhancement and attribute-based methods aim to enhance the identity-related information to improve the Re-ID task. In this work, we argue to simultaneously enhance the identity-related and weaken the identity-unrelated information.

In vehicle Re-ID, first, the images of different vehicles with similar attributes share a similar visual appearance (as shown in Fig. 1 (b, c, d)). This results in smaller distances between different vehicles in the feature space (as shown in Fig. 1 (f)), which is the key reason of inter-class similarity in vehicle Re-ID. Therefore, we argue that the feature gap of different vehicle images can be increased by enhancing their identity-related attribute information during feature learning. This is known as knowledge embedding [26] which has been commonly employed in many other computer vision problems [24, 27, 28]. Specifically, we propose an **attribute-based enhancement and expanding module to expand the attribute distribution and re-weight attribute features to the vehicle features to enhance the inter-class difference**. As shown in Fig. 1 (b, d), we enlarge the feature distance between manufacturers “KIA” and “Benz” to force larger feature distance between the two vehicles. Second, the images of the same vehicle (as shown in Fig. 1 (a, b)) under the different states generally present different visual appearance. This results in larger distances between the same vehicle images in the feature space (as shown in Fig. 1 (e)), which is the key reason of intra-class discrepancy in vehicle Re-ID. In the same way, we further argue to decrease the feature gaps between that the images of the same vehicle via state-based weakening during feature learning. Specifically, we propose a **state-based weakening and shrinking module to shrink the state distribution and re-weight state features to the vehicle features to weaken the intra-class difference**. As shown in Fig. 1 (a, b), we reduce the feature distance between camera 139 and camera 79 to encourage the smaller feature distance between the two images. By enforcing the attribute-based enhancement and state-based weakening constraints, identity-related attribute clues will be enhanced while the identity-

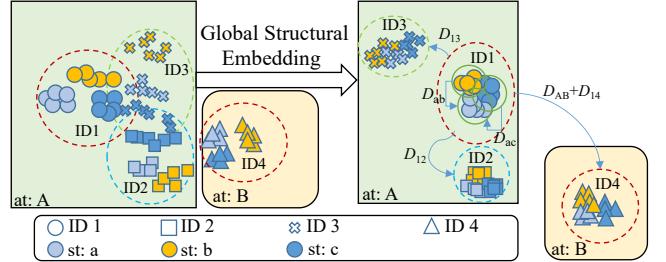


Fig. 2. Illustration of the global structural embedding module for vehicle Re-ID. These points denote the feature embeddings on 60 images from 4 identities in the VeRi-776 testing set. In the input image feature space, vehicle ID1 and ID3 share the same attribute present large overlap due to the large inter-class similarity. Meanwhile, the vehicle ID1 in different states appears sparse feature distribution due to the large intra-class discrepancy. After global structural embedding, images of the same vehicle have been compactly aggregated and the discrimination between different vehicles has been enhanced guided by their state discrepancy D_{ab} , D_{ac} , instance discrepancy D_{12} , D_{13} , D_{14} , and attribute discrepancy D_{AB} .

independent state factors will be weakened in the vehicle features.

Additionally, the deep metric learning methods, which utilize distance metric loss (*e.g.*, contrastive loss [29] and triplet loss [5]) rather than cross-entropy loss [29], aim to learn a deep feature embedding space by enforcing the distance between positive pairs smaller than that of negative pairs during learning. However, most exiting metric learning methods only focus on the appearance, which ignores the hierarchical structural relationships caused by the states and attributes. Concretely, different vehicle instances with similar appearance can be further distinguished based on their attribute diversity. Therefore it is effective to consider this relationship to increase the inter-class feature distance as shown in Fig. 2. Meanwhile, the images of the same vehicle instance with large appearance changes can be further recognized by their state information. Therefore, it is useful to decrease the intra-class feature distance between easy and hard positive samples as shown in Fig. 2. Herein, we propose a **global structural embedding module for all vehicle images to cluster images from the same vehicle compactly and enhance the discrimination between different vehicles guided by their state discrepancy, instance discrepancy and attribute discrepancy**.

In this work, we propose an attribute and state guided structural embedding network (ASSEN) towards enlarging the distance of vehicle inter-class features by all available vehicle attributes and reducing the distance of vehicle intra-class features by all available vehicle states. First, we construct an attribute-based enhancement and expanding module to obtain vehicle features enhanced by multiple attributes and design an attribute-based expanding loss to increase the vehicle inter-class gap. Then, we propose a state-based weakening and shrinking module to force the learned vehicle features to weaken state information that interferes with identity and design a state-based shrinking loss to reduce the vehicle intra-class gap. The above two modules encourage our ASSEN to be more focused on its identity-related information rather than identity-unrelated information. Finally, we construct a global structural embedding module to encourage vehicle features

159 to have a global structure related to instance discrepancy,
 160 state discrepancy and attribute discrepancy, which can bring
 161 hierarchical relationships into the feature embedding to obtain
 162 more discriminative vehicle features.

163 The contributions of this paper can be summarized as
 164 follows.

- 165 • We design an attribute-based enhancement and expanding
 166 module to obtain vehicle features enhanced by multiple
 167 attributes. Compare with previous attribute-based Re-ID
 168 methods, which use attribute labels to constrain identity
 169 features [22], or directly concatenating [23] or summing
 170 weighted [24, 25] identity features and attribute features
 171 to boost the Re-ID task. Our method utilizes the response
 172 relationship between attribute feature and identity feature
 173 to highlight the foreground area of the vehicle and expand
 174 the subtle differences between the same attribute.
- 175 • We propose a state-based weakening and shrinking mod-
 176 ule to weaken the influence of state information and re-
 177 duce the state change of the same vehicle. Different from
 178 previous work, we further divide the common attribute
 179 information into identity-related information and identity-
 180 unrelated information. Our key idea is to simultaneously
 181 enhance the identity-related and weaken the identity-
 182 unrelated information in a unified framework.
- 183 • We propose a global structural embedding module to
 184 consider hierarchical relationships related to instance
 185 discrepancy, state discrepancy and attribute discrepancy
 186 in the feature embedding to learn larger weights for hard
 187 negative (positive) samples with similar attributes (shar-
 188 ing different states). Existing metric learning methods
 189 only consider a small number of samples, or equally
 190 treat all samples. Our method adaptively assigns different
 191 weights to each sample pair.
- 192 • Comprehensive experiments on three large-scale vehicle
 193 Re-ID benchmark datasets with or without state and
 194 attribute information confirm the effectiveness and gen-
 195 eralization of the proposed model.

196 II. RELATED WORK

197 We briefly review the related works in the following two
 198 folds, i.e., vehicle Re-ID and deep metric learning.

199 A. Vehicle Re-identification

200 Due to wide applications in video surveillance and social
 201 security, the vehicle Re-ID task has gained more and more
 202 attention in recent years. Liu *et al.* [4] present a deep relative
 203 distance learning method to extract both model and instance
 204 differences. Features from the model and instance are concate-
 205 nated to learn the final vehicle feature with vehicle labels. Liu
 206 *et al.* [30] fuse color, texture, and deep features for vehicle Re-
 207 ID. They show that deep features outperform the others and
 208 feature fusion improves the Re-ID performance. Yan *et al.* [31]
 209 model the relationship of vehicle images as a multi-grain list
 210 to discriminate appearance-similar vehicles. By introducing
 211 multi-grain relationships, they force the deep model to learn
 212 the more discriminative feature between different grains over
 213 many images. Liu *et al.* [7] propose a spatial-temporal relation

214 model to re-rank vehicles to further improve the final results
 215 of vehicle Re-ID. Shen *et al.* [8] investigate spatial-temporal
 216 association for effectively regularizing vehicle Re-ID results.
 217 The spatial-temporal information along the candidate path is
 218 effectively incorporated to estimate the validness confidence
 219 of the path. Wang *et al.* [32] embed the spatial-temporal
 220 regularization into the orientation invariant module for ve-
 221 hicle Re-ID. With spatial-temporal regularization, the log-
 222 normal distribution is adopted to model the spatial-temporal
 223 constraints and the retrieval results can be refined.

224 Different from the above global feature based methods and
 225 path-based methods, He *et al.* [17] investigate vehicle local
 226 regions to learn part-regularized features for vehicle Re-ID.
 227 Khorramshahi *et al.* [18] present a dual-path adaptive attention
 228 model, to capture key-points related to parts for vehicle Re-ID.
 229 Meng *et al.* [19] propose a part perspective transformation on
 230 feature space to transform the deformed region to a unified
 231 perspective. Liu *et al.* [21] adopt the graph convolutional
 232 networks (GCNs) [33] to model the correlation among parts
 233 for vehicle Re-ID. However, the part-based approaches need
 234 additional part annotations, which takes extra costs. A part
 235 prediction network is also needed, which involves more training
 236 procedures and complicates the feature extraction model. In
 237 addition, identity-related part information is easily disturbed
 238 by identity-unrelated information, such as vehicle viewpoints.

239 To handle the viewpoint variation issue in vehicle Re-ID,
 240 Sochor *et al.* [34] learn a 3D orientation vector embedded
 241 into the feature map for vehicle recognition. They show that
 242 orientation information can decrease classification error and
 243 boost verification average precision. Zhou *et al.* [35] generate
 244 the opposite side features to handle the viewpoint problem.
 245 Zhou *et al.* [13] propose a viewpoint aware network that
 246 integrates features from viewpoint-based feature extractors
 247 with a GAN to create cross-view features for vehicle Re-
 248 ID. Zhou *et al.* [10] exploit the great advantages of DCNN
 249 and Long Short-Term Memory (LSTM) [36] to learn trans-
 250 formations across different viewpoints of vehicles. Lou *et*
 251 *al.* [14] propose an embedding adversarial learning network
 252 (EALN) to generate hard negative cross-view and same-view
 253 images for more robust training in vehicle Re-ID. Jin *et*
 254 *al.* [11] propose an Uncertainty-aware Multi-shot Teacher-
 255 Student (UMTS) Network to exploit the comprehensive in-
 256 formation of multi-view of the same vehicle for effective
 257 vehicle Re-ID. However, it is difficult to resolve the challenge
 258 of vehicle inter-class similarity with these viewpoint learning
 259 methods. Most of existing methods only reduce intra-class
 260 discrepancy by state (spatial-temporal, viewpoint) informa-
 261 tion or increase inter-class discrepancy by part information
 262 individually, while ignoring the global structural relationship
 263 related to states and attributes. We propose an attribute-based
 264 enhancement and state-based weakening framework, aiming to
 265 explore the global structural relationship to increase the inter-
 266 class discrepancy and simultaneously reduce the intra-class
 267 discrepancy.

