

Pedestrian Attribute Recognition: A Survey

Xiao Wang^{1,2}, Shaofei Zheng¹, Rui Yang¹,
Aihua Zheng¹, Zhe Chen³, Jin Tang¹, Bin Luo¹

1. School of Computer Science and Technology, Anhui University, Hefei, Anhui Province, China.

2. Peng Cheng Laboratory, Shenzhen, China

3. School of Computer Science, Faculty of Engineering, The University of Sydney, Australia

Abstract

Pedestrian Attribute Recognition (PAR) is an important task in computer vision community and plays an important role in practical video surveillance. The goal of this paper is to review existing works using traditional methods or based on deep learning networks. Firstly, we introduce the background of pedestrian attribute recognition, including the fundamental concepts and formulation of pedestrian attributes and corresponding challenges. Secondly, we analyze popular solutions for this task from eight perspectives. Thirdly, we discuss the specific attribute recognition, then, give a comparison between deep learning and traditional algorithm based PAR methods. After that, we show the connections between PAR and other computer vision tasks. Fourthly, we introduce the benchmark datasets, evaluation metrics in this community, and give a brief performance comparison. Finally, we summarize this paper and give several possible research directions for PAR. The project page of this paper can be found at: <https://sites.google.com/view/ahu-pedestrianattributes/>.

Keywords: Pedestrian Attribute Recognition; Multi-label Learning; Multi-task Learning; Deep Learning; CNN-RNN

1. Introduction

2 Pedestrian attributes, are humanly searchable semantic descriptions and can
3 be used as soft-biometrics in visual surveillance, with applications in person re-
4 identification, face verification and human identification. Pedestrian attribute recog-
5 nition (PAR) aims at mining the attributes of target person whose image is given.
6 Different from low-level features, such as HOG, LBP or deep features, attributes

7 can be viewed as high-level semantic information which is more robust to view-
8 point changes and viewing condition variations. Hence, many tasks in com-
9 puter vision integrate the attribute information into their algorithms to achieve
10 better performance, such as pedestrian detection (Chen et al., 2021), person re-
11 identification, action recognition and scene understanding. Although many works
12 have been proposed on this topic, however, PAR is still an unsolved problem due
13 to challenging factors, such as view point change, low illumination, low resolu-
14 tion.

15 Traditional pedestrian attribute recognition methods usually focus on devel-
16 oping robust feature representation from the perspectives of hand-crafted features,
17 powerful classifiers or attributes relations. Some milestones including HOG, SIFT,
18 SVM or CRF model. However, the reports on large-scale benchmark evaluations
19 suggest that the performance of these traditional algorithms is far from the re-
20 quirement of realistic applications. Over the past several years, deep learning has
21 achieved an impressive performance due to its success on automatic feature ex-
22 traction using multi-layer nonlinear transformation, especially in computer vision,
23 speech recognition and natural language processing. Many deep learning based
24 pedestrian attribute recognition algorithms have been proposed based on these
25 breakthroughs.

26 Although so many algorithms have been proposed, until now, there exists no
27 work to make a detailed survey, comprehensive evaluation and insightful analy-
28 sis on these attribute recognition algorithms. In this paper, we summarize exist-
29 ing works on pedestrian attribute recognition, including traditional methods and
30 popular deep learning based algorithms, to better understand this direction and
31 help other researchers to quickly capture main pipeline as well as latest research
32 frontier. Specifically speaking, we attempt to address the following several im-
33 portant issues: 1) What is the connection and difference between traditional and
34 deep learning-based pedestrian attribute recognition algorithms? We analyse tra-
35 ditional and deep learning based algorithms from different classification rules,
36 such as part-based, group-based or end-to-end learning; 2) How the pedestrian
37 attributes contribute to other related computer vision tasks? We also review some
38 person attributes guided computer vision tasks, such as person re-identification,
39 human detection, to fully demonstrate the effectiveness and widely applications
40 in many related tasks; 3) How to make better use of deep networks for pedestrian
41 attribute recognition and what is the future direction of the development on at-
42 tribute recognition? By analysing existing person attribute recognition algorithms
43 and some top-ranked baseline methods, we draw some useful conclusions and
44 provide some possible research directions.

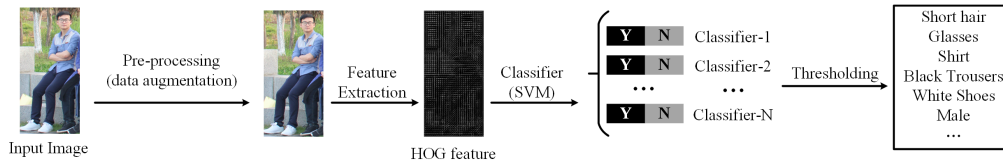


Figure 1: The regular pipeline of pedestrian attribute recognition.

45 2. Problem Formulation and Challenging Factors

46 Given a person image \mathcal{I} , pedestrian attribute recognition aims at predicting
 47 a group of attributes a_i to describe the characteristic of this person from a pre-
 48 defined attribute list $\mathcal{A} = \{a_1, a_2, \dots, a_L\}$. This task can be handled in different
 49 ways, such as multi-label classification and binary classification. As shown in
 50 Figure 1, the input images are usually processed with data augmentation to attain
 51 more training samples. Then, the features of processed images are extracted with
 52 deep learning methods or manual designed algorithms like HOG. With the feature
 53 representation and its labels, we can train the machine learning model, such as a
 54 classifier, for each attribute in a supervised way. In the testing phase, we can use
 55 this model to predict the response score of each attribute and assume that this input
 56 image has a corresponding attribute if its score is larger than the given threshold.
 57 In addition to such simultaneous attribute prediction, there are also algorithms that
 58 predict the attribute in a recurrent way, i.e., the attributes are predicted one after
 59 another.

60 Although good performance has been achieved based on deep learning mod-
 61 els, however, this task is still challenging due to the large intra-class variations in
 62 attribute categories (appearance diversity and appearance ambiguity (Deng et al.,
 63 2014)). We list challenging factors which may obviously influence the final recog-
 64 nize performance as follows: **1). Multi-views.** The images taken from different
 65 angles by the camera lead to the viewpoint issues for many computer vision tasks.
 66 Due to the body of human is not rigid, which further making the person attribute
 67 recognition more complicated. **2). Occlusion.** Partial occlusion of human body
 68 by other person or things increases the difficulty of person attributes recognition.
 69 Because the pixel values introduced by the occluded parts may make the model
 70 confused and lead to wrong predictions. **3). Unbalanced attribute distribution.**
 71 Each person have different attributes, therefore, the number of attributes are vari-
 72 able which leads to unbalanced data distribution. **4). Low resolution.** In practical
 73 scenarios, the resolution of images are rather low due to the high-quality cameras
 74 are rather expensive. **5). Illumination.** The images may taken from any time

75 in 24 hours. Hence, the light condition is variable at different time. The shadow
76 may also be taken in the person images and the images taken from night time
77 maybe totally ineffective. **6). Blur.** When person is moving, the images taken
78 by the camera may blur. Recognizing attributes in this situation is also a very
79 challenging task.

