
REINFORCED PEDESTRIAN ATTRIBUTE RECOGNITION WITH GROUP OPTIMIZATION REWARD

Zhong Ji, Zhenfei Hu, Yaodong Wang

School of Electrical and Information Engineering
Tianjin University
Tianjin

{jizhong, hzf0226, wangyaodong@tju.edu.cn}@tju.edu.cn

Shengjia Li

R&D Department
China Academy of Launch Vehicle Technology
Beijing

sjli@tju.edu.cn

ABSTRACT

Pedestrian Attribute Recognition (PAR) is a challenging task in intelligent video surveillance. Two key challenges in PAR include complex alignment relations between images and attributes, and imbalanced data distribution. Existing approaches usually formulate PAR as a recognition task. Different from them, this paper addresses it as a decision-making task via a reinforcement learning framework. Specifically, PAR is formulated as a Markov decision process (MDP) by designing ingenious states, action space, reward function and state transition. To alleviate the inter-attribute imbalance problem, we apply an Attribute Grouping Strategy (AGS) by dividing all attributes into subgroups according to their region and category information. Then we employ an agent to recognize each group of attributes, which is trained with Deep Q-learning algorithm. We also propose a Group Optimization Reward (GOR) function to alleviate the intra-attribute imbalance problem. Experimental results on the three benchmark datasets of PETA, RAP and PA100K illustrate the effectiveness and competitiveness of the proposed approach and demonstrate that the application of reinforcement learning to PAR is a valuable research direction.

Keywords Pedestrian Attribute Recognition · Intelligent Video Surveillance · Reinforcement Learning · Deep Q-learning

1 Introduction

Nowadays, intelligent video surveillance technology has been widely deployed [1, 2]. Pedestrian attributes, such as age, clothing style, gender, and accessory, are important soft-biometrics in video surveillance applications, such as person re-identification [3, 4, 5, 6], person search [7, 8], human parsing [9], and pedestrian detection [10, 11]. Thus, the recognition of them, called Pedestrian Attribute Recognition (PAR), has received great attention in recent years. Fig. 1 shows some examples from the popular PETA [12], RAP [13] and PA100K [14] datasets.

Generally, PAR has two key challenges: complex alignment relations between images and attributes, and imbalanced data distribution [15]. Since most pedestrian images are captured at far distance from real surveillance scenarios, they always show ambiguous appearance due to occlusion, low image resolution, and illumination. Meanwhile, the pedestrian attributes are quite diverse, for example, some of the “upperbody-Cotton” are in long style while some are in short; some of the “upperbody-Cotton” are zippered, some are not. Thus, the appearance ambiguity and attribute diversity make it quite hard to align the images and attributes. Some approaches employed attention-based sequence to sequence model to explore the complex alignment relations. For example, JRL [16] applies recurrent attention

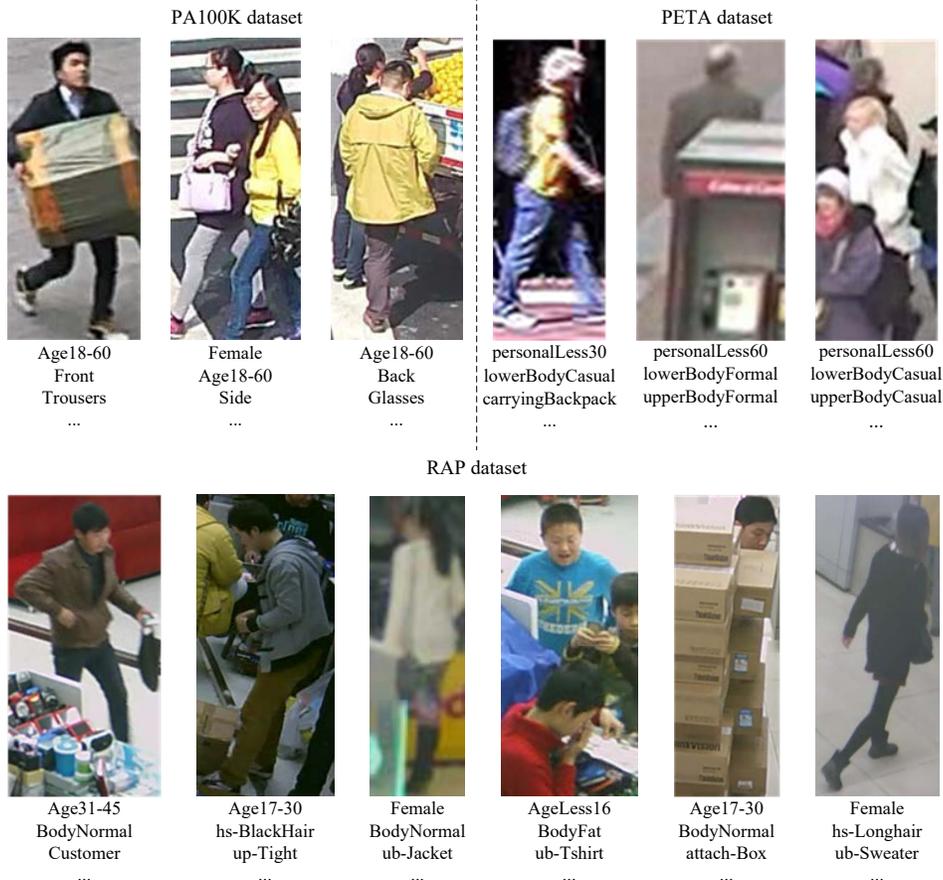


Figure 1: Example images and the corresponding attributes.

mechanism on the output of encoder and reformulates the attribute decoding algorithm to focus on relevant parts when predicting the current attribute.

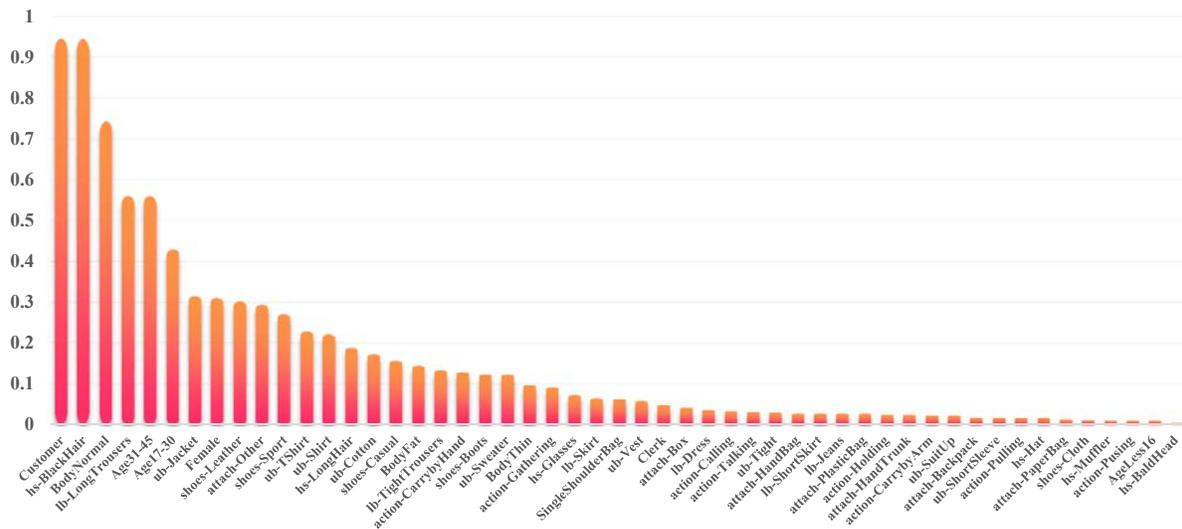
The second challenge is the imbalanced data distribution, which has always been an active research topic in computer vision filed. In PAR, this imbalance is manifested in two aspects, i.e., inter-attribute and intra-attribute imbalance. On the one hand, the inter-attribute distributions are imbalanced, that is to say, some attributes have a large proportion, while some are extremely few. Fig. 2 illustrates the imbalanced attribute distributions of the RAP and PA100K datasets. We could observe, for example, that the ratio of “hs-BlackHair” is 94.48%, while that of “hs-BaldHead” is only 0.38% in RAP dataset. Some approaches alleviate this type imbalance problem by designing special loss functions. For instance, IA²-Net [17] proposes a Focal Cross-Entropy (FCE) loss by combining the cross-entropy loss and the focal loss. Further, MTA-Net [15] proposes a Focal Balanced Loss (FBL) on the basis of FCE loss by increasing the cost of difficult-to-identify attributes. On the other hand, the presence or absence for a single attribute can be regarded as an intra-attribute imbalanced distribution. For example, the presence probability of “hs-Hat” is only 1.67%, while its absence probability is as high as 98.33%. For a specific attribute, such imbalanced distribution enforces the network pay more attention to the majority samples, while ignoring the other valuable few samples. However, this type of imbalance problem is rarely tackled in existing PAR approaches.