268 B. Attribute-based Re-identification

269 Recent works in person Re-ID [24, 37-39] adopt person
 270 attributes, such as gender and hair length, as important traits

271 to recognize pedestrians. Khamis *et al.* [37] jointly learn a dis-
 272 *273* 274 discriminative projection to a joint appearance-attribute subspace,
 275 by effectively leveraging the interaction between attributes
 276 and appearance for person Re-ID. Su *et al.* [38] propose
 277 a weakly supervised multi-type attribute learning framework
 278 based on the triplet loss by pre-training the attributes predictor
 279 on independent data. Lin *et al.* [24] simultaneously learn Re-
 280 ID embedding and pedestrian attributes, by sharing the same
 281 backbone and owning classification FC layers respectively.
 282 Sun *et al.* [39] train two different models for attribute and
 283 identity recognition tasks and concatenate two branches to one
 284 identity vector for Re-ID.

285 In vehicle Re-ID, Zheng *et al.* [25] propose a deep network
 286 architecture guided by meaningful attributes, including vehicle
 287 viewpoints, types, and colors, for vehicle Re-ID. Zhao *et*
 288 *al.* [23] collect a new vehicle dataset with 21 classes of
 289 structural attributes and proposed a region of interest (ROIs-
 290 based) vehicle Re-ID method. Qian *et al.* [22] propose a two-
 291 branch stripe-based and attribute-aware deep convolutional
 292 neural network (SAN) to learn the efficient feature embedding
 293 for vehicle Re-ID task. However, both attributes and vehicle
 294 images face challenges caused by appearance changes. Dif-
 295 ferent from previous work, we further divide the common
 296 attribute information into identity-related information (named
 297 attributes, such as color and type) and identity-unrelated
 298 information (named states, such as viewpoint and camera).
 299 Our key idea is to simultaneously enhance the identity-related
 and weaken the identity-unrelated information in a unified
 framework.

300 C. Deep Metric Learning

301 Deep metric learning aims to learn a deep feature embed-
 302 ding space, in which the samples of a same class are close to
 303 each other and the samples of different classes are far away.
 304 There are two fundamental types of loss functions for deep
 305 metric learning, *i.e.*, the contrastive loss [29] and the triplet
 306 loss [5], which have been widely used in both person and
 307 vehicle Re-ID [40-43]. However, the conventional contrastive
 308 loss or triplet loss based deep metric learning often suffers
 309 from slow convergence and poor local optima, since only a
 310 few samples are considered in each training batch.

311 There emerge many advances in more robust deep metric
 312 learning recently. Chen *et al.* [44] design a quadruplet loss to
 313 enforce a larger inter-class variation and a smaller intra-class
 314 variation compared to the triplet loss. Sohn *et al.* [45] propose
 315 an n -pair loss to generalize triplet loss by allowing joint com-
 316 parison among more than one negative example. He *et al.* [46]
 317 propose a triplet-center loss to learn a center for each class to
 318 enhance the discriminative power of the features. Ustinova *et*
 319 *al.* [47] propose a listwise loss to estimate two distributions
 320 of similarities between positive (matching) and negative (non-
 321 matching) pairs. Wang *et al.* [48] propose a ranked list loss to
 322 rank all positive points before the negative points and force a
 323 margin between them. Liu *et al.* [49] propose a Group-Group
 324 Loss (GGL) to accelerate the intra-group and inter-group
 325 feature learning and promote the discriminative ability. Wu
 326 propose [50] a margin loss that relaxes unnecessary constraints

327 from traditional contrastive loss and enjoys the flexibility of
 328 the triplet loss. However, all the images in positive/negative
 329 pairs are treated equally in existing metric learning approaches,
 330 which ignore the hierarchical relationships between vehicles.
 331 In this paper, we propose a global structural embedding
 332 loss to cluster images from the same vehicle compactly and
 333 enhance the discrimination between different vehicles guided
 334 by their state discrepancy, instance discrepancy and attribute
 335 discrepancy.

336 III. METHOD

337 To reduce the intra-class distance of vehicles and increase
 338 the inter-class distance of vehicles, we propose an Attribute
 339 and State guided Structural Embedding Network (ASSEN). It
 340 mainly consists of three modules: attribute-based enhancement
 341 and expanding, state-based weakening and shrinking, global
 342 structural embedding.

343 A. Baseline

344 In this work, our goal is to use the easily obtainable state
 345 and attribute information in real-world scenes together with
 346 the vehicle ID information to learn the discriminative vehicle
 347 identity features. Formally, we denote a vehicle input as
 348 $I = \{(x, y^{id}, y_i^{at}|_{i=1}^M, y_j^{st}|_{j=1}^N)\}$, where x and y^{id} denote the
 349 input training vehicle image and its associated vehicle identity
 350 label. y_i^{at} and y_j^{st} denote the i -th attribute label and the j -th state
 351 label of the image x respectively. M and N are the
 352 numbers of attribute and state respectively. **It's worth noting**
 353 **that, attribute/state labels are not essential during the training**
 354 **since we can use the pre-trained attribute/state branches when**
 355 **the attribute/state labels are absent.**

356 Given a deep backbone network $F(\cdot; \theta)$ with the input image
 357 $x \in R^{W \times H \times C}$, where θ represents the learnable parameters
 358 of the network. We adopt ResNet-50 [51] without final down-
 359 sampling as the backbone model followed by the state-of-the-
 360 art vehicle Re-ID methods, such as UMTS [11], PPT [19],
 361 FastReID [52], which is also a common setting in person Re-
 362 ID methods after PCB [53]. The corresponding vehicle feature
 363 tensor encoded by the network is denoted as $T = F(x; \theta) \in$
 364 $R^{w \times h \times c}$. Then the identity classification (cross-entropy) loss
 365 \mathcal{L}_{ce}^{id} is in the form of,

$$\mathcal{L}_{ce}^{id} = -y^{id} \log(FC(GAP(T))), \quad (1)$$

366 where GAP denotes a global average pooling operation, and
 367 FC denotes a Full Connected layer that predicts the result of
 368 classification. In this paper, we regard ResNet-50 with \mathcal{L}_{ce}^{id} as
 369 our baseline.

370 B. Attribute-based Enhancement and Expanding (AEE) Mod- 371 ule

372 Different from the previous attribute-based Re-ID meth-
 373 ods [22-25], which boost Re-ID tasks by concatenating or
 374 weighting attribute features. On the one hand, our AEE module
 375 hopes to enhance the image area corresponding to the attribute
 376 to improve the feature learning ability of a single sample.
 377 On the other hand, our AEE module hopes to expand the