80 **3. The Review of PAR Algorithms**

81 In this section, we will review existing pedestrian attribute recognition algo-
82 rithms from following eight aspects: global based, local parts based, visual at-
83 tention based, sequential prediction based, newly designed loss function based,
84 curriculum learning based, graphic model based and others algorithms. A brief
85 summary of these methods can be found in Table 2 and 3.

86 *3.1. Global Image-based Models*

87 (Sudowe et al., 2015) proposes multi-branch classification layers for each at-
88 tribute learning with convolutional network. They adopt a pre-trained AlexNet
89 as basic feature extraction sub-network, and replace the last fully connected layer
90 with one loss per attribute using the KL-loss (Kullback-Leibler divergence based
91 loss function). (Li et al., 2015) introduce deep neural network for PAR and at-
92 tempt to handle the following two issues existed in traditional methods: 1). hand
93 crafted features; 2). ignored correlations between attributes. Two algorithms
94 DeepSAR and DeepMAR are proposed in this paper. DeepSAR do not model the
95 correlations between human attributes which maybe the key to further improv-
96 ing the overall recognition performance. Therefore, they propose the DeepMAR
97 which takes human image and its attribute label vectors simultaneously and jointly
98 considers all the attributes via sigmoid cross entropy loss. In addition, they also
99 consider the unbalanced label distribution in practical surveillance scenarios and
100 propose an improved loss function which widely used in many subsequent deep
101 PAR works. (Abdulnabi et al., 2015) propose a joint multi-task learning algorithm
102 for attribute estimation using CNN, named MTCNN. The MTCNN lets the CNN
103 models share visual knowledge among different attribute categories. They adopt
104 multi-task learning on the CNN features to estimate corresponding attributes and
105 use decomposition method to obtain shareable latent task matrix and combination
106 matrix from total classifier weights matrix. Thus, they can achieve flexible global
107 sharing and competition between groups through learning localized features. The
108 Accelerate Proximal Gradient Descent algorithm is used for the optimization.

109 Many works adopt CNN-RNN framework to take advantage of the intra-group
110 mutual exclusion and inter-group correlation, but they ignore the prior knowledge
111 underlying the attribute dataset. (Kai Han, 2019) propose to explore the corre-
112 lation between different attributes by mining the attribute co-occurrence prior.
113 Specifically, they integrate the information from different predictions with an at-
114 tribute aware pooling method. Their model follows multi-branch architecture and
115 context information is gathered to improve the final recognition performance.

116 **Summary:** According to aforementioned algorithms, we can find that these
117 algorithms all take the whole images as input and conduct multi-task learning for
118 PAR. They all attempt to learn more robust feature representations using feature
119 sharing, end-to-end training or multi-task learning. The benefits of these models
120 are simple, intuitive and highly efficient which are very important for practical
121 applications. However, the performance of these models is still limited due to the
122 lack of consideration of fine-grained recognition.

123 3.2. Part-based Models

124 As is known to all, we can train attribute classifiers simpler if we could iso-
125 late image patches corresponding to the same body part from the same viewpoint.
126 However, direct use object detectors is not reliable for body parts localization
127 before the year of 2011 due to its limited ability. (Bourdev et al., 2011) adopt
128 the *poselets* to decompose the image into a set of parts, each capturing a salient
129 pattern corresponding to a given viewpoint and local pose. This provides a ro-
130 bust distributed representation of a person from which attributes can be inferred
131 without explicitly localizing different body parts. Specifically, they first detect the
132 poselets on given image and obtain their feature representations. Then, they train
133 multiple SVM classifiers which are used for *poselet-level*, *person-level*, *context-*
134 *level* attribute classification, respectively.

135 **RAD*** (ICCV-2013, (Joo et al., 2013)) proposes a part learning algorithm
136 from the perspective of appearance variance while previous works focus on han-
137 dling geometric variation which require manual part annotation, such as poselet
138 (Bourdev et al., 2011). They first divide the image lattice into a number of over-
139 lapping sub-regions (named *window*). A grid of size $W \times H$ is defined and any
140 rectangle on the grid containing one or more number of cells of the grid forms a
141 window. The proposed method is more flexible in shape, size and location of part
142 window while previous works (such as spatial pyramid matching structure, SPM
143 (Lazebnik et al., 2006)) recursively divide the region into four quadrants and make
144 all subregions are squares that do not overlap with each other at the same level.

145 With all these windows, they learn a set of part detectors that are spatially as-
146 sociated with that particular window. For each window, all corresponding image
147 patches are cropped from training images and represented by HOG and color his-
148 togram feature descriptors. Then, K-means clustering is conducted based on the
149 extracted features. Each obtained cluster denotes a specific appearance type of a
150 part. They also train a local part detector for each cluster by logistic regression
151 as a initial detector and iteratively refine it by applying it in the entire set again
152 and updating the best location and scale to handle the issue of noisy clusters. Af-
153 ter learning the parts at multi-scale overlapping windows, they follow the method
154 for attribute classification proposed in the Poselet-based approach (Bourdev et al.,
155 2011). Specifically, they aggregate the scores from these local classifiers with the
156 weights given by part detection scores for final prediction.

157 **PANDA (CVPR-2014, (Zhang et al., 2014))** find the signal associated with
158 some attributes is subtle and the image is dominated by the effects of pose and
159 viewpoint. For the attribute of *wear glasses*, the signal is weak at the scale of
160 the full person and the appearance varies significantly with the head pose, frame
161 design and occlusion by the hair. They think the key to accurately predicting
162 the underlying attributes lies on locating object parts and establishing their cor-
163 respondences with model parts. They propose to jointly use global image and
164 local patches for person attributes recognition. They first detect the poselets, then
165 adopt CNN to extract the feature representations of the local patches and whole
166 human image. They directly feed the combined local and global features into the
167 linear classifier which is a SVM (Support Vector Machine) for multiple attributes
168 estimation.

169 **AAWP (ICCV-2015, (Gkioxari et al., 2015))** is introduced to validate whether
170 parts could bring improvements on both action and attribute recognition. The
171 CNN features are computed on a set of bounding boxes which associated with the
172 instance to classify, i.e., the whole instance, the oracle or person detector provided
173 and poselet-like part detector provided. For the part detector module, they design
174 their network by following the object detection algorithm RCNN (Girshick et al.,
175 2015). Given the image and detected parts, they use CNN to obtain fc7 features
176 and concatenate them into one feature vector as its final representation. Therefore,
177 the action or attribute category can be estimated with pre-trained linear SVM clas-
178 sifier. This work further expanding and validating the effectiveness and necessity
179 of parts in a more wider way.