Interestingly, almost all the existing approaches regard PAR as a recognition task. Inspired by a large number of effective applications of reinforcement learning (RL) in computer vision, we address PAR with RL by regarding PAR as a decision-making process, that is, a determination whether an attribute exists. Actually, RL has attracted increasing attentions as the success of AlphaGo in 2015 [18]. Afterwards, it has been successfully applied to many directions in computer vision filed, such as video summarization [19], pedestrian tracking [20], and image classification [21, 22]. For example, in [21], ICMDP formulates binary imbalance classification problem as a sequential decision-making process,

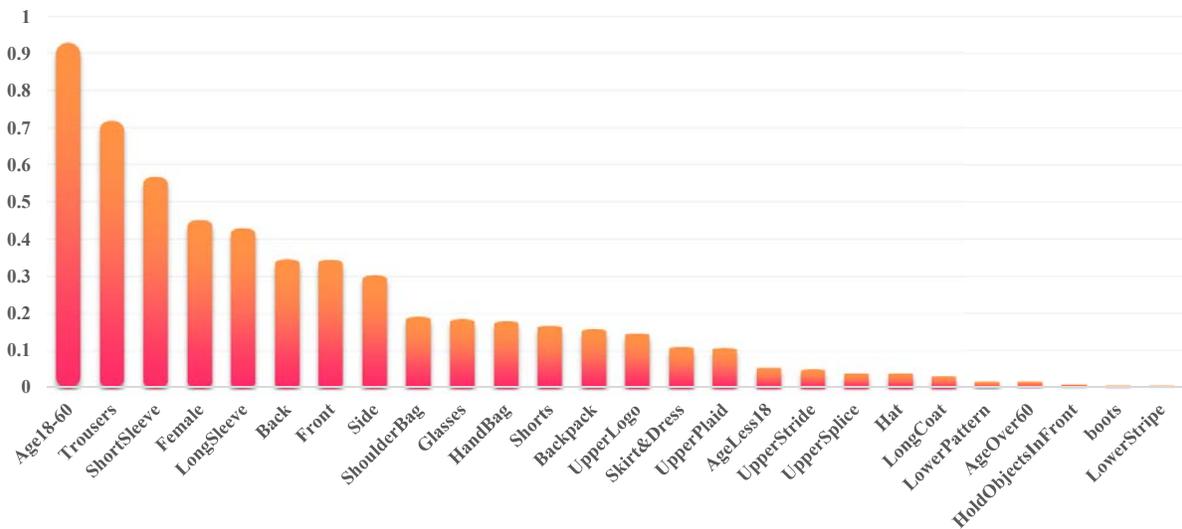
and solve it by Deep Q-learning. He *et al.* [22] introduced Deep Q-learning to curriculum learning for multi-label image classification, which explored the influence of label prediction order on the prediction results.

However, directly employing existing RL-based approaches in PAR has the following two limitations. First, it is hard to be applied to a dataset with too many labels for a single image. This is because as the number of labels increases, the corresponding action space will increasing, which will aggravate the overestimation generated by Deep Q-learning [23]. Secondly, the existing RL approaches cannot make use of the relationship among pedestrian attributes, nor can they effectively alleviate the imbalanced data distribution in PAR datasets.

To address the above challenges and limitations, we propose a novel group reinforcement learning framework. Moreover, we apply an attribute grouping strategy and a group optimization reward function to alleviate the imbalance data distribution problem.



(a) RAP dataset



(b) PA100K dataset

Figure 2: The imbalanced attribute distributions of the RAP and PA100K datasets.

The contributions are highlighted as follows:

- We address the PAR task in a Deep Q-learning framework, as shown in Fig. 3. Although PAR is an attractive paradigm, it has rarely previously been exploited in a deep reinforcement learning framework. Particularly, we define PAR as a Markov decision process by designing image features and encoded attribute labels as states, utilizing binary codes 0 and 1 as action space, and giving corresponding positive or negative rewards according to whether the prediction result is accepted. In this way, the complex alignment between the images and attributes are modeled implicitly. To our best knowledge, it is the first work that defines PAR as a Markov decision process.
- We apply an Attribute Grouping Strategy (AGS) based on the correlation among the attributes in the same pedestrian image region to alleviate the inter-attribute imbalance problem. Further, we propose a Group Optimization Reward (GOR) function to alleviate the intra-attribute imbalance problem, in which the reward function of each group is optimized according to the overall situation of each group of attributes.
- Extensive experiments on three benchmark datasets, i.e. PETA, RAP and PA100K, demonstrate the effectiveness and competitiveness of our proposed Rein-PAR approach.

The remaining sections of the paper are organized as follows. Section 2 reviews the related work. Section 3 introduces our proposed Rein-PAR in detail. Section 4 presents the experiments and ablation study, followed by the conclusion in Section 5.

2 Related work

Pedestrian attributes recognition has been broadly studied for many years. The applications of reinforcement learning in the field of computer vision are also received great attention. This section briefly introduces some works related to these two directions.

2.1 Pedestrian Attribute Recognition

Early PAR approaches applied hand-crafted features, such as Histogram of Oriented Gradients (HoG) [24] and texture histogram [25, 26]. With the renaissance of neural networks, deep learning based approaches have dominated the current PAR approaches. Generally, they can be divided into four groups, that are global-based, part-based, attention-based, and loss-based approaches.

Global-based approaches: As one of the pioneering deep PAR approaches, DeepMAR [27] takes the whole image as input and presents a weighted cross-entropy loss to handle the attribute imbalance challenge. GSR-MAR [28] converts low-resolution images into high-resolution images through the Global Super Resolution Network for PAR. Li *et al.* [29] proposed a CNN-RNN based sequential prediction model that takes global images as input, which can effectively encode scene context and inter-person social relations. JLAC [30] employs graph convolution network to explore the relationship between attributes, which effectively improves the performance of attribute recognition.

Part-based approaches: As PAR is actually a fine-grained recognition task, PGDM [31] recognizes pedestrian attributes by exploring human structure knowledge. It employs pre-trained pose estimation model to obtain keypoints of human image, and extracts part regions according to the keypoints for attribute recognition. Depend on the class activation maps, Liu *et al.* [32] proposed a Localization Guide Network (LGNet), which captures the activation box for each attribute by cropping the high response area of the corresponding activation map. DTM+AWK [33] leverages pose keypoints as auxiliary information to assist in positioning the attribute region. Tang *et al.* [34] proposed an Attribute localization module (ALM) which can discover the most discriminative regions adaptively. The application of auxiliary information has been proven to be effective, however, they [31, 32, 33] largely depend on the accuracy of positioning.

Attention-based approaches: HP-Net [14] applies the attention mechanism to PAR, which consists of two modules: Main Net (M-net) and Attentive Feature Net (AF-net). The AF-net contains multiple branches of multi-directional attention modules, which are applied to different semantic feature levels. In [35], DIAA introduces a multi-scale attention mechanism by directing the network to pay more attention to the spatial parts containing relevant information of the input image. Wu *et al.* [36] proposed a coarse-to-fine attention mechanism to reduce the irrelevant interference areas, which effectively improves the discriminant ability of attribute recognition. JRL [16] utilizes an encoder-decoder architecture to process image context and attribute correlation and applies attention mechanism to better focus on the local regions and obtain more accurate representation. CAS-SAL-FR [37] proposes a cascaded Split-and-Aggregate Learning (SAL) that captures the individuality and commonality of all attributes simultaneously at the feature map level and feature vector level with designed attribute-specific attention module (ASAM) and constrained losses. Although the visual attention mechanism has been successfully applied to PAR, the complexity of pedestrian images makes attention masks fail to obtain the position of a specific attribute.