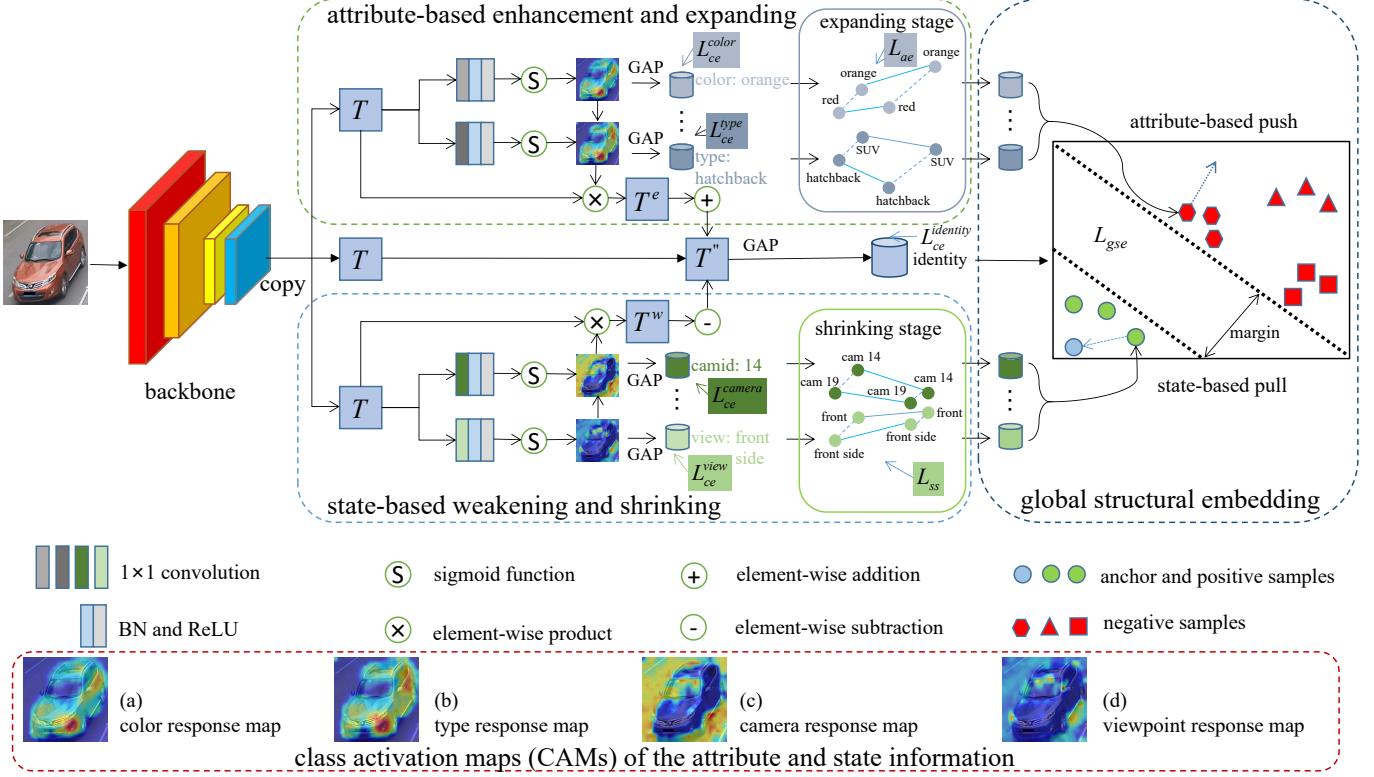


Fig. 3. Pipeline of Attribute and State guided Structural Embedding Network (ASSEN). Given the image x , we first extract the corresponding vehicle feature tensor T via the backbone. Next, we transform the feature tensor T into the attribute-based enhancement and expanding (AEE) module to obtain the enhanced feature tensor T^e . The AEE module is constrained by the attribute-related cross-entropy loss \mathcal{L}_{ce}^{at} and attribute-based expanding loss \mathcal{L}_{ae} . Then, we transform the feature tensor T into the state-based weakening and shrinking (SWS) module to obtain the weakened feature tensor T^w . The SWS module is constrained by the state-related cross-entropy loss \mathcal{L}_{ce}^{st} and state-based shrinking loss \mathcal{L}_{ss} . Followed by the combination T'' of the feature tensor T , the enhanced feature tensor T^e and the weakened feature tensor T^w to increase the identity-related information and simultaneously reduce the information that interferes with identity. Finally, the global structural embedding (GSE) module embeds instance discrepancy, attribute discrepancy and state discrepancy to obtain more discriminative vehicle features by a hierarchical structure. Note that ASSEN does not require attribute/state labels during the test. Furthermore, attribute/state labels are not essential during the training since we can use the pre-trained attribute/state branches when the attribute/state labels are absent.

378 distribution of attributes to increase the inter-class distance
379 of samples within a batch.

380 To obtain the attribute information of the vehicle, we
381 transform the vehicle feature tensor into the vehicle attribute
382 feature tensor. The i -th attribute feature tensor $T_i^{at} \in R^{w \times h \times c}$
383 can be formulated as:

$$T_i^{at}|_{i=1}^M = \text{ReLU}(BN(conv_i^{1 \times 1}(T))), \quad (2)$$

384 where $conv_i^{1 \times 1}$ denotes 1×1 convolutional operation about
385 the i -th attribute, BN denotes a Bath Normalize operation, and
386 ReLU denotes Rectified Linear Unit. $conv + BN + ReLU$
387 composes of a common convolutional block in DCNN.

388 Then the attribute classification loss \mathcal{L}_{ce}^{st} is in the form of,

$$\mathcal{L}_{ce}^{at} = - \sum_{i=1}^M y_i^{at} \log(FC(GAP(T_i^{at}))), \quad (3)$$

389 where M is the number of attributes, y_i^{at} denotes the i -th
390 attribute label of the image x .

391 The attribute tensor will be constrained by the cross-entropy
392 loss and the ground-truth attribute label. Our purpose here is
393 to use attribute labels to enable the output of vehicle features

to be guided by multiple attributes. The enhanced tensor can
be expressed as:

$$T^e = \frac{1}{M} \sum_{i=1}^M T \odot \text{Sigmoid}(T_i^{at}), \quad (4)$$

396 where $T^e \in R^{w \times h \times c}$ denotes the attribute enhanced tensor, the
397 **Sigmoid** function is used to control the value range of T_i^{at} in
398 the interval $[0, 1]$, and \odot is the element-wise product. Similar
399 to attention-based Re-ID methods [13, 54, 55], which aims
400 to re-weight the convolutional output of DCNN as a feature
401 combination. However, most of existing attention-based Re-ID
402 methods lack the guidance of identity-related annotations and
403 therefore fail to take advantage of the relationship among the
404 identity, color, and type of the same vehicle. We argue that
405 this intrinsic identity-related information is crucial in vehicle
406 Re-ID.

407 The overall attribute-based enhancement procedure can be
408 formulated as:

$$T' = T + \beta_1 T^e, \quad (5)$$

409 where T' denotes the vehicle feature tensor after attribute-
410 based enhancement operation, $\beta_1 = 0.05$ is a hyperparameter
411 used to balance the original feature and the enhanced feature.

412 We add the class activation maps (CAMs) [56] of the attribute
 413 (color and type) information, as shown in Fig. 3 (a, b). The
 414 color response map and type response map mainly respond
 415 to the foreground area related to the vehicle identity, which
 416 means that Eq. (5) tends to highlight the foreground area of
 417 the vehicle image.

418 In addition to attribute-based enhancement, we further pro-
 419 pose an attribute expanding operation to increase the inter-
 420 class attribute discrepancy. The global average pooling (GAP)
 421 is used to transfer the i -th attribute tensor $T_i^{at} \in R^{w \times h \times c}$
 422 into the i -th attribute feature vector $f_i^{at} \in R^c$. First, we calculate
 423 the i -th attribute standard deviation, which can be formulated
 424 as: $D_i^{at} = std(f_i^{at}, \bar{f}_i^{at})$, where $f_i^{at} = GAP(T_i^{at})$ denotes the
 425 i -th attribute feature vector about each image in a batch, \bar{f}_i^{at}
 426 denotes the i -th attribute mean vector about the whole batch-
 427 size. Our purpose here is to expand the feature distribution
 428 of the attribute under the premise of attribute classification,
 429 thereby increasing the inter-class attribute discrepancy. The
 430 attribute-based expanding loss can be formulated as:

$$\mathcal{L}_{ae} = \mathcal{L}_{ce}^{at} + \frac{1}{M} \sum_{i=1}^M \frac{1}{1 + exp(D_i^{at})}. \quad (6)$$

431 If there exist two samples that share the same color (or type)
 432 in a batch, their color (or type) feature distance will become
 433 larger under the premise of classification.

434 C. State-based Weakening and Shrinking (SWS) Module

435 Although attribute-based enhancement and expanding
 436 (AEE) module can enhance the inter-class difference by
 437 vehicle identity-related attribute information. These identity-
 438 related attribute information may be indistinguishable due to
 439 diverse state (e.g., camera views, vehicle viewpoints, capture
 440 time) changes. We argue that merely enhancing identity-related
 441 information is not sufficient for Re-ID, weakening the state
 442 information that interferes with identification is also crucial
 443 for vehicle Re-ID. Herein, we further consider weakening state
 444 information to reduce the intra-class feature gap for vehicle
 445 Re-ID.

446 The j -th state feature tensor $T_j^{st} \in R^{w \times h \times c}$ can be
 447 formulated as:

$$T_j^{st}|_{i=1}^N = ReLU(BN(conv_j^{1 \times 1}(T))), \quad (7)$$

448 where $conv_j^{1 \times 1}$ denotes 1×1 convolutional operation about
 449 the j -th state. Then the state classification loss \mathcal{L}_{ce}^{st} is in the
 450 form of,

$$\mathcal{L}_{ce}^{st} = - \sum_{j=1}^N y_j^{st} log(FC(GAP(T_j^{st}))), \quad (8)$$

451 where N is the number of states, and y_j^{st} denotes the j -th state
 452 label of the image x .

453 The state tensor will be constrained by the cross-entropy
 454 loss and the ground-truth state labels. Our goal is to make the
 455 learned vehicle feature tensor T alleviate the interference of
 456 multiple states as much as possible. The weakened tensor can
 457 be expressed as:

$$T^w = \frac{1}{N} \sum_{i=1}^N T \odot Sigmoid(T_j^{st}), \quad (9)$$

458 where $T^w \in R^{w \times h \times c}$ denotes the state weakened tensor, the
 459 *Sigmoid* function is used to control the value range of T_j^{st} to
 460 $[0, 1]$, and \odot is the element-wise product.

461 The overall state-based weakening procedure can be formu-
 462 lated as:

$$T'' = T' - \beta_2 T^w, \quad (10)$$

463 where T'' denotes the vehicle feature tensor after state-based
 464 weakening operation, T' denotes the vehicle feature tensor after
 465 attribute-based enhancement operation, $\beta_2 = 0.05$ is a
 466 hyperparameter used to balance the original feature and the
 467 state weakened feature. We add the class activation maps
 468 (CAMs) [56] of the state (camera and viewpoint) information
 469 as shown in Fig. 3 (c, d). The camera response map and
 470 viewpoint response map mainly respond to the background
 471 area of the vehicle image. Therefore Eq. (10) can suppress
 472 the background area of the vehicle image.

473 In addition to designing a state-based weakening procedure,
 474 we also added a state-based shrinking operation to reduce
 475 the intra-class state discrepancy. The global average pooling
 476 (GAP) is used to transfer the j -th state tensor $T_j^{st} \in R^{w \times h \times c}$
 477 into the j -th state feature vector $f_j^{st} \in R^c$. First, we calculate
 478 the j -th state standard deviation, which can be formulated as:
 479 $D_j^{st} = std(f_j^{st}, \bar{f}_j^{st})$, where $f_j^{st} = GAP(T_j^{st})$ denotes the j -th
 480 state feature vector about each image in a batch, \bar{f}_j^{st} denotes
 481 the j -th state mean vector about the whole batch-size. Our
 482 purpose here is to shrink the feature distribution of the state,
 483 thereby reducing the intra-class state discrepancy under the
 484 premise of state classification. The state-based shrinking loss
 485 can be formulated as:

$$\mathcal{L}_{ss} = \mathcal{L}_{ce}^{st} + \frac{1}{N} \sum_{j=1}^N \frac{exp(D_j^{st})}{1 + exp(D_j^{st})}, \quad (11)$$

486 If there exists one sample from different cameras (or view-
 487 points) in a batch, their camera (or viewpoint) feature distance
 488 will become smaller under the premise of classification.