180 **MLCNN (ICB-2015, (Zhu et al., 2015))** propose a multi-label convolutional
181 neural network to predict multiple attributes together in a unified framework. They
182 divide the whole image into 15 overlapping patches and use a convolutional net-

183 work to extract its deep features. They adopt corresponding local parts for spe-
184 cific attribute classification. They also use the predicted attributes to assist person
185 re-identification and their experiments validate the important role of attributes in
186 human related tasks.

187 **ALM (ICCV-2019, (Tang et al., 2019))** predict attributes in a hierarchical
188 manner and fuse these results with a simple voting scheme. More importantly,
189 they propose a weakly-supervised attribute localization module (ALM) based on
190 spatial transformer network for each branch. The ALM also contains a tiny channel-
191 attention module for feature augmentation. Their PAR network is trained with
192 deep supervision mechanism.

193 **ARAP (BMVC2016, (Luwei Yang and Tan, 2016))** adopts an end-to-end
194 learning framework for joint part localization and multi-label classification for
195 person attribute recognition. It mainly contains the initial convolutional feature
196 extraction layers, a key point localization network, an adaptive bounding box gener-
197 ator for each part, and the final attribute classification network for each part.
198 Their network contains three loss functions, i.e., the regression loss, aspect ratio
199 loss and classification loss. Specifically, they first extract the feature map of input
200 image, then conduct key points localization. Given the key points, they divide
201 human body into three main regions (including head, torso and legs) and obtain an
202 initial part bounding box. On the other hand, they also take previous fc7 layer’s
203 features as input and estimate the bounding box adjustment parameters. Given
204 these bounding box, they adopt bilinear sampler to extract corresponding local
205 features. Then, the features are fed into two fc layers for multi-label classifica-
206 tion.

207 **DeepCAMP (CVPR-2016, (Diba et al., 2016))** propose a novel CNN that
208 mines mid-level image patches for fine-grained human attributes recognition. Specif-
209 ically, they train a CNN to learn discriminative patch groups, named *DeepPattern*,
210 then, utilize regular contextual information and also deploy an iteration of feature
211 learning and patch clustering to purify the set of dedicated patches. The main
212 insight of this paper lies on that a better embedding can help improve the quality
213 of clustering algorithm in pattern mining algorithm. Therefore, they propose an
214 iteration algorithm where in each iteration, they train a new CNN to classify clus-
215 ter labels obtained in previous iteration to help improve the embedding. On the
216 other hand, they also concatenate features from both local patch and global human
217 bounding box to improve the clusters of mid-level elements.

218 **PGDM (ICME-2018, (Li et al., 2018))** is the first work which attempts to ex-
219 plore the structure knowledge of pedestrian body (i.e., pedestrian pose) for person
220 attributes learning. They first estimate the key points of given human image using

221 pre-trained pose estimation model. Then, they extract the part regions according
222 to these key points. The deep features of part regions and whole image are all ex-
223 tracted and used for attribute recognition independently. These two scores are then
224 fused together to achieve final attribute recognition. The attribute recognition al-
225 gorithm contains two main modules: i.e., the main net (AlexNet) and PGDM. The
226 introduced PGDM module is an existing pose estimation algorithm. They directly
227 train a regression network to predict the pedestrian pose with coarse ground truth
228 pose information which obtained from existing pose estimation model. Then,
229 they transform the key points into informative regions using spatial transformer
230 network, and use independent neural network for feature learning from each key
231 point related region. They jointly optimize the main net, PGDM and pose regres-
232 sion network.

233 **DHC (ECCV-2016, (Li et al., 2016))** propose to use *deep hierarchical con-*
234 *texts* to help person attribute recognition due to the background would sometimes
235 provide more information than target object only. Specifically, the *human-centric*
236 *context* and *scene context* are introduced in their network architecture. They first
237 construct input image pyramid and pass them all through VGG-16 to obtain multi-
238 scale feature maps. They extract features of four set of bounding box regions, i.e.,
239 the whole person, detected parts of target object, nearest neighbour parts from the
240 image pyramid and global image scene. The first two branches (the whole person
241 and parts) are regular pipeline for person attributes recognition algorithm. The
242 main contributions of this paper lie on the later two branches, i.e., the human-
243 centric and scene-level contexts help improve the recognition results. Once the
244 scores of these four branches are obtained, they sum up all the scores as final at-
245 tribute score. Due to the use of context information, this neural network needs
246 more external training data than regular pedestrian attribute recognition task. For
247 example, they need to detect the part of human body (head, upper and bottom
248 body regions) and recognize the style/scene of given image. They propose a new
249 dataset named *WIDER*, to better validate their ideas. Although the human attribute
250 recognition results can be improved significantly via this pipeline, however, this
251 model looks a little more complicated than other algorithms.

252 **LGNet (BMVC-2018, (Liu et al., 2018))** propose a Localization Guide Net-
253 work (LGNet) which can localize the areas corresponding to different attributes.
254 It also follows the local-global framework. Specifically, they adopt Inception-v2
255 as their basic CNN model for feature extraction. For global branch, they adopt
256 global average pooling layer (GAP) to obtain its global features. Then, a fc layer
257 is utilized to output its attribute predictions. For the local branch, they use 1×1
258 convolution layer to produce c class activation maps for each image. Then, they

259 capture an activation box for each attribute by cropping the high-response areas
260 of the corresponding activation map. They also use EdgeBoxes to generate region
261 proposals to obtain local features from the input image. In addition, they also
262 consider the different contributions of extracted proposals and different attributes
263 should focus on different local features. Therefore, they use the class active map
264 for each attribute to serve as a guide to determine the importance of the local
265 features to different attributes. Finally, the global and attended local features are
266 fused together by element-wise sum for PAR.

267 **Summary:** Based on the reviewed papers in this subsection, it is intuitive
268 to find that these algorithms all adopt both global and fine-grained local features.
269 The localization of body parts is achieved via an external part localization module,
270 such as part detection, pose estimation, poselets or proposal generation algorithm.
271 The use of part information improves the overall recognition performance signifi-
272 cantly. At the same time, it also brings some shortcomings as follows: Firstly, as
273 an operation in the middle phase, the final recognition performance heavily relies
274 on the accuracy of part localization. In another word, the inaccurate part detec-
275 tion results will bring the wrong features for final classification. Secondly, it also
276 needs more training or inference time due to the introducing of human body parts.
277 Thirdly, some algorithms need manual annotated labels for part location which
278 further increasing the cost of manpower and money.

279 3.3. Attention-based Models

280 **HydraPlus-Net (ICCV-2017, (Liu et al., 2017))** is introduced to encode
281 multi-scale features from multiple levels for pedestrian analysis using multi-directional
282 attention (MDA) modules. It contains two main modules, i.e., the Main Net (M-
283 net) which is a regular CNN and the Attentive Feature Net (AF-net) which in-
284 cludes multiple branches of multi-directional attention modules applied to differ-
285 ent semantic feature levels. The AF-net and M-net share same basic convolu-
286 tion architectures and their outputs are concatenated and fused by global average
287 pooling and fc layers. The output layer can be the attribute logits for attribute
288 recognition or feature vectors for person re-identification. In another word, it can
289 be used to minimize the cross-entropy loss and softmax loss for PAR and person
290 re-identification respectively.