Loss-based approaches: The design of innovative loss function to alleviate data imbalance in PAR task is a hot research direction. DeepMAR [27] proposes a weighted sigmoid cross entropy loss on the basis of the sigmoid cross entropy loss. Even though some attributes occupy a large proportion, they are still difficult to identify. Therefore, MTA-Net [15] proposes Focal Balance Loss (FBL) based on FCE Loss [17] by increasing the cost of those attributes difficult to recognize.

Different from the existing approaches that formulate PAR as a recognition process, our proposed Rein-PAR approach is a decision-making process, which formulates PAR in the reinforcement learning framework.

2.2 Reinforcement Learning in Computer Vision

The purpose of RL is to learn a policy for the agent from experimental trials by maximizing expected future rewards. It has been demonstrated to be effective in many computer vision tasks [38]. These approaches can be divided into two categories, one is partial-RL approaches and the other is full-RL approaches.

Partial-RL approaches: This line of approach refers to that RL is applied to deal with part of the problem in a task. For example, RL-RBN [39] applies RL successfully to learn an adaptive margin policy to mitigate bias among different races and learn more balanced features, thus effectively improving the racial fairness in face recognition task. In action recognition task, Dong *et al.* [40] found that irrelevant frames have an adverse impact on action prediction, and then applied RL and attention mechanisms to seek the most discriminate frames. In person Re-ID task, AAB [41] employs RL to design a novel Attribute Selection Module (ASM) to discard the noisy attributes and select the key ones. In person retrieval task, APN [42] utilizes policy gradient algorithm to optimize the agent to dynamically generate the optimal partitioning strategies for different images, which effectively reduces the human intervention problem. Wang *et al.* [43] modeled the online key decision process in dynamic video segmentation as a deep RL problem, and obtained an efficient and effective scheduling policy by leveraging expert information and appropriate training strategies. Zhang *et al.* [20] proposed a hierarchical RL framework for visual tracking called PACNet, which consists of the Policy and Actor-Critic Networks.

Full-RL approaches: They formulate a certain task as a Markov decision process in an RL framework, which is more challenging than Partial-RL approaches. For instance, Zhou *et al.* [19] formulated video summarization as a sequential decision making process, and proposed a label-free reward function that jointly explain the diversity and representativeness of generated summaries. Guo *et al.* [44] designed a deep RL based approach, which formulates the object tracking problem as a Markov decision process where state consisted of image regions extracted by the bounding box and eight actions. ICMDP [21] models the binary classification problem as a Markov decision process, and applies the Deep Q-learning algorithm to train the agent. He *et al.* [22] combined Deep Q-learning with curriculum learning to enable the agent to sequentially predict labels based on image feature and previously predicted labels. In [45], the authors casted the object detection problem as a Markov decision process. The agent is utilized to find a region of interest in the image first, and then reduce the region of interest to find the smaller region based on the previously selected region, thus forming a hierarchy. Our Rein-PAR falls into this group, which formulates PAR as a Markov decision process in a Deep Q-learning framework.

3 Method

Fig. 3 is the overall framework of our Rein-PAR. In this section, we first introduce in detail how we formulate PAR as a Markov decision process. Then, we present how to train the agent using Deep Q-learning [46]. After that, the grouping of attributes is detailedly described. Finally, we introduce our proposed group optimization reward function.

3.1 PAR as Markov Decision Process

We formulate PAR as a Markov decision process, which can be addressed with reinforcement learning algorithms. The agent is used to determine whether attributes exist, whose goal is to learn a policy to maximize the total discount reward, a larger discount reward means a higher the accuracy of recognition. We define the feedbacks from the environment as reward to the agent, and regard the recognition process of an image as an episode.

Markov decision process can be described by a five-tuple $(S, A, R, \mathcal{T}, \gamma)$, where S is the state space, A is a limited set of actions, R represents the reward brought by the transition from state s to state s' via action a , \mathcal{T} represents the transition from the previous state to the next state, and γ is a discount factor. Their details are as follows.

State: State refers to the information of the observed environment. In practical situations, the information that can be observed is very limited due to cognitive limitations. In our approach, we represent it as a tuple, including image features and attribute information, expressed as :

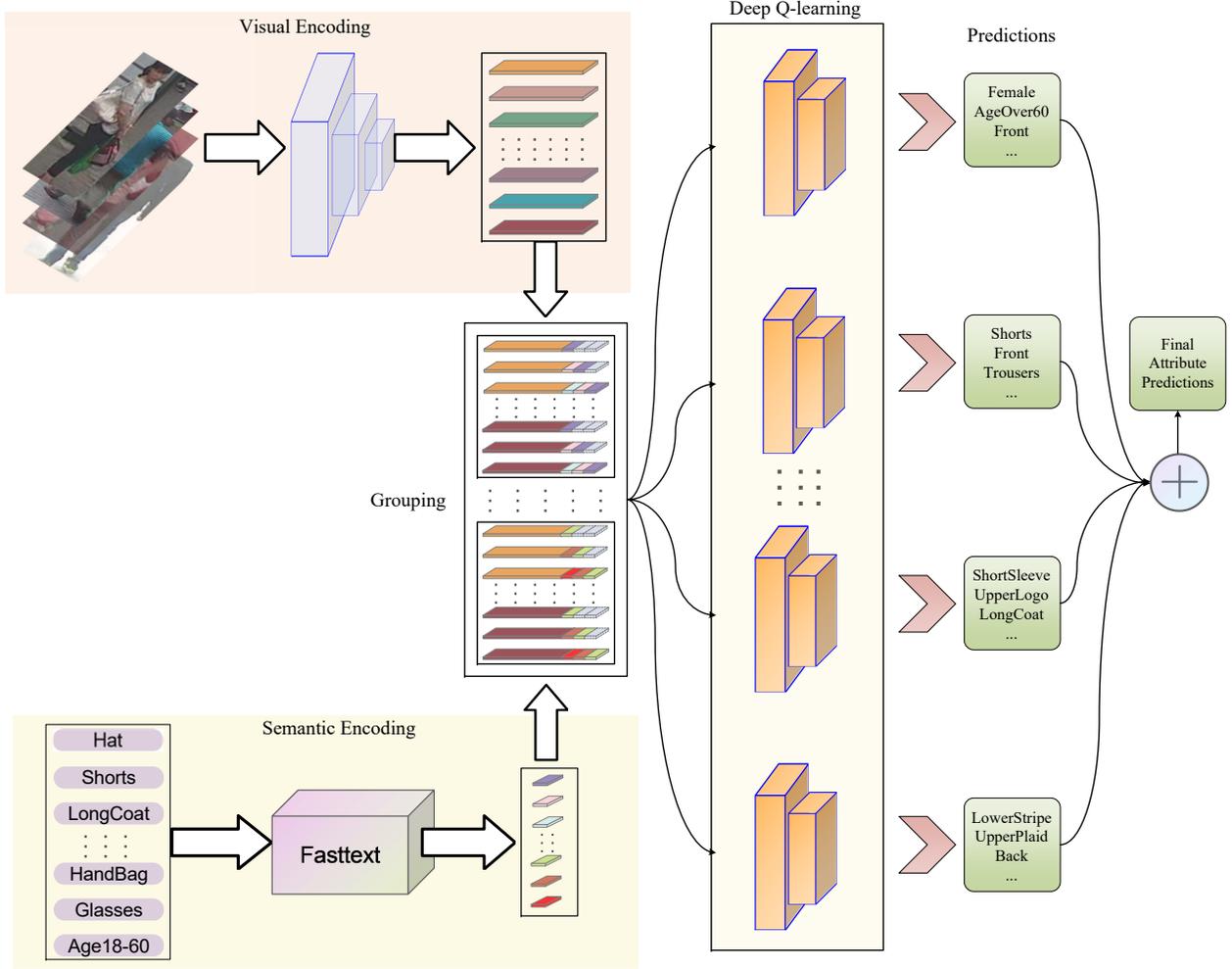


Figure 3: Architecture of our proposed Rein-PAR model.

$$s = (f, v), \quad (1)$$

where $s \in S$ summarizes the information of the observed image and the attribute, f is the feature of the current image, and v is the attribute vector. The 2048-dimensional image features are extracted by the popular Resnet-50 model [47], which is pre-trained model on ImageNet [48]. Since the object categories in ImageNet are different from those in PAR task, we fine-tuned it on each PAR training dataset.