489 D. Global Structural Embedding (GSE) Module

490 After attribute-based enhancement and state-based weaken-
 491 ing operations, we can obtain a final vehicle feature tensor
 492 $T'' \in R^{w \times h \times c}$. Followed by a global average pooling (GAP)
 493 on this tensor, the final vehicle feature vector $f \in R^c$ can
 494 be expressed as $f = GAP(T'')$. The idea of AEE and SWS
 495 is to embed attribute and state information respectively in the
 496 training stage to help learn more discriminative identity feature
 497 f , which is the feature used in the testing stage.

498 Although the vehicle feature f can be trained through the
 499 cross-entropy loss in Eq. (1), the training and testing of vehicle
 500 Re-ID include completely different classes. Therefore it is
 501 insufficient to solely rely on the cross-entropy loss. Addition-
 502 ally, the metric learning methods utilize distance metric loss
 503 (e.g., contrastive loss [29] and triplet loss [5]) to learn a deep
 504 feature embedding space where the samples of a same class
 505 are close to each other and the samples of different classes
 506 are far away. Wu *et al.* [50] propose a simple margin loss that
 507 relaxes unnecessary constraints from traditional contrastive
 508 loss and enjoys the flexibility of the triplet loss. Based on

509 the margin loss [50], we design a new GSE loss to pay more
 510 attention to the hard negative and positive samples by their
 511 state discrepancy and attribute discrepancy.

512 Given a batch of vehicle images $x_i|_{i=1}^B$, B is batch size, we
 513 can get a batch of vehicle feature vectors $f_i|_{i=1}^B$. The margin
 514 loss [50] aims to push its negative samples farther than an
 515 upper boundary u and pull its positive samples closer than a
 516 lower boundary l . Thus $u - l$ is the margin between two
 517 boundaries. Mathematically,

$$\mathcal{L}_m = y_{ij} \max(d_{ij} - l, 0) + (1 - y_{ij}) \max(u - d_{ij}, 0), \quad (12)$$

518 where $y_{ij} = 1$ if $y_i = y_j$, $y_{ij} = 0$ otherwise. $d_{ij} = \|f_i - f_j\|_2$
 519 is the Euclidean distance between two samples.

520 It can be seen from Eq. (12) that margin loss only considers
 521 the instance difference d_{ij} between sample pairs, but ignores
 522 the hierarchical relationship between sample pairs. Concretely,
 523 different vehicle instances with similar appearance can be
 524 further distinguished based on their attribute diversity, we con-
 525 sider this attribute relationship to help the feature embedding
 526 of negative sample pairs:

$$\mathcal{L}_m^- = \exp(-d_{ij}^{at})(1 - y_{ij}) \max(u - d_{ij}, 0), \quad (13)$$

527 where d_{ij}^{at} denotes the mean Euclidean distance of the
 528 attributes between two negative samples in a batch. It worth
 529 noticing that the gradient magnitude concerning any negative
 530 embedding is different in Eq. (13). Mathematically,

$$\|\frac{\partial \mathcal{L}_m^-}{\partial f_j}\|_2 = \exp(-d_{ij}^{at}), \text{if } y_i \neq y_j, \quad (14)$$

531 which means that our GSE module encourages negative
 532 samples with smaller attribute differences to obtain greater
 533 gradient magnitude. If a negative sample pair has the same
 534 attribute, the $d_{ij}^{at} \approx 0$, then $\exp(-d_{ij}^{at})d_{ij} \approx d_{ij}$, which
 535 denotes the feature embedding mainly depends on the instance
 536 difference d_{ij} .

537 In the same way, since the images of the same vehicle
 538 instance with large appearance changes can be further recog-
 539 nized by their state information, we consider this relationship
 540 to help the feature embedding of positive sample pairs:

$$\mathcal{L}_m^+ = \exp(-\frac{1}{\tilde{d}_{ij}^{st}})y_{ij} \max(d_{ij} - l, 0), \quad (15)$$

541 where $\tilde{d}_{ij}^{st} = d_{ij}^{st} + \epsilon$, $\epsilon = 0.000001$ is a small value to avoid
 542 zero denominators, d_{ij}^{st} is the mean Euclidean distance of the
 543 states between two positive samples in a batch. $\exp(-\frac{1}{\tilde{d}_{ij}^{st}})$
 544 can be considered as a gradient magnitude of positive embed-
 545 ding, which means that our GSE module encourages positive
 546 samples with larger state differences to obtain greater gradient
 547 magnitude.

548 The state and attribute guided global structural embedding
 549 loss is:

$$\mathcal{L}_{gse} = S_{ij}y_{ij} \max(d_{ij} - l, 0) + W_{ij}(1 - y_{ij}) \max(u - d_{ij}, 0), \quad (16)$$

550 where $S_{ij} = \exp(-\frac{1}{\tilde{d}_{ij}^{st} + \epsilon})$ and $W_{ij} = \exp(-d_{ij}^{at})$ construct
 551 a global structure for the whole batch-size vehicle images.
 552 If $S_{ij} = W_{ij} = 1$, \mathcal{L}_{gse} is equivalent to margin loss [50].

553 $S_{ij} \in [0, 1]$ and $W_{ij} \in [0, 1]$ can be regarded as state-related
 554 weights and attribute-related weights respectively.

555 In GSE module, the designed loss can be explained as giving
 556 larger weights for hard negatives and positives. Note that the
 557 attribute and state features are imposed into the loss function.
 558 The corresponding gradients are as following:

$$\begin{aligned} \|\frac{\partial \mathcal{L}_m^-}{\partial f_j}\|_2 &= \exp(-d_{ij}^{at})(u - d_{ij}), \text{if } y_i \neq y_j, \\ \|\frac{\partial \mathcal{L}_m^+}{\partial f_j}\|_2 &= \exp(-1/\tilde{d}_{ij}^{st})(d_{ij} - l)/(\tilde{d}_{ij}^{st} * \tilde{d}_{ij}^{st}), \text{else,} \end{aligned} \quad (17)$$

559 which means that our GSE module encourages negative sam-
 560 ples with smaller instance differences and attribute differences
 561 to obtain greater gradient magnitude of the attribute. Even if t-
 562 wo negative samples have the same attributes, the gradient still
 563 exists as $\|\frac{\partial \mathcal{L}_m^-}{\partial f_j}\|_2 = (u - d_{ij})$. Homologous, our GSE module
 564 encourages positive samples with larger instance differences
 565 and state differences to obtain greater gradient magnitude of
 566 the state, until the distance between the positive samples is
 567 less than the lower boundary.

568 To reduce hand-tuned hyperparameters, we reconsider the
 569 goals of attribute-based expanding and state-based shrinking,
 570 and design a new loss function \mathcal{L}_{aess} to replace the original
 571 loss function \mathcal{L}_{ae} and \mathcal{L}_{ss} . Mathematically,

$$\mathcal{L}_{aess} = \alpha(\mathcal{L}_{ce}^{at} + \mathcal{L}_{ce}^{st}) + \frac{\frac{1}{N} \sum_{j=1}^N \exp(D_j^{st})}{\frac{1}{M} \sum_{i=1}^M \exp(D_i^{at}) + \frac{1}{N} \sum_{j=1}^N \exp(D_j^{st})}, \quad (18)$$

572 where $\alpha = \frac{2}{M+N}$ is an adaptive parameter inversely propor-
 573 tional to the number of annotations. D_i^{at} (D_j^{st}) represents the
 574 i -th attribute (j -th state) standard deviation. Under the premise
 575 of attribute/state classification, the attribute difference of all
 576 samples is enlarged, while the state difference is reduced. The
 577 final objective function for our ASSEN model rewrite as:

$$\mathcal{L}_{total} = \mathcal{L}_{ce}^{id} + \mathcal{L}_{aess} + \eta \mathcal{L}_{gse}, \quad (19)$$

578 where only η is used to balance the classification learning and
 579 metric learning.

IV. EXPERIMENT

580 To validate the superiority of the proposed Attribute and S-
 581 tate guided Structural Embedding Network (ASSEN) method,
 582 it is compared with state-of-the-art vehicle Re-ID approaches
 583 on three large-scale databases.

A. Datasets

585 **VeVi-776 dataset** [7] consists of 49357 images of 776 distinct
 586 vehicles captured in 20 non-overlapping cameras with various
 587 orientations and lighting conditions, where 576 identities with
 588 37778 images and 200 identities with 11579 images are
 589 assigned as training and testing respectively. Furthermore,
 590 1678 images from 200 identities have been selected as the
 591 queries from the testing set. The original VeVi-776 [7] contains
 592 the labels of the vehicle IDs, camera IDs, color IDs and type
 593 IDs, while Zheng *et al.* [25] have annotated the viewpoint
 594

595 information, including *front*, *front_side*, *side*, *rear_side*,
 596 and *rear*. We use two kinds of state information (camera,
 597 viewpoint) and two kinds of attribute information (color, type)
 598 in VeRi-776 dataset [7].

599 **VERI-Wild dataset** [6] is a newly released dataset. Different
 600 from VeRi-776 [7] captured at day, VERI-Wild [6] are
 601 captured at both day and night. The training subset consists
 602 of 277797 images of 30671 vehicles. Besides, there are three
 603 different scale testing subsets, *i.e.*, Test3000 (Small), Test5000
 604 (Medium), and Test10000 (Large). Except for vehicle ID
 605 information, VERI-Wild [6] contains various labels of camera,
 606 color, type, and manufacturer annotations. Furthermore, we
 607 have annotated the time labels according to the acquisition
 608 hour of each image. For example, the image captured at
 609 22:15:29 is annotated as 22, and there are 24 time IDs in
 610 total. We use two kinds of state information (camera, time) and
 611 three kinds of attribute information (color, type, manufacturer)
 612 in VERI-Wild dataset [6].