291 **VeSPA (arXiv-2017, (Sarfranz et al., 2017))** takes the view cues into consid-
292 eration to better estimate corresponding attribute. Because the authors find that the
293 visual cues hinting at attributes can be strongly localized. Besides, the inference
294 of person attributes such as hair, backpack, shorts, are highly dependent on the ac-
295 quired view of the pedestrian. The image is fed into the Inceptions networks and

296 its feature representation can be obtained. The view-specific unit is introduced to
297 mapping the feature maps into coarse attribute prediction. Then, a view predictor
298 is used to estimate the view weights. The attention weights are used to multi-
299 ply view-specific predictions and obtain the final multi-class attribute prediction.
300 The view classifier and attribute predictors are trained with separate loss function.
301 The whole network is an unified framework and can be trained in an end-to-end
302 manner.

303 **DIAA (ECCV-2018, (Sarafianos et al., 2018))** can be seen as an ensemble
304 method for person attribute recognition. Their model contains a multi-scale visual
305 attention and a weighted focal loss for deep imbalanced classification. For the
306 multi-scale visual attention, the authors adopt feature maps from different layers.
307 They propose the weighted focal loss function to measure the difference between
308 predicted attribute vectors and ground truth. In addition, they also propose to
309 learn the attention maps in a weakly supervised manner (only the attribute labels,
310 no specific bounding box annotation) to improve the classification performance
311 by guiding the network to focus its resources to those spatial parts that contain
312 information relevant to the input image. The attention sub-network takes the fea-
313 ture map as input and output an attention mask. The output is then fed to attention
314 classifier to estimate the pedestrian attributes.

315 **CAM (PRL-2017, (Guo et al., 2017))** propose to use and refine attention
316 map to improve the performance of PAR. Their model contains two main mod-
317 ules, i.e., the multi-label classification sub-network and attention map refinement
318 module. The adopted CAM net also follows the category-specific framework, in
319 another word, different attribute classifiers have different parameters for the fc
320 layer. They use the parameters in fc layer as weights to linearly combine the fea-
321 ture maps from the last convolutional layer to get the attention of each category.
322 However, this naive implementation of attention mechanism could not focus on
323 the right regions all the time due to low resolution, over-fitting training. To handle
324 this issue, they exploring refine the attention map by tuning CAM network. They
325 measure the appropriateness of an attention map based on its concentration and
326 attempt to make the attention map to highlight a smaller but concentrated region.
327 Specifically, they introduce a weighted average layer to obtain attention map first.
328 Then, they use average pooling to down-sample its resolution to capture the im-
329 portance of all the potential relevant regions. After that, they also adopt softmax
330 layer to transform the attention map into a probability map. Finally, the maxi-
331 mum probability can be obtained via the global average pooling layer. On the
332 basis of the maximum probability, the authors propose the *exponential loss func-*
333 *tion* to measure the appropriateness of the attention heat map. For the training of

334 the network, the authors first pre-training the CAM network only by minimizing
335 classification loss; then, they adopt joint loss functions to fine-tuning the whole
336 network.

337 **JLPLS-PAA (TIP-2019, (Tan et al., 2019))** explore multiple attention mech-
338 anisms to select important and discriminative regions or pixels to handle the issues
339 such as large pose variations, clutter background. Different from regular spatial,
340 temporal or channel-view, they propose the parsing attention, label attention and
341 spatial attention. Specifically, the parsing model is used to locate the specific body
342 regions at pixel-level in a split-and-aggregate way. The label attention is formu-
343 lated by assigning several attention maps for each label under image-level super-
344 visions. The spatial attention is also considered to locate the most discriminative
345 image regions for all attributes with image-level supervisions. It is worthy to note
346 that this work is the first attempt to jointly learn multiple attention mechanisms in
347 a multi-task-like learning manner.

348 **IA²-Net (PRL-2019, (Ji et al., 2019))** propose an image-attribute reciprocal
349 guidance representation (RGR) method to investigate image-guided feature and
350 attribute-guided feature. Their method is developed based on the following obser-
351 vation: some attributes are concrete, such as “Hair Style, Shoes Style”, but some
352 are abstract attributes (For example, “Age Range, Role Types”). They also de-
353 velop a fusion attention mechanism to assign different attentions to different RGR
354 features. Besides, they combine the focal loss and cross-entropy loss to handle the
355 attribute imbalance problem.

356 **Da-HAR (AAAI-2020, (Wu et al., 2019))** attempt to recognize the human
357 attributes based on coarse-to-fine framework with self-mask operator. Their self-
358 mask block is trained on MS-COCO dataset and used for person segmentation.
359 With the help of a mask, their model is insensitive to distraction and clutter back-
360 ground. Hierarchical features from various layers of backbone network are fused
361 with 1×1 operator and attention module. The predictions from such side branch
362 are fused with the main branch for final decision making.

363 **CAS (ICME-2020, (Zeng et al., 2020))** A Co-Attentive Sharing module
364 is proposed by (Zeng et al., 2020) based on soft-sharing structure in multi-task
365 learning, which could mine discriminative channels and spatial regions for more
366 effective feature sharing. More detail, synergistic branch, attentive branch and
367 task-specific branch are explored for each layer, then, the results of three branches
368 are aggregated as the input features for the subsequent layer of each task.

369 (Zhang et al., 2019) propose the task-aware attention mechanism (named TAN)
370 to explore the importance of each position across different tasks. They first use
371 a cloth detector to crop out the target region, then, extract its feature with CNN.

372 The spatial attention and task attention modules are employed to learn feature
373 maps and the t-distribution Stochastic Triplet Embedding (t-STE) loss function is
374 used for the optimization.

375 **Summary:** Visual attention is a hot research topic in current deep learning
376 era and has been widely used in many domains. Generally speaking, attention
377 is the behavioral and cognitive process of selectively concentrating on a discrete
378 aspect of information, whether deemed subjective or objective, while ignoring
379 other perceivable information ¹. Pedestrian attribute recognition also follows this
380 framework and aforementioned works also validate the effectiveness of attention
381 mechanism. However, the works integrate with attention mechanism are still lim-
382 ited. How to design new attention models or directly borrow existing attention
383 algorithms from other domains is still unexplored.

384 3.4. Sequential Prediction based Models

385 **CNN-RNN (CVPR-2016, (Wang et al., 2016))** Regular multi-label image
386 classification framework learn independent classifier for each category and em-
387 ploy ranking or threshold on the classification results, fail to explicitly exploit
388 the label dependencies in an image. This paper first adopts RNNs to address
389 this problem and combine with CNNs to learn a joint image-label embedding to
390 characterize the semantic label dependency as well as the image-label relevance.
391 This model can model the label co-occurrence dependencies in the joint embed-
392 ding space by sequentially linking the label embeddings. For the inference of
393 CNN-RNN model, they attempt to find the sequence of labels that maximize the
394 prior probability. The training of the CNN-RNN model can be achieved by cross-
395 entropy loss function and back-propagation through time (BPTT) algorithm.