The attribute information refers to the attributes to be predicted at the current time t and those have been predicted at the time $t - 1$ and $t - 2$. It is represented by a fasttext-encoded vector [49]. For each attribute, we encode it into an L -dimensional vector, where L is the number of attributes in the dataset.

Action: Action refers to the behavior issued by the agent and the interaction between the agent and the environment. RMIC [22] made an attempt by directly utilizing label set as the action space, which is likely to cause overestimation problem. Usually, a smaller action space avoids the overestimation problem. By doing so, we define the agent's action space as a set $A = \{0, 1\}$, which determines whether the corresponding attribute exists. Particularly, the element 0 represents that there is no corresponding attribute, and that of 1 represents the existence of the corresponding attribute. In the training stage, the agent interacts with the environment after taking action, and obtains the corresponding positive or negative rewards feedback from the environment. For each image, the number of actions depends on the number of labels, and all actions corresponding to each image are the final predicted labels of the image.

Reward: Reward is a significant research direction in RL, which refers to the value that an agent receives after conducting an action to describe whether its behavior is good or bad. In our approach, a basic reward could be set as +1

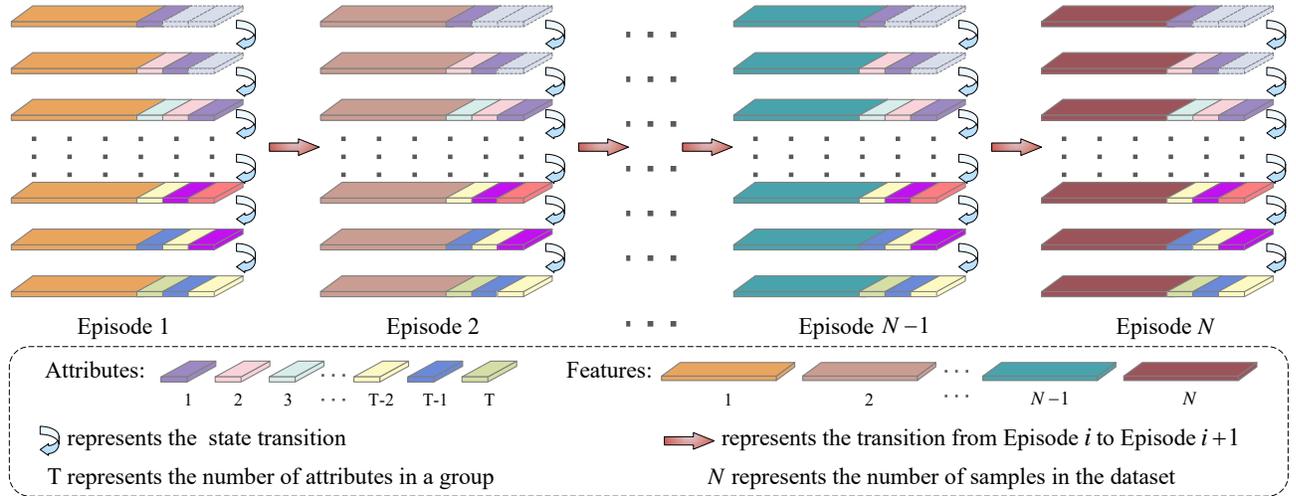


Figure 4: The state transition process of Rein-PAR. The concept of "group" will be presented in the attribute grouping strategy.

and -1 , that is, when the value of the action is the same as the label, a positive reward is given, otherwise, a negative reward is given. It is expressed as:

$$r = \begin{cases} +1 & \text{if } a_t = l_t \\ -1 & \text{if } a_t \neq l_t \end{cases}, \quad (2)$$

where a_t is the action of the agent in state s_t , l_t is the label of the attribute on the sample in state s_t .

However, this design ignores the intra-attribute imbalance problem, an elaborate reward will be introduced in the fourth part of this section.

Transitions: After an action is performed, the current state is transformed to the next state. In our approach, the MDP transitions are deterministic, that is, for each state there is a specified next state, the action has no affect on the next state. The specific state transfer formula is as follows:

$$\mathcal{T}(s, a) = \mathcal{T}((f, v), a) = (f, v') = s', \quad (3)$$

The state transition process is shown in Fig. 4.

Discount factor: Discount factor $\gamma \in [0, 1]$ is to balance the relationship between immediate reward and future reward. When its value approaches 0, the agent focuses more on short-term return, while more long-term return is considered when it is closed to 1.

3.2 Attributes Grouping Strategy

There are strong correlations among the pedestrian attributes, especially those in the same region are often mutually exclusive. For example, "BoldHair" and "LongHair" cannot appear at the same time. We apply the grouping idea to group attributes according to their position and characteristic. Then, we make a separate prediction for each group. Through grouping, the correlation among attributes in the same group is more fully utilized. Moreover, the grouping strategy enables the distribution of attributes in the same group relatively balanced, which alleviates the inter-imbalance attributes distribution problem to some extent. The attributes groups are shown in Tables 1, 2 and 3:

3.3 Group Optimization Reward Function

Reinforcement learning algorithms are usually sensitive to the reward function, which is one of the key parts in RL. The typical one shown in Eq. (2) treats all attributes equally, which does not meet the practical PAR requirement. We could observe from Fig. 2 that there are only five attributes in RAP dataset and three ones in PA100K dataset are common (the

Table 1: The groups of the PETA dataset, which are obtained according to the position and characteristic of the attributes.

Group	Attribute
Age, Gender	personalLess30, personalLess45, personalMale, etc.
Head	Hat, Muffler, Sunglasses, hairLong
Upper Body	Casual, Formal, Jacket, Logo, Plaid, ShortSleeve, etc.
Lower Body	Formal, Jeans, Shorts, ShortSkirt, Trousers, etc.
Footwear	LeatherShoes, Sandals, Shoes, Sneaker
Carrying	Backpack, Other, MessengerBag, Nothing, PlasticBags

Table 2: The groups of the RAP dataset, which are obtained according to the position and characteristic of the attributes.

Group	Attribute
Age, Gender	AgeLess16, Age17-30, Age31-45, Female, BodyFat,
Bodyshape, Role	BodyNormal, BodyThin, Customer, Clerk
Head	BaldHead, LongHair, BlackHair, Hat, Glasses, Muffler
Upper Body	Shirt, Sweater, Vest, TShirt, Cotton, Jacket, SuitUp, etc.
Lower Body	LongTrousers, Skirt, ShortSkirt, Dress, Jeans, etc.
Footwear	Leather, Sport, Boots, Cloth, Casual
Accessory	Backpack, SingleShoulderBag, HandBag, Box, etc.
Action	Calling, Talking, Gathering, Holding, Pushing, etc.

Table 3: The groups of the PA100K dataset, which are obtained according to the position and characteristic of the attributes.

Group	Attribute
Age, Gender	AgeOver60, Age18-60, AgeLess18, Female
Location, Head	Front, Side, Back, Hat, Glasses
Upper Body	ShortSleeve, LongSleeve, UpperStride, UpperLogo, etc.
Lower Body, Foot	LowerStripe, LowerPattern, LongCoat, Trousers, Shorts, Skirt&Dress, boots
Accessory	HandBag, ShoulderBag, Backpack, HoldObjectsInFront

presence ration larger than 0.5), while most attributes in the PAR are less common. Thus, the appearance of an attribute is actually an imbalanced data distribution problem. This intra-attribute imbalanced distribution enforces the training of the network biased towards the common attributes. To alleviate this problem, we propose a Group Optimization Reward (GOR) function to pay more attention to those less common attributes, which optimizes the reward function on the basis of grouping. It is designed as follows:

$$r' = \begin{cases} 1 & a_t = 1 \& l_t = 1 \\ -1 & a_t = 0 \& l_t = 1 \\ -\rho & a_t = 1 \& l_t = 0 \\ \rho & a_t = 0 \& l_t = 0 \end{cases}, \quad (4)$$

$$\rho = \begin{cases} 0.15 & c \in [0, 0.05) \\ 0.25 & c \in [0.05, 0.25) \\ 0.35 & c \in [0.25, 0.35) \\ 0.45 & c \in [0.35, 0.45) \\ 0.55 & c \in [0.45, 1) \end{cases}, \quad (5)$$

$$c = \frac{\sum_T n_T}{T * N - \sum_T n_T}, \quad (6)$$

where c refers to the imbalanced coefficient, T is the number of attributes in the group, N is the size of the dataset, n_T is the number of images containing a certain attribute, a_t is the action of the agent in state s_t , and l_t is the label of the attribute on the sample in state s_t . It should be noticed that c is the ratio of the sum of the attributes presence to the sum of the attributes absence in a group, which reflects the overall situation of each group of attributes.