613 **VehicleID dataset** [4] is composed of 221567 images from
 614 26328 unique vehicles. Half of the identities, *i.e.*, 13164,
 615 serves for training while the other half for testing evaluation.
 616 There are 6 testing splits with various gallery sizes as 800,
 617 1600, 2400, 3200, 6000, and 13164. Following the protocol
 618 in [14, 17, 18], we use the first three splits Test800 (Small),
 619 Test1600 (Medium) and Test2400 (Large) for testing. This
 620 procedure is repeated ten times and the averaged metrics. **Note**
 621 **that VehicleID [4] only contains ID information without any**
 622 **attribute or state information. Therefore, we use the attribute**
 623 **and state branch parameters pre-trained on VERI-Wild [6] to**
 624 **obtain state and attribute information for VehicleID [4].**

625 *B. Evaluation Metrics*

626 Following the general evaluation protocols in the Re-ID
 627 field [1, 53, 57], the Rank-1 identification rate (R-1), Rank-5
 628 identification rate (R-5), and mean average precision (mAP)
 629 are used as performance metrics. Rank-score is an estimation
 630 of finding the correct match in the Rank-K returned results.
 631 The mAP is a comprehensive index that considers both the
 632 precision and recall of the results. To evaluate the stability
 633 of our model, we train the model in 10 random trials on
 634 each dataset and take the average result as our performance.
 635 The corresponding standard deviation values are updated in
 636 Table I - IV.

637 *C. Implementation Details*

638 **Network Architecture.** We adopt ResNet-50 [51] as the
 639 backbone model in our experiments. In our implementation, all
 640 the input images are resized to $W \times H \times C = 256 \times 256 \times 3$.
 641 Follow [53], we remove the last spatial down-sampling operation
 642 in ResNet-50 [51]. After the backbone model, the size of the feature tensor is $w \times h \times c = 16 \times 16 \times 2048$.
 643 For classifiers, we use a batch normalization layer [58] and
 644 a fully connected layer followed by a softmax function. For
 645 data augmentation, the images are augmented with random
 646 horizontal flipping, padding 10 pixels, random cropping, and
 647 random erasing [59]. The Adam optimizer [60] is used with a
 648 batch size of 64. We further evaluate our method on a stronger
 649

TABLE I
 COMPARISON RESULTS OF OUR METHOD AGAINST THE
 STATE-OF-THE-ART METHODS ON VERI-776 DATASET (IN %).

Methods	mAP	Rank-1	Rank-5	Reference
(1)	BOW-CN [1]	12.2	33.9	53.7 ICCV 2015
	LOMO [2]	9.6	25.3	CVPR 2015
	GoogLeNet [3]	17.9	52.3	CVPR 2015
	FACT [30]	18.8	52.2	ICME 2016
	FDA-Net [6]	55.5	84.3	CVPR 2019
(2)	FastReID [52]	80.4	96.5	arXiv 2020
	OIFE [32]	48.0	65.9	- ICCV 2017
	SCPL [8]	58.3	83.5	ICCV 2017
(3)	NuFACT [9]	48.5	76.9	TMM 2018
	VAMI [13]	50.1	77.0	CVPR 2018
	EALN [14]	57.4	84.4	TIP 2019
(4)	UMTS [11]	75.9	95.8	AAAI 2020
	RAM [16]	61.5	88.6	ICME 2018
	AAVER [18]	61.2	89.0	ICCV 2019
(5)	PRN [17]	74.3	94.3	CVPR 2019
	PPT [19]	80.6	96.5	MM 2020
DF-CVTC [25]	61.1	91.3	95.8	TETCI 2021
	SAN [22]	72.5	93.3	MST 2020
ASSEN	81.3 ± 0.2	96.9 ± 0.1	98.7 ± 0.1	Ours
Fast_ASSEN	81.7 ± 0.2	97.3 ± 0.1	98.8 ± 0.1	Ours

baseline FastReID [52]. Note that due to the GPU memory
 650 limitations, we implement FastReID [52] with the same batch-
 651 size as our method in $16 \text{ ids} * 4 \text{ imgs}$ for fair comparison. The
 652 new architecture is named Fast_ASSEN in the experiments.
 653

654 **Hyper Parameters.** In Attribute-based Enhancement and Ex-
 655 panding (AEE) module, β_1 is used to balance the original
 656 tensor and the enhanced tensor and set as 0.05. In State-
 657 based Weakening and Shrinking (SWS) module, β_2 is used
 658 to balance the original tensor and the weakened tensor and
 659 set as 0.05. In Global Structural Embedding (GSE) module,
 660 we empirically fix the upper and lower boundaries in the GSE
 661 module to 1 and 0.3, following the commonly used margin
 662 loss [50]. In the final objective function, the weight parameter
 663 $\eta = 0.3$. These hyperparameters will be discussed in detail in
 664 Table VI. We run our experiments on two Tesla P100 GPU
 665 with 16 GB RAM. **Our model requires about 13.5 GB of RAM**
 666 **and 348 minutes of training time on VeRi-776 dataset [7].** The
 667 base learning rate is 3.5×10^{-4} and the learning rate decays
 668 to 3.5×10^{-5} and 3.5×10^{-6} at the 40-th epoch and the 70-
 669 th epoch respectively. Our model is trained in a total of 120
 670 epochs.

671 **Compared Methods.** We compare our method with some
 672 state-of-the-art methods which mainly fail into four categories.
 673

674 (1) *Global feature based methods.* E.g., Bag-of-Words +
 675 Color Names (BOW-CN) [1], Local Maximal Occurrence
 676 (LOMO) [2], GoogLeNet [3], Fusion of Attributes and Color
 677 feaTures (FACT) [30], Feature Distance Adversarial Network
 678 (FDA-Net) [6], Deep Relative Distance Learning (DRDL) [4],
 679 Triplet [5], Softmax [7], Hard-aware Deeply Cascaded embed-
 680 ding (HDC) [61], Unlabeled-GAN [62], Group-sensitive Triplet
 681 Embedding (GSTE) [41].

682 (2) *Path based methods.* E.g., Orientation Invariant Feature
 683 Embedding (OIFE) [32], Siamese-CNN + Path + LSTM
 684 (SCPL) [8], Null space base Fusion of Attribute and Color
 685 feaTures (NuFACT) [9].

686 (3) *Viewpoint based methods.* E.g., Viewpoint-aware Attentive

TABLE II

COMPARISON RESULTS OF OUR METHOD AGAINST THE STATE-OF-THE-ART METHODS ON VEHICLEID DATASET (IN %).

Methods	Small		Medium		Large		
	R-1	R-5	R-1	R-5	R-1	R-5	
(1) DRDL [4]	BOW-CN [1]	13.1	22.7	12.9	21.1	10.2	17.9
	LOMO [2]	19.7	32.1	19.0	29.5	15.3	25.6
	GoogLeNet [3]	47.9	67.4	43.5	63.5	38.2	59.5
	FACT [30]	48.9	66.7	46.4	64.4	41.0	60.0
	FDA-Net [6]	-	-	59.8	77.1	55.5	74.7
	FastReID [52]	82.3	95.5	80.7	72.7	77.8	90.1
(2) OIFE [32]	-	-	-	-	67.0	82.9	
	NuFACT [9]	48.9	69.5	43.6	65.3	38.6	60.7
(3) EALN [14]	VAMI [13]	63.1	83.3	52.9	75.1	47.3	70.3
	UMTS [11]	75.1	88.1	71.8	83.9	71.0	69.3
	RAM [16]	80.9	-	78.8	-	76.1	-
(4)	AAVER [18]	75.2	91.5	72.3	87.0	67.7	84.5
	PRN [17]	74.7	93.8	68.6	90.0	63.5	85.6
	PPT [19]	78.4	92.3	75.0	88.3	74.2	86.4
	DF-CVTC [25]	79.6	92.3	76.0	89.4	74.8	87.0
(5)	ROIVR [23]	75.2	88.1	72.2	84.4	70.5	82.1
	SAN [22]	76.1	91.2	73.1	87.5	71.2	84.7
	ASSEN	79.7	94.3	78.4	91.3	75.6	88.3
ASSEN		85.2\pm0.2	97.7\pm0.1	82.7	95.7	80.9	93.9
Fast_ASSEN		86.0\pm0.3	97.8\pm0.1	84.5	96.0	82.4	94.3

687 Multi-view Inference (VAMI) [13], Embedding Adversarial
688 Learning (EALN) [14], Uncertainty-aware Multi-shot Teacher-
689 Student Network (UMTS) [11].

690 (4) *Local information enhancement methods.* E.g., Region-
691 aware deep Model (RAM) [16], Adaptive Attention Model
692 for Vehicle Re-identification (AAVER) [18], Part-regularized
693 Near-duplicate (PRN) [17], Part Perspective Transformation
694 (PPT) [19].

695 (5) *Attribute based methods.* E.g., Jointly learns Deep Feature
696 representations, Camera Views, vehicle Types and Colors (DF-
697 CVTC) [25], Two-branch Stripe-based and Attribute-aware
698 Network (SAN) [22], Region of Interests-based Vehicle Re-
699 identification (ROIVR) [23].

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

701 **Evaluation Results on VeVi-776.** Table I reports the performance
702 comparison of our method against the state-of-the-art
703 methods on VeVi-776 dataset [7]. From which we can see, the
704 local information enhancement method PPT [19] has higher
705 performance on VeVi-776 [7] compared with the method
706 UMTS [11] based on viewpoint learning. The reason may be
707 because the viewpoint change of VeVi-776 [7] is not too drastic,
708 challenges mainly come from similar vehicles. Compared
709 with the method based on local information enhancement and
710 viewpoint-based methods, our approach significantly beats the
711 state-of-the-art methods as 81.3% and 96.9% on mAP and
712 the Rank-1 respectively. Although the second-best method
713 PPT [19] achieves 80.6% and 96.5% on mAP and Rank-1
714 respectively. PPT [19] propose a part perspective transform
715 module to map key points related to part regions to a unified
716 viewpoint on feature space. However, keypoint extraction
717 usually requires a large amount of annotated data which is time
718 and labor consuming, and inaccurate results of keypoint would
719 affect the performance of vehicle Re-ID greatly. Our ASSEN
720 significantly surpasses the most competitive attribute-based

TABLE III

COMPARISON RESULTS ON MAP OF OUR METHOD AGAINST THE STATE-OF-THE-ART METHODS ON VERI-WILD DATASET (IN %).