396 **JRL (ICCV-2017, (Wang et al., 2017))** firstly analyse existing learning is-
397 sues in the pedestrian attribute recognition task, e.g., poor image quality, appear-
398 ance variation and little annotated data, and propose to explore the interdepend-
399 ency and correlation among attributes and visual context as extra information
400 source to assist attribute recognition. Hence, the JRL model is proposed to joint
401 recurrent learning of attribute context and correlation, as its name shows. To better
402 mine these extra information for accurate person attribute recognition, the authors
403 adopt *sequence-to-sequence* model to handle aforementioned issues. They first
404 divide the given person image into multiple horizontal strip regions and form a
405 region sequences in top-bottom order. The obtained region sequences can be seen

¹<https://en.wikipedia.org/wiki/Attention>

406 as the input sentence in natural language processing, and can be encoded with the
407 LSTM network in a sequential manner. In decoding phase, the decoder LSTM
408 takes both *intra-person attribute context* and *inter-person similarity context* as in-
409 put and output variable-length attributes over time steps. The attribute prediction
410 in this paper can also be seen as a generation scheme. To better focus on local
411 regions of person image for specific attributes and obtain more accurate repre-
412 sentation, they also introduce the attention mechanism to attend the intra-person
413 attribute context.

414 **GRL (IJCAI-2018, (Zhao et al., 2018))** is developed based on JRL which
415 also adopts the RNN model to predict the human attributes in a sequential man-
416 ner. Different from JRL, GRL is formulated to recognize human attributes by
417 group, and gradually pay attention to both intra-group and inter-group relation-
418 ships. They divide the whole attribute list into many groups because the attributes
419 in intra-group are mutual exclusive and also correlated between inter-group. For
420 example, *BoldHair* and *BlackHair* cannot occur on the same person image, but
421 they are both related to the head-shoulder region of a person and can be in the
422 same group to be recognized together. It is an end-to-end single model algorithm
423 with no need for preprocessing and it also exploits more latent intra-group and
424 inter-group dependency among grouped pedestrian attributes.

425 **JCM (arXiv-2018, (Liu et al., 2018))** Existing sequential prediction based
426 PAR algorithms, such as JRL, GRL, may be easily influenced by different man-
427 ual division and attributes orders due to the weak alignment ability of RNN. This
428 paper proposes a joint CTC-Attention model (JCM) to conduct attribute recogni-
429 tion, which could predicts multiple attribute values with arbitrary length at a time
430 avoiding the influence of attribute order in the mapping table.

431 JCM is actually a multi-task network which contains two tasks: the attribute
432 recognition and person re-identification. They use ResNet-50 as the basic model
433 to extract features for both tasks. For the attribute recognition, they adopt the
434 Transformer as their attention model for the alignment of long attribute sequence.
435 And the connectionist temporal classification (CTC) loss and cross entropy loss
436 functions are used for the training of network. For the person re-ID, they directly
437 use two fully connected layers to obtain feature vectors and use softmax loss func-
438 tion to optimize this branch. In the test phase, the JCM could simultaneously pre-
439 dict the person identity and a set of attributes. They also use beam search for the
440 decoding of attribute sequence. Meanwhile, they extracts the features from the
441 CNN in base model to classify pedestrians for person re-ID task.

442 **RCRA (AAAI-2019, (Xin Zhao and Yan, 2019))** propose two models, i.e.,
443 Recurrent Convolutional (RC) and Recurrent Attention (RA) for pedestrian at-

444 tribute recognition. The RC model is used to explore the correlations between
445 different attribute groups with Convolutional-LSTM model and the RA model
446 takes the advantage of the intra-group spatial locality and inter-group attention
447 correlation to improve the final performance. Specifically, they first divide all the
448 attributes into multiple attribute groups, similar with GRL. For each pedestrian
449 image, they use CNN to extract its feature map and feed it to ConvLSTM layer
450 group by group. Then, new feature map for each time step can be obtained by
451 adding a convolutional network after ConvLSTM. Finally, the features are used
452 for attribute classification on current attribute group. Based on aforementioned
453 RC model, they also introduce visual attention module to highlight the region of
454 interest on the feature map. The attended feature maps are used for final classifi-
455 cation. The training of this network is also based on weighted cross-entropy loss
456 function proposed in WPAL-network.

457 **Summary:** As we can see from this subsection, these algorithms all adopt
458 the sequential estimation procedure. Because the attributes are correlated to each
459 other, and they also have various difficulties. Therefore, it is an interesting and in-
460 tuitive idea to adopt the RNN model to estimate the attributes one by one. Among
461 these algorithms, they integrate different neural networks, attribute groups, multi-
462 task learning into this framework. Compared with CNN based methods, these al-
463 gorithms are more elegant and effective. The disadvantage of these algorithms is
464 the time efficiency due to the successive attribute estimation. In the future works,
465 more efficient algorithms for the sequential attributes estimation are needed.

466 3.5. Newly Designed Loss Function based Models

467 **WPAL-network (BMVC-2017, (Zhou et al., 2017))** is proposed to simulta-
468 neously recognize and locate the person attributes in a weakly-supervised man-
469 ner (i.e., only person attribute labels, no specific bounding box annotation). The
470 GoogLeNet is adopted as their basic network for feature extraction. They fuse fea-
471 tures from different layers and feed them into Flexible Spatial Pyramid Pooling
472 layer (FSPP). The outputs of each FSPP are fed into fully connected layers and
473 output a vector whose dimension is same as the number of pedestrian attributes. In
474 addition, the authors also introduce a novel weighted cross entropy loss function
475 to handle the extremely imbalanced distribution of positive and negative samples
476 of most attribute categories.

477 **AWMT (MM-2017, (He et al., 2017))** As is known to all, the learning dif-
478 ficulty of various attributes is different. However, most of existing algorithms
479 ignore this situation and share relevant information in their multi-task learning
480 framework. This will leads to *negative transfer*, in another word, the inadequate

481 brute-force transfer may hurt the learner’s performance when two tasks are dis-
482 similar. AWMT proposes to investigate a shared mechanism that is possible of
483 *dynamically* and *adaptively* coordinating the relationships of learning different
484 person attribute tasks. Specifically, they propose an adaptively weighted multi-
485 task deep framework to jointly learn multiple person attributes, and a validation
486 loss trend algorithm to automatically update the weights of weighted loss layer.

487 They use ResNet-50 as backbone network and take both train and val images
488 as input. The basic network will output its predicted attribute vectors for both
489 train and val images. Hence, the train loss and val loss can be obtained simulta-
490 neously. The val loss is used to update the weight vectors which are then utilized
491 to weight different attributes learning. They propose the validate loss trend algo-
492 rithm to adaptively tuning the weight vector. The intuition behind their algorithm
493 is, when learning multiple tasks simultaneously, the “important” tasks should be
494 given higher weights to increase the scale of loss of the corresponding tasks.