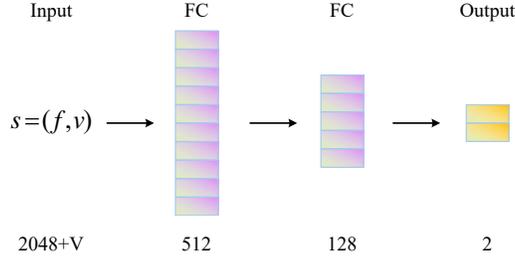


Figure 5: The architecture of our Deep Q-learning network.

It could be observed that the reward function is assigned a “+1” or “-1” feedback in case the attribute is present, and a ρ or $-\rho$ feedback for the opposite case. In this way, the less common attributes generally receive a smaller feedback in the case of its absence in a sample. Since the absence of the less common attributes predominate in the samples, it alleviates the intra-attribute imbalance problem existing in most attributes. However, it should be pointed out that the reward function will have a negative influence on the recognition of the common attributes, but for the overall recognition result, its influence is positive.

3.4 Deep Q-learning for Rein-PAR

The optimal policy to maximize discount rewards can be obtained via RL. In Rein-PAR, we resort to the Deep Q-learning [46] algorithm to train our agent, as shown in Fig. 5. It consists of three fully connected layers, 512, 128 and 2, respectively. The input is the state composed of image features and attribute-encoded vectors. Our training goal is to enable the agent to recognize attributes as accurately as possible.

During the training process, we store the transition information in the replay memory in the form of (s, a, r, s') . It should be noted that the use of replay memory allows the agent to train without considering the order of images and attributes. When the policy network is updated, a mini batch randomly sampled from the replay memory is applied for training. The loss function is the (Mean Square Error) MSE loss as follow:

$$L(\theta) = \mathbb{E}_{(s,a,r,s') \sim D} \left[\left(r + \hat{Q}(s', a'; \theta^-) - Q(s, a; \theta) \right)^2 \right], \quad (7)$$

where $r + \hat{Q}(s', a'; \theta^-)$ is from the target network, θ^- is the parameter in the target network, $Q(s, a; \theta)$ is the output value in the policy network, and θ is the weights in the policy network. We use Adam optimization algorithm to update network parameters.

The Deep Q-learning [46] for Rein-PAR algorithm is shown as Algorithm 1. This process can be understood as the process of judging all attributes of each image as a game (episode), and each judging whether an attribute exists in an image is regarded as a step in the game, and the epoch refers to how many times a dataset is trained. The purpose of training is to enable the agent to judge all attributes of an image as correctly as possible.

4 Experiment

We first briefly depict the experiment setup, including the datasets, implementation details and evaluation metrics. Then we show and analyze the experimental results. Finally, we conduct ablation experiments to analyze the impacts of attributes grouping strategy and group optimization reward function.

4.1 Experiment Setup

Datasets: We verify the effectiveness of the proposed approach on three benchmark datasets, namely PETA [12], PA100K [14] and RAP [13]. Concretely, PETA consists of 8,705 pedestrians and 19,000 images (resolution ranges from 17×39 to 169×365). The training, validation and test set of this dataset include 9,500, 1,900 and 760 images respectively. PA100K is the largest dataset used for PAR so far, which contains a total of 100,000 pedestrian images collected from outdoor surveillance cameras, each image has 26 popular attributes. The entire dataset is randomly divided into 80,000 training images, 10,000 validation images and 10,000 test images. RAP contains 41,585 pedestrian images from indoor scenes, of which 33,268 images are employed for training and 8,317 for testing.

Algorithm 1 Deep Q-learning for Rein-PAR**Initialization:**

Initialize Replay Memory D to capacity B ;
 Initialize the whole action space A ;
 Initialize action-value function Q with random weights θ ;
 Initialize target action-value function \hat{Q} with random weights $\theta^- = \theta$;

Train: the training process for Rein-PAR;

for epoch = 1, M **do**

for each image **do**

 Initialize the current state s_0 with the current image
 feature f and semantic vector v_i ;

for t = 1, T **do**

 Select action a_t from A with ε - greedy policy

 Execute a_t , observe reward r_t , next state s_{t+1}

 Store transition (s_t, a_t, r_t, s_{t+1}) in D

 Sample random minibatch of transitions

(s_j, a_j, r_j, s_{j+1}) from D

 Set $y_j = \begin{cases} r_j & \text{for terminal } s_{j+1} \\ r_j + \gamma \max_{a'} \hat{Q}(s_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$

 Perform a gradient descent step on $(y_j - Q(s_{j+1}, a_j; \theta))^2$ with respect to the network parameters θ ;

 Update $s_{t+1} \leftarrow s_t$

 Update $\hat{Q} = Q$ every C step

end for

end for

end for

Implementation Details: We train the network for 15 epochs, and set the capacity of replay memory D to 2000, update frequency C of the target network to 100 and the mini batch size to 64. The discount factor γ is set to 0.9. The ε -greedy strategy is utilized on action selection. During the training stage, the probabilities of exploration are gradually decreased from 0.9 to 0.05. During the test stage, we set the probability of exploration is to 0. The network is optimized by Adam optimizer and the default parameter Settings are adopted. Our approach is implemented on the publicly available Pytorch platform on a single NVIDIA GeForce GTX 1080Ti GPU with 12GB of memory.

Evaluation Metrics: We evaluate the performance of our approach on five metrics, as shown in Eqs.(8)-(12). Among them, mean accuracy (mA) is a label-based metric, which treats each attribute independently. It first calculates the classification accuracy of positive and negative samples for each attribute, and then averages their values as the recognition result for the attributes. Accuracy (Acc), precision (Prec), recall rate (Rec) and F1 score are instance-based metrics.

$$mA = \frac{1}{2L} \sum_{i=1}^L \left(\frac{TP_i}{P_i} + \frac{TN_i}{N_i} \right), \quad (8)$$

$$ACC = \frac{1}{N} \sum_{i=1}^N \frac{|Y_i \cap f(x_i)|}{|Y_i \cup f(x_i)|}, \quad (9)$$

$$Prec = \frac{1}{2N} \sum_{i=1}^N \frac{|Y_i \cap f(x_i)|}{|f(x_i)|}, \quad (10)$$

$$Rec = \frac{1}{2N} \sum_{i=1}^N \frac{|Y_i \cap f(x_i)|}{|Y_i|}, \quad (11)$$

$$F1 = \frac{2 * Prec * Rec}{Prec + Rec}, \quad (12)$$

where L is the number of attributes, N is the number of samples, P_i and N_i are the number of positive and negative samples, TP_i and TN_i are the number of correctly predicted positive and negative samples, Y_i is the ground truth positive labels of the i -th sample, and $f(x_i)$ is the predicted positive labels of the i -th sample.

4.2 Comparison with State-of-the-Art Methods

Table 4: Comparison results on PETA dataset. The best results are marked in boldface.

Method	mA	Acc	Prec	Rec	F1
ACN [50]	81.15	73.66	84.06	81.26	82.64
DeepMAR [27]	82.89	75.07	83.68	83.14	83.41
SR [51]	82.83	-	82.54	82.76	82.65
CTX [29]	80.13	-	79.68	80.24	79.68
WPAL [52]	85.50	76.98	84.07	85.78	84.90
HP-Net [14]	81.77	76.13	84.92	83.24	84.07
VeSPA [53]	83.45	77.73	86.18	84.81	85.49
PGDM [31]	82.97	78.08	86.86	84.68	85.76
IA ² -Net [17]	84.13	78.62	85.73	86.07	85.88
JLPLS-PAA [54]	84.88	79.46	87.42	86.33	86.87
MT-CAS [55]	83.17	78.78	87.49	85.35	86.41
MTA-Net [15]	84.62	78.80	85.67	86.42	86.04
Rein-PAR (Ours)	85.51	78.45	84.08	88.77	85.91

Table 5: Comparison results on RAP dataset. The best results are marked in boldface.