Methods	Small		Medium		Large		Reference
	R-1	R-5	R-1	R-5	R-1	R-5	
(1) HDC [61]	GoogLeNet [3]	24.3	24.2	21.5	CVPR 2015		
	Triplet [5]	15.7	13.3	9.9	CVPR 2015		
	Softmax [7]	26.4	22.7	17.6	ECCV 2016		
	DRDL [4]	22.5	19.3	14.8	CVPR 2016		
	Unlabeled-GAN [62]	29.1	24.8	18.3	ICCV 2017		
	GSTE [41]	31.4	26.2	19.5	TMM 2018		
(3) AAVER [18]	FDA-Net [6]	35.1	29.8	22.8	CVPR 2019		
	FastReID [52]	81.9	75.7	66.7	arXiv 2020		
	UMTS [11]	72.7	66.1	54.2	AAAI 2020		
(4) PPT [19]	ASSEN	80.6\pm0.2	74.5\pm0.1	66.2\pm0.1	Ours		
	Fast_ASSEN	84.3\pm0.3	78.7\pm0.2	70.1\pm0.1	Ours		

TABLE IV

COMPARISON RESULTS ON RANK SCORE OF OUR METHOD AGAINST THE STATE-OF-THE-ART METHODS ON VERI-WILD DATASET (IN %).

Methods	Small		Medium		Large		
	R-1	R-5	R-1	R-5	R-1	R-5	
(1) HDC [61]	GoogLeNet [3]	57.2	75.1	53.2	71.1	44.6	63.6
	Triplet [5]	44.7	63.3	40.3	59.0	33.5	51.4
	Softmax [7]	53.4	75.0	42.2	69.9	37.9	59.9
	DRDL [4]	57.0	75.0	51.9	71.0	44.6	61.0
	HDC [61]	57.1	78.9	49.6	72.3	44.0	64.9
	Unlabeled-GAN [62]	58.1	79.6	51.6	74.4	43.6	65.5
(3) AAVER [18]	GSTE [41]	60.5	80.1	52.1	74.9	45.4	66.5
	FDA-Net [6]	64.0	82.8	57.8	78.3	49.4	70.5
	FastReID [52]	96.3	99.2	94.5	98.7	91.1	97.6
	UMTS [11]	84.5	-	79.3	-	72.8	-
(4) PPT [19]	ASSEN	94.9\pm0.1	98.3\pm0.1	91.7	96.5	88.8	94.7
	Fast_ASSEN	97.1\pm0.1	99.7\pm0.1	95.6	99.2	93.9	98.4

method SAN [22] by +8.8% and +3.6% in mAP and Rank-1 accuracies respectively. The key reason is SAN [22] only considers the enhancement of attributes while ignoring the state diversity. By jointly considering the enhancement of attributes, the weakening of states and the hierarchical relationships in the vehicle Re-ID network, our ASSEN learns more robust feature representation on VeVi-776 dataset [7] comparing to the state-of-the-art methods. Fast_ASSEN further boosts the performance in both mAP and ranking scores.

Evaluation Results on VehicleID. Table II shows the comparison results on VehicleID [4] on three different testing sets. The vehicle images in VehicleID [4] only contain two viewpoints, e.g., front and rear, which result in drastic viewpoint changes. As reported in Table II, the method UMTS [11] based on viewpoint learning has higher performance than the local information enhancement method PPT [19] on VehicleID [4] compared with VeVi-776 dataset [7]. This implies that it is necessary to consider joint learning from different viewpoints in VehicleID [4]. In addition to the viewpoint factor similar as UMTS [11], our ASSEN also considers the time factor and the camera factor, as well as the attribute information to enhance the discrimination ability. As shown in Table II, the Rank-1 accuracies of our approach improve 4.3%, 3.9% and 4.8% than UMTS. Note that our methods, ASSEN and Fast_ASSEN, without any attribute and state annotation

TABLE V
ABLATION STUDY ON VeRI-776, VeRI-WILD AND VEHICLEID (IN %).

Variant	VeRI-776				VehicleID						VeRI-Wild					
			Small		Medium		Large		Small		Medium		Large			
	mAP	R-1	mAP	R-1												
(a1) baseline (L_{ce})	74.3	94.8	85.0	78.6	81.3	74.6	79.6	72.0	72.3	89.5	64.3	84.9	53.6	80.7		
(a2) + AEE	76.8	95.5	85.9	82.4	82.0	77.6	80.7	74.7	73.5	92.6	66.4	87.1	57.3	81.1		
(a3) + SWS	77.3	95.2	86.1	79.5	83.0	75.4	81.4	72.2	75.6	90.3	69.6	85.2	62.9	80.8		
(a4) + GSE	78.9	95.9	88.9	83.4	85.2	80.7	82.8	77.4	77.1	93.2	72.7	89.6	63.8	85.1		
(a5) + AEE + SWS	78.3	95.6	87.0	82.6	83.8	79.2	81.3	76.6	76.9	93.5	71.8	89.9	63.1	83.3		
(a6) + AEE + SWS + GSE	81.3	96.9	90.4	85.2	88.0	82.7	85.5	80.9	80.6	94.9	74.5	91.7	66.2	88.8		
(b1) FastReID ($L_{ce} + L_{tri}$)	80.4	96.5	85.8	82.3	83.6	80.7	82.6	77.8	81.9	96.3	75.7	94.5	66.7	91.1		
(b2) + AEE	80.0	96.8	86.0	83.0	84.0	81.6	82.6	78.2	81.9	96.4	75.2	94.9	66.3	92.2		
(b3) + SWS	80.5	96.5	86.6	82.4	84.3	80.7	83.0	77.9	82.2	96.4	76.2	94.6	67.3	91.2		
(b4) + AEE + SWS	81.2	96.9	88.1	85.2	86.9	82.6	85.5	81.0	83.0	96.4	78.0	94.9	69.1	92.6		
(b5) + AEE + SWS + GSE	81.7	97.3	90.9	86.0	89.1	84.5	87.2	82.4	84.3	97.1	78.7	95.6	70.1	93.9		
(c1) baseline ($L_{ce} + L_{tri}$)	76.6	95.7	85.0	80.2	82.9	77.5	79.7	73.8	76.2	91.8	68.0	87.3	57.8	83.5		
(c2) + AEE	76.9	96.3	85.7	81.9	84.0	78.3	80.0	76.1	76.9	93.2	68.9	88.5	59.3	86.6		
(c3) + SWS	77.8	95.9	87.3	81.6	85.9	78.3	81.3	75.7	77.9	92.8	72.2	87.8	63.1	84.2		
(c4) + AEE + SWS	79.8	96.5	88.7	83.4	87.1	81.0	82.8	78.0	79.0	93.9	73.6	91.0	64.8	88.8		
(c5) + AEE + SWS + GSE	81.3	97.0	90.4	85.4	88.6	83.6	85.9	81.2	81.0	95.4	75.2	91.9	66.8	90.2		

on VehicleID [4], still significantly beats the state-of-the-art attribute-based methods, especially comparing SAN [22] and ROIVR [23] with additional attribute annotations. This further verifies the generality of our method of leveraging the attribute and state information on more general scenarios.

Evaluation Results on VeRI-Wild. As shown in Table III and Table IV, our ASSEN achieves competitive results on all of the testing subsets on the VeRI-Wild dataset [6]. Specifically, the Rank-1 accuracies of our approach achieve 94.9%, 91.7% and 88.8% on Test3000 (small), Test5000 (middle) and Test10000 (large) respectively, which improve 3.0%, 2.6% and 4.0% than the second-best method PPT [19]. Meanwhile, the mAP of our method achieve 80.6%, 74.5% and 66.2% on Test3000 (small), Test5000 (middle) and Test10000 (large) respectively, which improve 6.4%, 7.0% and 6.9% than the second-best method PPT [19]. The data size of VeRI-Wild dataset [6] is about 6 times that of VeRI-776 dataset [7]. Although our Re-ID performance is very close to PPT [19] on VeRI-776 [7], our performance on VeRI-Wild dataset [6] is much higher than that of PPT [19], which implies the promising performance in potential large-scale applications. Integrating our method into FastReID [52] consistently improves the performance both mAP and ranking scores.

E. Ablation Study

Component Study. To verify the contribution of the components in our model, we implement several variants of our method on the three datasets, as reported in Table V. Our baseline is ResNet-50 with \mathcal{L}_{ce} . By progressively introducing the attribute-based enhancement and expanding module (AEE), state-based weakening and shrinking module (SWS), and global structural embedding module (GSE) into the baseline, both mAP, and Rank-1 scores significantly increase on all the three datasets with different test settings. This verifies the contribution of each component in our model.

Analysis of Different Baselines. To further validate the effectiveness of our method, we evaluate the component of two stronger baselines, (1) FastReID [52], which is a strong baseline for vehicle Re-ID as shwon in Table V (b1-b5),

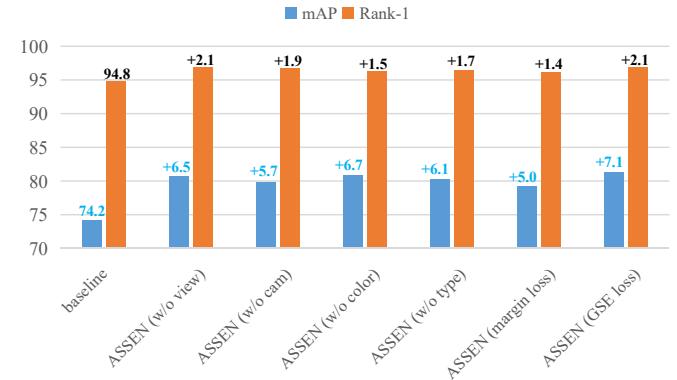


Fig. 4. Subcomponent analysis on VeRI-776.