495 **ArXiv-2019, (Yaghoubi et al., 2020)** is the first work which utilize the *hard*
496 attention to address the influence of background using binary mask predicted by
497 mask R-CNN. Then, they train their network based on the multi-task learning to
498 capture the semantic dependencies between most of the labels. The authors define
499 a weighted sum loss function to consider various contributions of each category
500 in the loss value.

501 **HFE (CVPR-2020, (Yang et al., 2020))** introduces external person ID con-
502 straints for hierarchical feature embedding (HFE) based on newly designed HFE
503 loss. This loss function is extended from triplet loss function and consists of
504 inter-triplet loss, intra-triplet loss and absolute boundary regularization. There-
505 fore, each class could gather more compactly, leading to a more distinct boundary
506 between classes.

507 Meanwhile, (Jia et al., 2020) argue that existing setting of PAR is not practical
508 because of the large number of identical pedestrian identities in train and test set.
509 They re-divide the dataset to ensure that the images with the same person ID do not
510 occur in train and test set simultaneously, and implement a strong baseline method
511 based on this setting. Their experimental results demonstrate that existing PAR
512 algorithms are overclaimed. They think distinguish the fine-grained attributes in
513 the same area (such as sandals *vs.* sneakers) is more important than locating the
514 area of the specific attribute.

515 (Ji et al., 2020) propose the **MTA-Net** to address complex relations between
516 images and attributes, and imbalanced distribution of pedestrian attributes. They
517 jointly use the knowledge of previous, current and next time steps based on CNN-
518 RNN framework. Besides, the focal balance loss (FBL) function is proposed to

519 handle the second issue.

520 **Summary:** There are few works focus on designing new loss functions for
521 pedestrian attribute recognition. WPAL-network (Zhou et al., 2017) consider the
522 unbalanced distribution of data and propose a weighted cross-entropy loss func-
523 tion according to the proportion of positive labels over all attribute categories in
524 the training dataset. This method seems a little tricky but has been widely used in
525 many PAR algorithms. AWMT (He et al., 2017) propose an adaptive weighting
526 mechanism for each attribute learning to make the network focus more on han-
527 dling the “hard” tasks. These works full demonstrate the necessity of designing
528 novel loss functions to better train the PAR network.

529 3.6. Curriculum Learning based Algorithms

530 **MTCT (WACV-2017, (Dong et al., 2017))** proposes a multi-task curriculum
531 transfer network to handle the issue on the lack of manually labelled training data.
532 Their algorithm contains multi-task network and curriculum transfer learning. For
533 the multi-task network, they adopt five stacked Network-In-Network (NIN) con-
534 volutional units and N parallel branches, with each branch representing a three
535 layers of fully connected sub-network for modelling one of the N attributes re-
536 spectively. Softmax loss function is adopted for the model training.

537 Cognitive studies suggest that a better learning strategy adopted by human/animals
538 is to start with learning easier tasks before gradually increasing the difficulties
539 of the tasks, rather than blindly learn randomly organised tasks. Therefore, they
540 adopt curriculum transfer learning strategy for clothing attribute modelling. Specif-
541 ically, it is consisted of two main stages. In the first stage, they use the clean
542 (easier) source images and their attribute labels to train the model. In the sec-
543 ond stage, they embed cross-domain image pair information and simultaneously
544 append harder target images into the model training process to capture harder
545 cross-domain knowledge. They adopt t-STE (t-distribution stochastic triplet em-
546 bedding) loss function to train the network

547 **CILICIA (ICCV-2017, (Sarafianos et al., 2017))** Similar with MTCT (Dong
548 et al., 2017), CILICIA also introduces the idea of curriculum learning into person
549 attribute recognition task to learn the attributes from easy to hard. They explore
550 the correlations between different attribute learning tasks and divide such correla-
551 tions into strongly and weakly correlated tasks. Specifically, under the framework
552 of multi-task learning, they use the respective Pearson correlation coefficients to
553 measure the strongly correlated tasks. For the multi-task network, they adopt the
554 categorical cross-entropy function (Zhu et al., 2017) to measure the difference be-
555 tween predictions and targets. To weight different attribute learning tasks, one

556 intuitive idea is to learn another branch network for weights learning. They adopt
557 the *supervision transfer* learning technique to help attribute learning in weakly
558 correlated group.

559 They also propose CILICIA-v2 (Sarafianos et al., 2018) by introducing an
560 effective method to obtain the groups of tasks using hierarchical agglomerative
561 clustering. It can be any number and not just only two groups (i.e., strong/weakly
562 correlated).

563 **DCL (ICCV-2019, (Wang et al., 2019))** introduces an unified framework,
564 named dynamic curriculum learning, to online adaptively adjust the sampling
565 strategy and loss learning in a batch to handle the issues caused by imbalanced
566 data distribution. Specifically, they design two level curriculum schedulers: sam-
567 pling scheduler and loss scheduler. The first one aims at finding the most mean-
568 ingful samples in one batch to learn from imbalanced to balanced distribution and
569 easy to hard. The second one is used to achieve a good trade-off between clas-
570 sification and metric learning loss. They achieve new state-of-the-art recognition
571 performance on two attribute datasets.

572 **Summary:** Inspired by recent progress of cognitive science, the researchers
573 also consider using such “easy” to “hard” learning mechanism for PAR. They
574 introduce existing curriculum learning algorithm into their learning procedure to
575 model the relations between each attribute. This makes the PAR algorithms look
576 more intelligent due to the ability of estimating the “easier” attributes first just
577 like humans. Some other algorithms such as self-paced learning are also used to
578 model the multi-label classification problem or other computer vision tasks. It is
579 also worthy to introduce more advanced works of cognitive science to guide the
580 learning of PAR. In addition, the meta-learning has shown its ability to “learning
581 to learn” in many tasks, such as fine-grained classification, few-shot learning. It
582 will also be an interesting research direction to integrate this learning framework
583 for PAR.

584 3.7. *Graphic Model based Algorithms*

585 Graphic models are commonly used to model structure learning in many ap-
586 plications. Similarly, there are also some works to integrate these models into the
587 PAR task.

588 **DCSA* (ECCV-2012, (Chen et al., 2012))** propose to model the correlations
589 between human attributes using conditional random field (CRF). They first esti-
590 mate the pose information and locate the local parts of upper body only. Then, four
591 types of base features are extracted from these regions. These features are fused to
592 train multiple attribute classifiers via SVM. The key idea of this paper is to apply

593 the fully connected CRF to explore the mutual dependencies between attributes.
594 They treat each attribute function as a node of CRF and the edge connecting every
595 two attribute nodes reflects the joint probability of these two attributes. The belief
596 propagation is adopted to optimize the attribute label cost.