Method	mA	Acc	Prec	Rec	F1
ACN [50]	69.66	62.61	80.12	72.26	75.98
DeepMAR [27]	73.79	62.02	74.92	76.21	75.56
CTX [29]	70.13	-	71.03	71.20	70.23
JRL [16]	77.81	-	78.11	78.98	78.58
SR [29]	74.10	-	75.11	76.52	75.83
WPAL [52]	81.25	50.30	57.17	78.39	66.12
HP-Net [14]	76.12	65.39	77.33	78.79	78.05
VeSPA [53]	77.70	67.35	79.51	79.67	79.59
PGDM [31]	74.31	64.57	78.86	75.90	77.35
GSR-MAR [28]	67.76	63.44	82.27	71.82	76.69
IA ² -Net [17]	77.44	65.75	79.01	77.45	78.03
JLPLS-PAA [54]	81.25	67.91	78.56	81.45	79.98
MSE-Net [56]	71.12	62.43	78.82	72.94	75.77
MTA-Net [15]	77.62	67.17	79.72	78.44	79.07
HR-Net [57]	81.10	45.70	51.48	78.56	62.20
VALA [58]	78.33	67.48	79.81	80.84	80.32
Rein-PAR (Ours)	81.67	66.24	73.24	85.80	78.68

We choose eighteen state-of-the-art approaches for comparison, which include global-based models such as ACN [50], DeepMAR [27], GSR-MAR [28], SR [51], CTX [29], MSE-Net [56], and HR-Net [57]; part-based models such as PGDM [31] and LGNet [32]; attention-based models such as JRL [16], HP-Net [14], VeSPA [53], VALA [58], JLPLS-PAA [54], and MT-CAS [55]; and loss function based models such as WPAL [52], MTA-Net [15], and IA²-Net [17]. The results on the three datasets are shown in Table 4, Table 5 and Table 6, respectively. We have the following observations:

Table 6: Comparison results on PA100K dataset. The best results are marked in boldface.

Method	mA	Acc	Prec	Rec	F1
DeepMAR [27]	72.70	70.39	82.24	80.42	81.32
HP-Net [14]	74.21	72.19	82.97	82.09	82.53
VeSPA [53]	76.32	73.00	84.33	81.49	83.20
PGDM [31]	74.95	73.08	84.36	82.24	83.29
LGNet [32]	76.96	75.55	86.99	83.17	85.04
GSR-MAR [28]	72.43	73.46	87.68	79.94	83.63
MT-CAS [55]	77.20	78.09	88.46	84.86	86.62
VALA [58]	80.08	78.14	87.60	86.73	87.16
Rein-PAR (Ours)	80.55	77.20	84.76	87.67	85.70

- 1) It could be observed that our Rein-PAR achieves a competitive performance on the three datasets. It should be noted that most approaches reported their results on only two datasets due to the diverse data challenge, and it is really hard to achieve the best performance on all the five metrics.
- 2) Compared with the competitors, Rein-PAR is the best on the metrics of mA and Rec on the three datasets. For example, on the RAP dataset, the mA reaches 81.67%, 0.42% higher than the second-ranked approach WPAL, and 4.35% higher than the second-ranked VeSPA approach on the Rec metric. The highest mA indicates that our approach is better than the competitors in the recognition accuracy of a single attribute. The highest Rec indicates that our approach has a higher accuracy in discriminating samples that are indeed positive. The main reason is that we apply GOR to pay more attention to those less common attributes.
- 3) Rein-PAR is also comparable with the competitors on the other metrics. Specifically, our approach outperforms most approaches on the metrics Acc and F1. Our approach is inferior on the Prec metric, whose reason lies in that our approach pays less attention to negative samples. This is the direction that we need to improve in our future work.

Fig. 6 shows some recognition results on the PETA dataset. It consists of the image, the corresponding Ground Truth, and the recognition results. Firstly, we can observe that the number of attributes of each pedestrian image in the dataset is not necessarily the same. Secondly, it can be observed that the recognition errors are often caused by illumination, low resolution, blur, multiple pedestrian targets, etc. Take Fig. 6 (e) as an example, our approach wrongly considers that the female has three attributes, namely "carryingNothing", "lowerBodytrousers" and "footwearSneaker" due to the interference of surrounding pedestrians. Thirdly, the redundant attributes are often closely related to pedestrian images. For example, the male in Fig. 6 (d) may actually has the attributes of "carryingOther" and "lowerBodyShorts", although they are not annotated in the ground-truth. It should be noted that these unannotated categories will be regarded as negative categories during training, which will interfere the training stage.

4.3 Ablation Studies

We analyze the effectiveness of each part of our proposed Rein-PAR in Table 7. We first consider the following variants:

Baseline refers to the approach that does not apply the approaches of Attributes Grouping Strategy (AGS) and Group Optimization Reward (GOR).

Baseline+AGS is the approach applying AGS for prediction based on RPAR.

Rein-PAR is the proposed approach, which incorporates both AGS and GOR.

It can be observed that even the Baseline approach could achieve satisfactory results. For example, it is better than GSR-MAR and ACN on the RAP dataset, and also better than DeepMAR, HP-Net and GSR-MAR on the PA100K dataset. Furthermore, it is observed that the Baseline+AGS approach outperforms the Baseline by 3.08% on PETA, 5.78% on RAP and 3.9% on PA100K on mA metric. And on the F1 metric, the Baseline+AGS approach has improvements of 3.29% on PETA, 2.68% on RAP and 1.32% on PA100K compared with Baseline. It proves the effectiveness of the grouping attributes strategy. Additionally, Rein-PAR further applies group optimization reward to alleviate the intra-attribute imbalance problem, which brings improvements of 2.27%, 5.99%, and 1.99% on mA on the three datasets compared with Baseline+AGS. In addition to the significant improvement under mA metric, other metrics such as Acc and Rec have also been significantly improved. For instance, the Acc rises from 77.86% to 78.45% on PETA. It's worth noting that the utilization of GOR reduces attention on negative samples, resulting in a decline on the Prec metric.

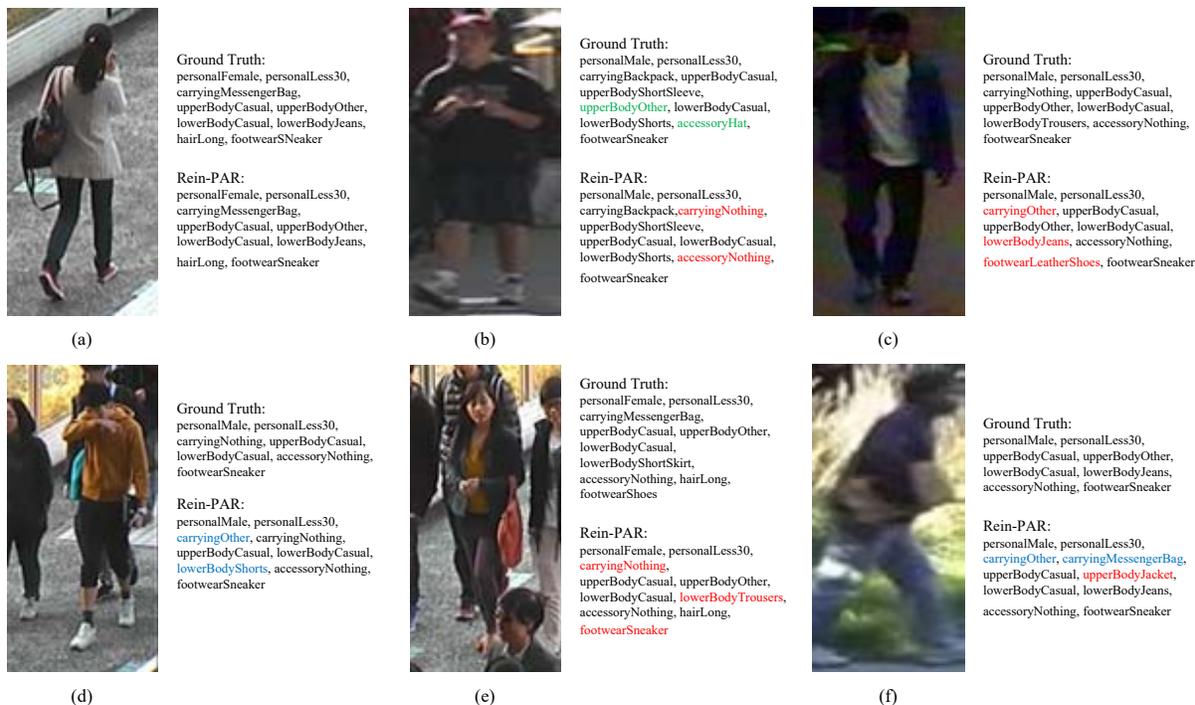


Figure 6: Pedestrian attribute recognition results by Rein-PAR on the PETA dataset, where the black font indicates the correct recognition result, the red indicates the incorrectly recognized attribute, the blue is the redundantly recognized attribute, and the green represents the unrecognized attribute. Among the six pedestrian images, (a) is an example where the recognition results are completely correct, and (b-f) are examples where the recognition results are defective. Note that the images are blurred due to the low image resolution in the dataset.