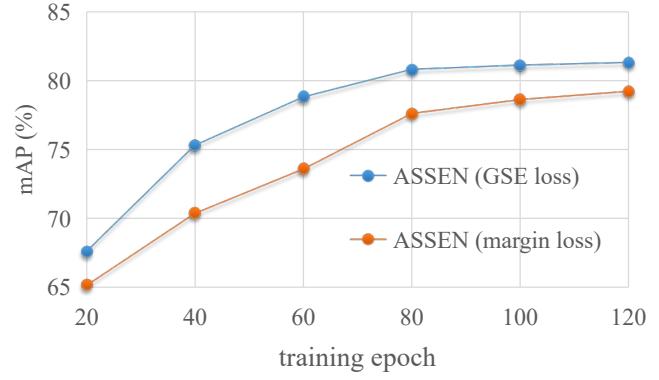


Fig. 5. The mAP performance against the number of training epochs using global structural embedding loss and margin loss [50] on VeRI-776.

and (2) the baseline in the state-of-the-art methods such as UMTS [11], PPT [19], FastReID [52], with both cross-entropy loss and triplet loss (baseline ($L_{ce} + L_{tri}$)), as shown in Table V (c1-c5). Note that due to the GPU memory limitations, we implement FastReID [52] with the same batch size as our method in 16 $ids \times 4$ $imgs$ for fair comparison. Consistently, all

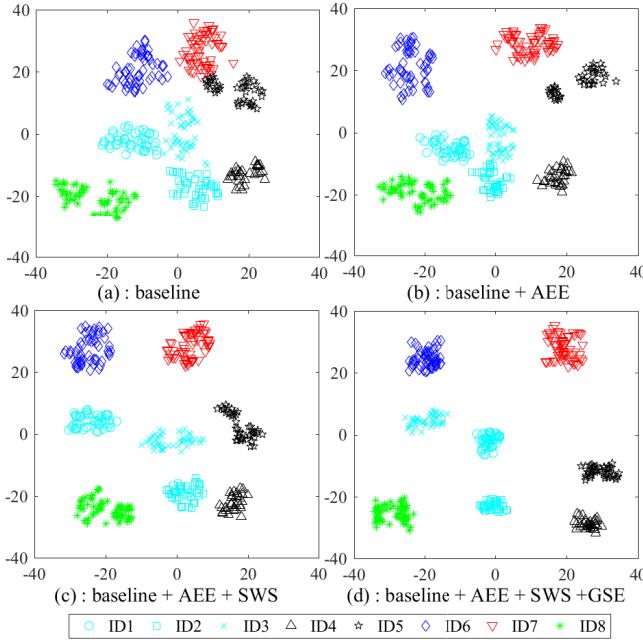


Fig. 6. T-SNE [63] visualization of the learned feature embeddings on 329 images from 8 identities in the VeRi-776 testing set. The points with the same shape indicate the same identity, while the different colors represent different attributes. These points contain the samples (of ID1, ID2, ID3, ID4) in Fig. 2.

790 the AEE, SWS, and GSE modules make effective contributions
791 in our method on the new baselines.

792 Furthermore, Fig. 6 visualizes the feature map during the
793 ablation study. The AEE module increases the inter-class
794 distance of different attributes, while the SWS module reduces
795 the intra-class distance and increase the inter-class distance
796 with the same attribute. GSE module can further reduce the
797 intra-class gap and increase the inter-class gap.

798 **Subcomponent Study.** To further evaluate the contribution of
799 each state and attribute, we evaluate our method by removing a
800 certain attribute or state as shown in Fig. 4. It is clear that each
801 attribute or state information contributes to our ASSEN model.
802 In addition, we compare the performance and convergence
803 of the global embedding loss (ASSEN (GSE loss)) with the
804 margin loss [50] (ASSEN (margin loss)) as shown in Fig. 4
805 and in Fig. 5, respectively. ASSEN (margin loss) denotes
806 *baseline + AEE + SWS + margin loss* and has the same
807 hyperparameters as ASSEN. As shown in Eq. (16), the margin
808 loss [50] can be seen as a special form of our global em-
809 bedding loss without weight. By considering the hierarchical
810 relationships (inter-class attribute discrepancy and intra-class
811 state discrepancy) between vehicles, our global embedding
812 loss converges faster and achieves better performance.

813 F. Parameter Analysis

814 There are five important parameters in our model. β_1 and
815 β_2 balances the contribution of the enhanced feature and
816 the weakened feature respectively, while u and l control
817 the margin between positive samples and negative samples
818 respectively. In the final loss function, η control the weight
819 of classification learning and metric learning. We empirically

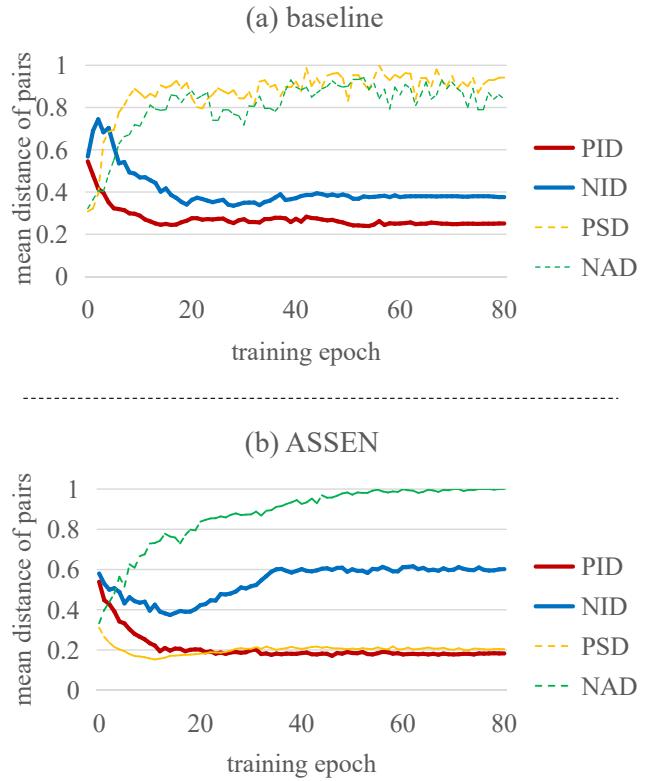


Fig. 7. Feature distance discrepancy of the baseline and ASSEN. Distance discrepancy mainly includes the instance distance between positive sample pairs (PID), the instance distance between negative sample pairs (NID), the state distance between positive sample pairs (PSD) and the attribute distance between negative sample pairs (NAD).

TABLE VI
PARAMETER ANALYSIS ON VERI-776 (IN %).

Parameter	Setting	mAP	R-1	Parameter	Setting	mAP	R-1
β_1	0	79.6	95.9	η	0.1	80.8	96.8
	0.05	81.3	96.9		0.2	80.5	97.0
	0.1	80.6	97.0		0.3	81.3	96.9
β_2	0	80.8	96.8	l	0.2	80.7	96.7
	0.05	81.3	96.9		0.3	81.3	96.9
	0.1	80.1	96.5		0.6	80.6	96.6
u	0.8	79.8	96.1		0.2	80.7	96.7
	1.0	81.3	96.9		0.3	81.3	96.9
	1.2	80.3	96.6		0.4	81.0	96.5

set $\beta_1 = 0.05$, $\beta_2 = 0.05$, $u = 1$, $l = 0.3$ and $\eta = 0.3$. The
820 parameter analysis results with diverse parameter changes on
821 VeRi-776 [7] are shown in Table VI, which demonstrates that
822 our model is not sensitive to the parameters.
823

824 G. Analysis of Distance Discrepancy

825 To further verify the ability of handling the inter-class
826 similarity and intra-class discrepancy of our method, we
827 visualize the instance distance of positive sample pairs (PID),
828 the instance distance of negative sample pairs (NID), the
829 state distance of positive sample pairs (PSD) and the attribute
830 distance of negative sample pairs (NAD). **We first average**
831 **the PID/NID/PSD/NAD of each anchor in a batch, and then**

average over all batches in an epoch. As shown in Fig. 7, our ASSEN significantly shortens the state distance of positive samples (PSD), while increasing the attribute distance of negative samples (NAD), which shortens the instance distance of positive samples (PID) and enlarges the instance distance of negative samples (NID). It shows that weakening the state information can help reduce the intra-class distance, and enhancing the attribute information can help enlarge the inter-class distance. They are both effective ways to improve the discrimination of the vehicle Re-ID network.

V. CONCLUSION

To our best knowledge, this is the first work to solve the problem of Re-ID by enhancing attribute information and weakening state information. In this paper, we first argue the factors that cause the challenge of vehicle Re-ID into state factors and attribute factors. We have contributed an attribute and state guided structural embedding network (ASSEN), followed by three novel modules: attribute-based enhancement and expanding, state-based weakening and shrinking, global structural embedding. Comparing with state-of-the-art vehicle Re-ID methods, extensive experiments demonstrate the promising performance of the proposed method. Although our method requires additional state information and attribute information, this information is easy to obtain and has strong generalization capabilities. In the future, we will consider applying the idea of reducing state discrepancy and increasing attribute discrepancy to other recognition tasks (pedestrians, animals) and unsupervised vehicle Re-ID problems.