597 **A-AOG*** (TPAMI-2018, (Park et al., 2018)) is short for attribute And-Or
598 grammar, which is proposed explicitly to represent the decomposition and articulation
599 of body parts, and account for the correlations between poses and attributes.
600 This algorithm is developed based on And-Or graph and the and-nodes
601 denote decomposition or dependency; the or-nodes represent alternative choices
602 of decomposition or types of parts. Specifically speaking, it mainly integrates the
603 three types of grammars: *phrase structure grammar*, *dependency grammar* and
604 an *attribute grammar*. They use deep CNN to generate proposals for each part
605 and adopt greedy algorithm based on the beam search to optimize aforementioned
606 objective function.

607 **VSGR** (AAAI-2019, (HUANG, 2019)) propose to estimate the pedestrian
608 attributes via visual-semantic graph reasoning (VSGR). They argue that the accuracy
609 of person attribute recognition is heavily influenced by: 1). only local parts
610 are related with some attributes; 2). challenging factors, such as pose variation,
611 viewpoint and occlusion; 3). the complex relations between attributes and different
612 part regions. Therefore, they propose to jointly model spatial and semantic
613 relations of region-region, attribute-attribute, and region-attribute with a graph-
614 based reasoning framework.

615 This algorithm mainly contains two sub-networks, i.e., the visual-to-semantic
616 sub-network and semantic-to-visual sub-network. For the first module, it first
617 divides the human image into a fixed number of local parts. They construct a
618 graph whose node is the local part and edge is the similarity of different parts.
619 Different from regular relation modelling, they adopt both the similarity relations
620 between parts and topological structures to connect one part with its neighbour
621 regions. The two sub-graphs are combined to compute the output of spatial graph.
622 The semantic-to-visual sub-network can also be processed in similar manner and it
623 also outputs sequential attribute prediction. The outputs of these two sub-networks
624 are fused as the final prediction and can be trained in an end-to-end way.

625 **JLAC** (AAAI-2020, (Tan et al., 2020)) propose the JLAC (Joint Learning
626 of Attribute and Contextual relations) for PAR which contains two main modules:
627 Attribute Relation Module (ARM) and Contextual Relation Module (CRM). The
628 ARM module is used to explore the correlations among multiple attributes based
629 on an attribute graph with attribute-specific features. For the CRM, the authors
630 construct a graph projection scheme that targets at project the 2-D feature map

631 into a set of nodes from different image regions. This module fully explored
632 the contextual relations among those regions. The GCN is adopted to mine the
633 graph structured features for the two modules and the whole architecture can be
634 optimized in an end-to-end manner.

635 **BCRNNs (CVPR-2018, (Wang et al., 2018))** propose to use Bidirectional
636 Convolutional Recurrent Neural Networks (BCRNNs) to address the problem of
637 visual fashion analysis based on their defined grammar topologies. Specifically,
638 their proposed dependency grammar could capture kinematics-like relations, and
639 symmetry grammar can accounting for the bilateral symmetry of clothes.

640 **Summary:** Due to the relations existed in multiple attributes, many algo-
641 rithms are proposed to discover such information for PAR. Therefore, the Graphic
642 models are easily introduced into the learning pipeline, such as Markov Random
643 Field, Conditional Random Field, And-Or-Graph or Graph Neural Networks. The
644 works reviewed in this subsection are the outputs by integrating the graphic mod-
645 els with PAR. Maybe the other graphic models can also be used for PAR to achieve
646 better recognition performance. Although these algorithms have so many advan-
647 tages on model the relations between pedestrian attributes, however, these algo-
648 rithms seem more complex than others. The efficiency issue is also needs to be
649 considered in practical scenarios.

650 3.8. Other Algorithms

651 This subsections are used to demonstrate algorithms that not suitable for afore-
652 mentioned categories, including: PatchIt (Sudowe and Leibe, 2016), FaFS (Lu
653 et al., 2017), GAM (Fabbri et al., 2017) and IFSL (Liuyu Xiang, 2019).

654 PatchIt proposes a self-supervised pre-training approach, named PatchTask,
655 to obtain weight initializations for the PAR. It’s key insight is to leverage data
656 from the same domain as the target task for pre-training and it only relies on
657 automatically generated rather than human annotated labels.

658 FaFS is proposed to design compact multi-task deep learning architecture au-
659 tomatically. This algorithm starts with a thin multi-layer network and dynamically
660 widens it in a greedy manner during training. This will create a tree-like deep ar-
661 chitecture by repeating above widening procedure and similar tasks reside in the
662 same branch until at the top layer.

663 GAM proposes to handle the issue of occlusion and low resolution of pedes-
664 trian attributes using deep generative models. Specifically, their overall algorithm
665 contains three sub-networks, i.e., the attribute classification network, the recon-
666 struction network and super-resolution network. For the attribute classification
667 network, they also adopt joint global and local parts for final attribute estimation.

668 To handle the occlusion and low-resolution problem, they introduce the deep gen-
669 erative adversarial network (Mirza and Osindero, 2014) to generate re-constructed
670 and super-resolution images. And use the pre-processed images as input to the
671 multi-label classification network for attribute recognition.

672 (Liuyu Xiang, 2019) propose the IFSL to handle the few-shot pedestrian at-
673 tribute recognition problem. Because most previous PAR algorithms are designed
674 for a fixed set of attributes and unable to handle the incremental few-shot learning
675 scenario. This work introduces an extra module named attribute prototype gen-
676 erator, which can be seen as a high-level meta-learner that extracts the multiple-
677 attribute information from the feature embedding. And it can produce discrimina-
678 tive attribute prototype embedding and therefore provide the classification weights
679 for the novel attributes.

680 (Zhang et al., 2020) propose the TS-FashionNet, i.e. the Texture and Shape
681 biased Two-Stream Networks, for fashion image analysis. Specifically, the shape-
682 biased stream contains a landmark branch to help extract shape features; while
683 the texture-biased stream is used to emphasize on the extraction of texture fea-
684 tures. Then, these two branches are concatenated together to predict the clothing
685 attributes and classify the clothes categories.

686 (Jia et al., 2021) argue that current evaluation for PAR is not consistent with
687 practical scenarios and advocate zero-shot pedestrian identity setting. They pro-
688 pose two new dataset *PETA_{ZS}* and *RAP_{ZS}* for the evaluation.

689 4. Discussion

690 In this section, we will first discuss the specific attribute recognition in this
691 section, then, we will give a comparison between deep learning and traditional
692 algorithm based PAR methods. After that, we will show the connections between
693 PAR and other computer vision tasks.

694 4.1. Specific Attribute Recognition

695 In addition to the attribute recognition on whole body, there are also some
696 attribute recognition algorithms focus on local parts of people, for example, face
697 attribute recognition (e.g., gender, age, race). In this subsection, we will give
698 a brief review on specific attribute recognition algorithms. For a more detailed
699 introduction for face attribute recognition, please refer to the (Zheng et al., 2018)
700 and (Fasel and Luetin, 2003).