Table 7: Ablation studies of Rein-PAR on PETA, RAP and PA100K datasets. The best results are marked in boldface.

Dataset	Method	mA	Acc	Prec	Rec	F1
PETA	Baseline	80.06	72.63	84.28	81.54	82.20
	Baseline+AGS	83.14	77.86	86.02	85.55	85.49
	Rein-PAR	85.51	78.45	84.08	88.77	85.91
RAP	Baseline	69.90	63.71	79.53	75.05	76.35
	Baseline+AGS	75.68	67.26	80.94	78.02	79.03
	Rein-PAR	81.67	66.24	73.24	85.80	78.68
PA100K	Baseline	74.66	74.49	85.78	82.87	83.72
	Baseline+AGS	78.56	76.41	86.63	84.38	85.04
	Rein-PAR	80.55	77.20	84.76	87.67	85.70

4.4 Impact of ρ in the Reward Function

It is well-known that different reward functions have different impacts on performance. We select two groups of attributes on the PETA and RAP datasets and test the impact of different values of ρ in the reward function. Table 8 shows the two groups of attributes and the corresponding attribute ratios.

We test the performance when the ρ is from 0.05 to 1, with the interval is 0.05. The impact of ρ is visualized in terms of Acc, F1 and mA, as shown in Fig. 7. It should be noted that, because the difficulty of recognizing each group of attributes is different, some of the metrics in some groups are lower. For example, as shown in Fig. 7(b), the attribute group selected from the RAP dataset are mostly uncommon attributes, which are more difficult to recognize, so the Acc and F1 metrics are low. We can observe that value of ρ has a significant effect on performance, and the optimal ρ values corresponding to different groups are different. We could also observe that the influence trends of ρ on the three metrics

Table 8: Two selected groups of attributes in the PETA and RAP datasets and their corresponding attribute ratios.

PETA	Attribute	LeatherShoes	Sandals	Shoes	Sneaker
	Ratio	0.296	0.02	0.363	0.216
RAP	Attribute	Calling	Talking	Gathering	Holding
	Ratio	0.034	0.032	0.092	0.025
	Attribute	Pushing	Pulling	CarrybyArm	CarrybyHand
	Ratio	0.01	0.018	0.023	0.129



(a) on PETA dataset



(b) on RAP dataset

Figure 7: The impact of ρ on Acc, F1 and mA.

of Acc, F1 and mA are almost the same. The value of the appropriate ρ appears at the peak of the evaluation metrics curves.

5 Conclusion

This paper has proposed the Rein-PAR approach for addressing the PAR task. Different from previous approaches, Rein-PAR defines PAR as a Markov decision-making process for the first time, and employs the Deep Q-learning algorithm to train the network to recognize attributes. Moreover, an attribute grouping strategy is applied to alleviate inter-attribute imbalance problem, and a group optimization reward function is further developed to alleviate the intra-attribute imbalance problem. Experimental results on PETA, RAP and PA100K datasets have demonstrated the effectiveness of our approach. Rein-PAR is a successful attempt to apply reinforcement learning on PAR, which demonstrates that reinforcement learning has considerable potential in PAR. In the future, we consider to construct a

more appropriate Markov decision process, a reward function that is more in line with the actual situation of PAR task, and utilize more advanced reinforcement learning algorithms for training.

References

- [1] Lixuan Yi, Qian Zhao, Wei Wei, and Zongben Xu. Robust online rain removal for surveillance videos with dynamic rains. *Knowledge-Based Systems*, 222:107006, 2021.
- [2] Fanglin Chen, Weihang Wang, Huiyuan Yang, Wenjie Pei, and Guangming Lu. Multiscale feature fusion for surveillance video diagnosis. *Knowledge-Based Systems*, page 108103, 2022.
- [3] Yinsong Xu, Zhuqing Jiang, Aidong Men, Haiying Wang, and Haiyong Luo. Multi-view feature fusion for person re-identification. *Knowledge-Based Systems*, 229:107344, 2021.
- [4] Qiang Liu, Xiaohai He, Mozhi Zhang, Qizhi Teng, Bo Li, and Linbo Qing. Feature separation and double causal comparison loss for visible and infrared person re-identification. *Knowledge-Based Systems*, 239:108042, 2022.
- [5] Shanshan Wang, Lei Zhang, Weihua Chen, Fan Wang, and Hao Li. Refining pseudo labels for unsupervised domain adaptive re-identification. *Knowledge-Based Systems*, page 108336, 2022.
- [6] Honghu Pan, Yang Bai, Zhenyu He, and Chunkai Zhang. Aagcn: Adjacency-aware graph convolutional network for person re-identification. *Knowledge-Based Systems*, 236:107300, 2022.
- [7] Zhong Ji and Shengjia Li. Multimodal alignment and attention-based person search via natural language description. *IEEE Internet of Things Journal*, 7(11):11147–11156, 2020.
- [8] Yaqing Zhang, Xi Li, and Zhongfei Zhang. Efficient person search via expert-guided knowledge distillation. *IEEE Transactions on Systems, Man, and Cybernetics*, pages 1–12, 2019.
- [9] Beibei Yang, Changqian Yu, Jin-Gang Yu, Changxin Gao, and Nong Sang. Pose-guided hierarchical semantic decomposition and composition for human parsing. *IEEE Transactions on Cybernetics*, pages 1–12, 2021.
- [10] Lu Ding, Yong Wang, Robert Laganière, Dan Huang, Xinbin Luo, and Huanlong Zhang. A robust and fast multispectral pedestrian detection deep network. *Knowledge-Based Systems*, 227:106990, 2021.
- [11] Zhi-Ri Tang, Ruihan Hu, Yanhua Chen, Zhao-Hui Sun, and Ming Li. Multi-expert learning for fusion of pedestrian detection bounding box. *Knowledge-Based Systems*, page 108254, 2022.
- [12] Yubin Deng, Ping Luo, Chen Change Loy, and Xiaoou Tang. Pedestrian attribute recognition at far distance. In *ACM International Conference on Multimedia*, pages 789–792, 2014.
- [13] Dangwei Li, Zhang Zhang, Xiaotang Chen, Haibin Ling, and Kaiqi Huang. A richly annotated dataset for pedestrian attribute recognition. *arXiv preprint arXiv:1603.07054*, 2016.
- [14] Xihui Liu, Haiyu Zhao, Maoqing Tian, Lu Sheng, Jing Shao, Shuai Yi, Junjie Yan, and Xiaogang Wang. Hydraplusnet: Attentive deep features for pedestrian analysis. In *International Conference on Computer Vision*, pages 350–359, 2017.
- [15] Zhong Ji, Zhenfei Hu, Erlu He, Jungong Han, and Yanwei Pang. Pedestrian attribute recognition based on multiple time steps attention. *Pattern Recognition Letters*, 138:170–176, 2020.
- [16] Jingya Wang, Xiatian Zhu, Shaogang Gong, and Wei Li. Attribute recognition by joint recurrent learning of context and correlation. In *International Conference on Computer Vision*, pages 531–540, 2017.
- [17] Zhong Ji, Erlu He, Haoran Wang, and Aiping Yang. Image-attribute reciprocally guided attention network for pedestrian attribute recognition. *Pattern Recognition Letters*, 120:89–95, 2019.
- [18] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneshelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- [19] Kaiyang Zhou, Yu Qiao, and Tao Xiang. Deep reinforcement learning for unsupervised video summarization with diversity-representativeness reward. In *AAAI Conference on Artificial Intelligence*, volume 32, pages 7582–7589, 2018.
- [20] Dawei Zhang, Zhonglong Zheng, Riheng Jia, and Minglu Li. Visual tracking via hierarchical deep reinforcement learning. In *AAAI Conference on Artificial Intelligence*, volume 35, pages 3315–3323, 2021.
- [21] Enlu Lin, Qiong Chen, and Xiaoming Qi. Deep reinforcement learning for imbalanced classification. *Applied Intelligence*, 50(8):2488–2502, 2020.