REFERENCES

- [1] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian, “Scalable person re-identification: A benchmark,” in *IEEE International Conference on Computer Vision*, 2015, pp. 1116–1124.
- [2] S. Liao, Y. Hu, X. Zhu, and S. Z. Li, “Person re-identification by local maximal occurrence representation and metric learning,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 2197–2206.
- [3] L. Yang, P. Luo, C. C. Loy, and X. Tang, “A large-scale car dataset for fine-grained categorization and verification,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3973–3981.
- [4] H. Liu, Y. Tian, Y. Wang, L. Pang, and T. Huang, “Deep relative distance learning: Tell the difference between similar vehicles,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2167–2175.
- [5] F. Schroff, D. Kalenichenko, and J. Philbin, “Facenet: A unified embedding for face recognition and clustering,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 815–823.
- [6] 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 *IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3235–3243.
- [7] 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*, 2016, pp. 869–884.
- [8] 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 *IEEE International Conference on Computer Vision*, 2017, pp. 1918–1927.
- [9] 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, 2018.
- [10] 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. 3275–3278, 2018.
- [11] X. Jin, C. Lan, W. Zeng, and Z. Chen, “Uncertainty-aware multi-shot knowledge distillation for image-based object re-identification,” in *The Thirty-Fourth AAAI Conference on Artificial Intelligence*, 2020, pp. 11 165–11 172.
- [12] A. Porrello, L. Bergamini, and S. Calderara, “Robust re-identification by multiple views knowledge distillation,” in *European Conference on Computer Vision*, 2020, pp. 93–110.
- [13] Y. Zhou and L. Shao, “Viewpoint-aware attentive multi-view inference for vehicle re-identification,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 6489–6498.
- [14] Y. Lou, Y. Bai, J. Liu, S. Wang, and L. Duan, “Embedding adversarial learning for vehicle re-identification,” *IEEE Transactions on Image Processing*, vol. 28, no. 8, pp. 3794–3807, 2019.
- [15] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” in *Advances in Neural Information Processing Systems*, 2014, pp. 2672–2680.
- [16] X. Liu, S. Zhang, Q. Huang, and W. Gao, “Ram: A region-aware deep model for vehicle re-identification,” in *International Conference on Multimedia and Expo*, 2018, pp. 1–6.
- [17] B. He, J. Li, Y. Zhao, and Y. Tian, “Part-regularized near-duplicate vehicle re-identification,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3997–4005.
- [18] 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,” in *IEEE International Conference on Computer Vision*, 2019, pp. 6132–6141.
- [19] D. Meng, L. Li, S. Wang, X. Gao, Z.-J. Zha, and Q. Huang, “Fine-grained feature alignment with part perspective transformation for vehicle reid,” in *ACM international conference on Multimedia*, 2020, pp. 619–627.
- [20] P. Khorramshahi, N. Peri, J.-C. Chen, and R. Chellappa, “The devil is in the details: Self-supervised attention for vehicle re-identification,” in *European Conference on Computer Vision*, 2020, pp. 369–386.
- [21] X. Liu, W. Liu, J. Zheng, C. Yan, and T. Mei, “Beyond the parts: Learning multi-view cross-part correlation for vehicle re-identification,” in *ACM international conference on Multimedia*, 2020, pp. 907–915.
- [22] J. Qian, W. Jiang, H. Luo, and H. Yu, “Stripe-based and attribute-aware network: A two-branch deep model for vehicle re-identification,” *Measurement Science and Technology*, vol. 31, no. 9, p. 095401, 2020.
- [23] 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*, vol. 21, no. 2, pp. 723–734, 2019.
- [24] Y. Lin, L. Zheng, Z. Zheng, Y. Wu, and Z. Hu, “Improving person re-identification by attribute and identity learning,” *Pattern Recognition*, vol. 95, pp. 151–161, 2019.
- [25] H. Li, X. Lin, A. Zheng, C. Li, B. Luo, R. He, and A. Hussain, “Attributes guided feature learning for vehicle re-identification,” *IEEE Transactions on Emerging Topics in Computational Intelligence*, pp. 1–12, 2021.
- [26] T. Chen, L. Lin, R. Chen, Y. Wu, and X. Luo, “Knowledge-embedded representation learning for fine-grained image recognition,” in *International Joint Conference on Artificial Intelligence*, 2018, pp. 627–634.
- [27] X. Liu, J. Wang, S. Wen, E. Ding, and Y. Lin, “Localizing by describing: Attribute-guided attention localization for fine-grained recognition,” in *The Thirty-First AAAI Conference on Artificial Intelligence*, 2017, pp. 4190–4196.
- [28] L. Lin, L. Huang, T. Chen, Y. Gan, and H. Cheng, “Knowledge-guided recurrent neural network learning for task-oriented action prediction,” in *IEEE International Conference on Multimedia and Expo (ICME)*, 2017.
- [29] R. Hadsell, S. Chopra, and Y. Lecun, “Dimensionality reduction by learning an invariant mapping,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2006, pp. 1735–1742.
- [30] X. Liu, W. Liu, H. Ma, and H. Fu, “Large-scale vehicle re-identification in urban surveillance videos,” in *International Conference on Multimedia and Expo*, 2016, pp. 1–6.
- [31] K. Yan, Y. Tian, Y. Wang, W. Zeng, and T. Huang, “Exploiting multi-grain ranking constraints for precisely searching visually-similar vehicles,” in *IEEE International Conference on Computer Vision*, 2017, pp. 562–570.
- [32] 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 *IEEE International Conference on Computer Vision*, 2017, pp. 379–387.
- [33] T. Kipf and M. Welling, “Semi-supervised classification with graph convolutional networks,” in *International Conference on Learning Representations*, 2017, pp. 1–14.

- [34] J. Sochor, A. Herout, and J. Havel, "Boxcars: 3d boxes as cnn input for improved fine-grained vehicle recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 3006–3015.
- [35] Y. Zhou and L. Shao, "Cross-view gan based vehicle generation for re-identification," in *British Machine Vision Conference*, vol. 1, 2017, pp. 1–12.
- [36] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [37] S. Khamis, C.-H. Kuo, V. K. Singh, V. D. Shet, and L. S. Davis, "Joint learning for attribute-consistent person re-identification," in *European Conference on Computer Vision*, 2014, pp. 134–146.
- [38] C. Su, S. Zhang, J. Xing, W. Gao, and Q. Tian, "Multi-type attributes driven multi-camera person re-identification," *Pattern Recognition*, vol. 75, pp. 77–89, 2018.
- [39] C. Sun, N. Jiang, L. Zhang, Y. Wang, W. Wu, and Z. Zhou, "Unified framework for joint attribute classification and person re-identification," in *International Conference on Artificial Neural Networks*, 2018, pp. 637–647.
- [40] Y. Cho and K. Yoon, "Improving person re-identification via pose-aware multi-shot matching," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1354–1362.
- [41] Y. Bai, Y. Lou, F. Gao, S. Wang, Y. Wu, and L. Duan, "Group-sensitive triplet embedding for vehicle reidentification," *IEEE Transactions on Multimedia*, vol. 20, no. 9, pp. 2385–2399, 2018.
- [42] S. Zhou, J. Wang, D. Meng, Y. Liang, Y. Gong, and Z. Nanning, "Discriminative feature learning with foreground attention for person re-identification," *IEEE Transactions on Image Processing*, vol. 28, no. 9, pp. 4671–4684, 2019.
- [43] H.-X. Yu and W.-S. Zheng, "Weakly supervised discriminative feature learning with state information for person identification," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 5527–5537.
- [44] W. Chen, X. Chen, J. Zhang, and K. Huang, "Beyond triplet loss: A deep quadruplet network for person re-identification," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1320–1329.
- [45] K. Sohn, "Improved deep metric learning with multi-class n-pair loss objective," in *Advances in Neural Information Processing Systems*, 2016, pp. 1857–1865.
- [46] X. He, Y. Zhou, Z. Zhou, S. Bai, and X. Bai, "Triplet-center loss for multi-view 3d object retrieval," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1945–1954.
- [47] E. Ustinova and V. Lempitsky, "Learning deep embeddings with histogram loss," in *Neural Information Processing Systems*, 2016, pp. 4170–4178.
- [48] X. Wang, Y. Hua, E. Kodirov, G. Hu, R. Garnier, and N. Robertson, "Ranked list loss for deep metric learning," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5207–5216.
- [49] X. Liu, S. Zhang, and X. Wang, "Group-group loss-based global-regional feature learning for vehicle re-identification," *IEEE Transactions on Image Processing*, vol. 29, pp. 2638–2652, 2020.
- [50] C.-Y. Wu, R. Manmatha, A. J. Smola, and P. Krahenbuhl, "Sampling matters in deep embedding learning," in *IEEE International Conference on Computer Vision*, 2017, pp. 2840–2848.
- [51] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [52] L. He, X. Liao, W. Liu, X. Liu, P. Cheng, and T. Mei, "Fastreid: A pytorch toolbox for general instance re-identification," *arXiv preprint arXiv:2006.02631*, 2020.
- [53] 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 *European Conference on Computer Vision*, 2018, pp. 501–518.
- [54] W. Li, X. Zhu, and S. Gong, "Harmonious attention network for person re-identification," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2285–2294.
- [55] 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*, vol. 28, no. 9, pp. 4328–4338, 2019.
- [56] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba, "Learning deep features for discriminative localization," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2921–2929.
- [57] H. Luo, W. Jiang, Y. Gu, F. Liu, X. Liao, S. Lai, and J. Gu, "A strong baseline and batch normalization neck for deep person re-identification," *IEEE Transactions on Multimedia*, vol. 22, no. 10, pp. 2597–2609, 2020.
- [58] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International Conference on Machine Learning*, 2015, pp. 448–456.
- [59] Z. Zhong, L. Zheng, G. Kang, S. Li, and Y. Yang, "Random erasing data augmentation," in *The Thirty-Fourth AAAI Conference on Artificial Intelligence*, 2020, pp. 13 001–13 008.
- [60] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [61] Y. Yuan, K. Yang, and C. Zhang, "Hard-aware deeply cascaded embedding," in *IEEE International Conference on Computer Vision*, 2017, pp. 814–823.
- [62] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *IEEE International Conference on Computer Vision*, 2017, pp. 2223–2232.
- [63] L. V. Der Maaten and G. E. Hinton, "Visualizing data using t-sne," *Journal of Machine Learning Research*, vol. 9, pp. 2579–2605, 2008.