- [22] Shiyi He, Chang Xu, Tianyu Guo, Chao Xu, and Dacheng Tao. Reinforced multi-label image classification by exploring curriculum. In *AAAI Conference on Artificial Intelligence*, pages 3183–3190, 2018.
- [23] Hado Van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-learning. In *AAAI Conference on Artificial Intelligence*, volume 30, pages 2094–2100, 2016.
- [24] Gaurav Sharma and Frederic Jurie. Learning discriminative spatial representation for image classification. In *British Machine Vision Conference*, pages 1–11, 2011.
- [25] Ryan Layne, Timothy M Hospedales, Shaogang Gong, and Q Mary. Person re-identification by attributes. In *British Machine Vision Conference*, volume 2, page 8, 2012.
- [26] Ryan Layne, Timothy M Hospedales, and Shaogang Gong. Towards person identification and re-identification with attributes. In *European Conference on Computer Vision*, pages 402–412. Springer, 2012.
- [27] Dangwei Li, Xiaotang Chen, and Kaiqi Huang. Multi-attribute learning for pedestrian attribute recognition in surveillance scenarios. In *Asian Conference on Pattern Recognition*, pages 111–115, 2015.
- [28] Thomhert Suprpto Siadari, Mikyong Han, and Hyunjin Yoon. Gsr-mar: Global super-resolution for person multi-attribute recognition. In *IEEE International Conference on Computer Vision Workshops*, pages 1098–1103, 2019.
- [29] Yao Li, Guosheng Lin, Bohan Zhuang, Lingqiao Liu, Chunhua Shen, and Anton van den Hengel. Sequential person recognition in photo albums with a recurrent network. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1338–1346, 2017.
- [30] Zichang Tan, Yang Yang, Jun Wan, Guodong Guo, and Stan Z Li. Relation-aware pedestrian attribute recognition with graph convolutional networks. In *AAAI Conference on Artificial Intelligence*, volume 34, pages 12055–12062, 2020.
- [31] Dangwei Li, Xiaotang Chen, Zhang Zhang, and Kaiqi Huang. Pose guided deep model for pedestrian attribute recognition in surveillance scenarios. In *IEEE International Conference on Multimedia and Expo*, pages 1–6, 2018.
- [32] Pengze Liu, Xihui Liu, Junjie Yan, and Jing Shao. Localization guided learning for pedestrian attribute recognition. *arXiv preprint arXiv:1808.09102*, 2018.
- [33] Jiajun Zhang, Pengyuan Ren, and Jianmin Li. Deep template matching for pedestrian attribute recognition with the auxiliary supervision of attribute-wise keypoints. *arXiv preprint arXiv:2011.06798*, 2020.
- [34] Chufeng Tang, Lu Sheng, Zhaoxiang Zhang, and Xiaolin Hu. Improving pedestrian attribute recognition with weakly-supervised multi-scale attribute-specific localization. In *International Conference on Computer Vision*, pages 4997–5006, 2019.
- [35] Nikolaos Sarafianos, Xiang Xu, and Ioannis A Kakadiaris. Deep imbalanced attribute classification using visual attention aggregation. In *European Conference on Computer Vision*, pages 680–697, 2018.
- [36] Mingda Wu, Di Huang, Yuanfang Guo, and Yunhong Wang. Distraction-aware feature learning for human attribute recognition via coarse-to-fine attention mechanism. In *AAAI Conference on Artificial Intelligence*, volume 34, pages 12394–12401, 2020.
- [37] Yang Yang, Zichang Tan, Prayag Tiwari, Hari Mohan Pandey, Jun Wan, Zhen Lei, Guodong Guo, and Stan Z Li. Cascaded split-and-aggregate learning with feature recombination for pedestrian attribute recognition. *International Journal of Computer Vision*, 129(10):2731–2744, 2021.
- [38] Ngan Le, Vidhiwar Singh Rathour, Kashu Yamazaki, Khoa Luu, and Marios Savvides. Deep reinforcement learning in computer vision: A comprehensive survey. *arXiv preprint arXiv:2108.11510*, 2021.
- [39] Mei Wang and Weihong Deng. Mitigating bias in face recognition using skewness-aware reinforcement learning. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 9322–9331, 2020.
- [40] Wenkai Dong, Zhaoxiang Zhang, and Tieniu Tan. Attention-aware sampling via deep reinforcement learning for action recognition. In *AAAI Conference on Artificial Intelligence*, volume 33, pages 8247–8254, 2019.
- [41] Jianfu Zhang, Li Niu, and Liqing Zhang. Person re-identification with reinforced attribute attention selection. *IEEE Transactions on Image Processing*, 30:603–616, 2020.
- [42] Yuxuan Shi, Zhen Wei, Hefei Ling, Ziyang Wang, Pengfei Zhu, Jialie Shen, and Ping Li. Adaptive and robust partition learning for person retrieval with policy gradient. *IEEE Transactions on Multimedia*, 14:1–14, 2020.
- [43] Yujiang Wang, Mingzhi Dong, Jie Shen, Yang Wu, Shiyang Cheng, and Maja Pantic. Dynamic face video segmentation via reinforcement learning. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 6959–6969, June 2020.

- [44] Minghao Guo, Jiwen Lu, and Jie Zhou. Dual-agent deep reinforcement learning for deformable face tracking. In *European Conference on Computer Vision*, pages 768–783, 2018.
- [45] Míriam Bellver, Xavier Giró Nieto, Fernando Marqués Acosta, Jordi Torres, et al. Hierarchical object detection with deep reinforcement learning. *Deep Learning for Image Processing Applications*, 31(164):3, 2017.
- [46] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fiedjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- [47] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- [48] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015.
- [49] Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. Bag of tricks for efficient text classification. *arXiv preprint arXiv:1607.01759*, 2016.
- [50] Patrick Sudowe, Hannah Spitzer, and Bastian Leibe. Person attribute recognition with a jointly-trained holistic cnn model. In *IEEE International Conference on Computer Vision Workshops*, pages 87–95, 2015.
- [51] Feng Liu, Tao Xiang, Timothy M Hospedales, Wankou Yang, and Changyin Sun. Semantic regularisation for recurrent image annotation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 2872–2880, 2017.
- [52] Yang Zhou, Kai Yu, Biao Leng, Zhang Zhang, Dangwei Li, Kaiqi Huang, Bailan Feng, and Chunfeng Yao. Weakly-supervised learning of mid-level features for pedestrian attribute recognition and localization. pages 1–12. *British Machine Vision Conference*, 2017.
- [53] M Saquib Sarfraz, Arne Schumann, Yan Wang, and Rainer Stiefelhagen. Deep view-sensitive pedestrian attribute inference in an end-to-end model. *arXiv preprint arXiv:1707.06089*, 2017.
- [54] Zichang Tan, Yang Yang, Jun Wan, Hanyuan Hang, Guodong Guo, and Stan Z Li. Attention-based pedestrian attribute analysis. *IEEE transactions on image processing*, 28(12):6126–6140, 2019.
- [55] Haitian Zeng, Haizhou Ai, Zijie Zhuang, and Long Chen. Multi-task learning via co-attentive sharing for pedestrian attribute recognition. In *IEEE International Conference on Multimedia and Expo*, pages 1–6. IEEE, 2020.
- [56] Miaomiao Lou, Zhenxia Yu, Feng Guo, and Xiaoqiang Zheng. Mse-net: Pedestrian attribute recognition using mlsc and se-blocks. In *International Conference on Artificial Intelligence and Security*, pages 217–226. Springer, 2019.
- [57] Haoran An, Hai-Miao Hu, Yuanfang Guo, Qianli Zhou, and Bo Li. Hierarchical reasoning network for pedestrian attribute recognition. *IEEE Transactions on Multimedia*, 23:268–280, 2021.
- [58] Wei-Chen Chen, Xin-Yi Yu, and Lin-Lin Ou. Pedestrian attribute recognition in video surveillance scenarios based on view-attribute attention localization. *Machine Intelligence Research*, pages 1–16, 2022.