701 (Rodríguez et al., 2017) is proposed to discover the most informative and reli-
702 able parts of a given face for improving age and gender classification. Specifically,

703 it is a feedforward attention mechanism and mainly consists of three modules: an
704 attention CNN, a patch CNN and a multi layer perceptron (MLP). The two CNN
705 modules are used to predict the best attention grid to perform the glimpses and
706 evaluate the higher resolution patches based on their importance predicted by the
707 attention grid, respectively. The MLP module is used to integrate the informa-
708 tion from both CNNs and make the final classifications. (Li et al., 2017) propose
709 cumulative hidden layer and comparative ranking layer to combat the sample im-
710 balance problem and learn more effective aging features. The cumulative hidden
711 layer is supervised by a point-wise cumulative signal which encodes the target
712 age labels continuously. The comparative ranking layer is supervised by a pair-
713 wise comparative signal, in another word, who is older. This is inspired by the
714 observation that it is easier to tell which one is older given two faces than tell
715 its accurately age. (Xing et al., 2017) conduct a comprehensively diagnose on
716 the training and evaluating procedures of deep leaning methods for age estima-
717 tion. They achieve state-of-the-art performance by following previous work with
718 appropriate problem formulation and loss function. They also consider various
719 factors to build a better age estimation model based on multi-task learning frame-
720 work, such as the strategies to incorporate information like race and gender. Their
721 studies are helpful to get better understandings of a deep age estimation algorithm.
722 (Antipov et al., 2017) shed light on some open questions of human demograph-
723 ics estimation to improve the existing CNN-based approaches for gender and age
724 prediction. Their work analyse four important factors of the CNN training: the
725 target age encoding and loss function, the CNN depth, the pre-training, the train-
726 ing strategy. Then, they deign their model based on these experiments and achieve
727 state-of-the-art performance. (Liu et al., 2017) propose a group-aware deep fea-
728 ture learning approach for facial age estimation. Specifically, they split ordinal
729 ages into a set of discrete groups and learn deep feature transformations across
730 age groups to project each face pair into the new feature space. They simultane-
731 ously minimize the intra-group variances of positive face pairs and maximize the
732 inter-group variances of negative face pairs. (Chen et al., 2017) propose an ap-
733 proach to automatically discover “spectral attributes” which avoids manual work
734 required for defining hand-crafted attribute representations. (Fasel and Luettn,
735 2003) conduct an review on automatic facial expression analysis including: facial
736 motion, deformation extraction approaches and classification methods. (Hadid
737 and Pietikäinen, 2009) investigate the combination of facial appearance and mo-
738 tion for face analysis in videos. They are inspired by the psychophysical finds
739 which state that facial movements can provide valuable information to face anal-
740 ysis. They design an extended set of volume local binary patterns as well as a

741 boosting scheme for spatio-temporal face and gender recognition from videos.

742 There are also some works focusing on backpack detection given a human
743 image, for example, (Branca et al., 2002), (Ghadiri et al., 2019), (Damen and
744 Hogg, 2011), (Ghadiri et al., 2016). The regular pipeline of these methods is
745 to detect the human body first, then segment the carried object in a fine-grained
746 manner.

747 4.2. Comparison between Deep Learning and Traditional based Algorithm

748 Before the deep neural network based algorithms take over the PAR com-
749 munity, most of traditional approaches follow a standard pipeline, which can be
750 found in Fig. 1. Usually, we need to first conduct some pre-processing to augment
751 the dataset, such as flip, rotation, scale variation, crop, translation, add Gaussian
752 noise. Then, manual designed features (for example, HOG or SIFT features) are
753 extracted to represent the person image. After that, multiple classifiers are trained
754 to discriminate all the pedestrian attributes, such as support vector machine. In the
755 test phase, we need to set a threshold to give an estimation whether corresponding
756 attribute exists or not.

757 According to aforementioned PAR algorithms including traditional methods
758 and deep learning based approaches, we can find the following observations: 1).
759 Both methods all attempt to handle the PAR from the fine-grained perspective,
760 such as estimate the attributes from local human body. The major difference lies
761 on how to locate these regions: traditional methods rely on object detector, while
762 deep learning methods employ more advanced object detector, visual attention
763 mechanisms or some other information obtained from auxiliary task (for example,
764 pose estimation). 2). Both methods all need the powerful feature representation
765 of pedestrian images. Traditional approaches use the manual designed features,
766 while deep learning based algorithms could learn the deep features automatically
767 from given training dataset. This is also one of the most unique characteristics
768 of deep learning based PAR algorithms. 3). Both methods all attempt to utilize
769 the prior information or relations between human attributes to augment the final
770 recognition performance. Traditional methods usually adopt graphical models
771 such as conditional random field, markov random field as post-processing, while
772 deep learning based algorithms can integrate such relations into their pipeline and
773 learning in an end-to-end manner based on graph neural networks.

774 Generally speaking, traditional and deep learning based PAR algorithms all
775 share similar ideas, but deep learning methods always achieve better recognition
776 accuracy than traditional algorithms. We think one of the most important and
777 intuitive reasons is the powerful deep features which can learn from large scale

778 datasets. Another reason is that many challenges of PAR are hard to be modelled
779 with traditional algorithms, but this is easy to be implemented with deep neural
780 networks. The third reason is that deep neural networks can be integrated with
781 traditional methods, i.e., the mode of “deep + X”. This will further extending the
782 applications of deep neural networks.

783 *4.3. Connections between PAR and Other Tasks*

784 Visual attributes can be seen as a kind of mid-level feature representation
785 which may provide important information for high-level human related tasks, such
786 as person re-identification, pedestrian detection, person tracking, person retrieval,
787 human action recognition and scene understanding.

788 For the pedestrian detection, regular algorithms treat it as a single binary clas-
789 sification task, while (Tian et al., 2015) propose to jointly optimize person detec-
790 tion with semantic tasks to address the confusion of positive and hard negative
791 samples. They use existing scene segmentation dataset to transfer attribute infor-
792 mation to learn high-level features from multiple tasks and dataset sources.

793 For the person re-identification, pedestrian attributes can be seen as a kind
794 of middle-level representation and share a common target at the pedestrian de-
795 scription with person re-ID. PAR focuses on local information mine while person
796 re-identification usually capture the global representations of a person. There are
797 already many works attempting to integrate the PAR into their person re-ID sys-
798 tem. For example, (Lin et al., 2019) propose an attribute-person recognition net-
799 work, a multi-task network which learns a re-ID embedding and predicts person
800 attributes simultaneously. (Han et al., 2018) propose an attribute-aware attention
801 model to learn local attribute and global category representation simultaneously
802 in an end-to-end fashion. (Su et al., 2016) also propose to integrate the mid-level
803 attributes into person re-identification framework and train the attribute model in
804 a semi-supervised manner. Specifically, they first pre-train the deep CNN on an
805 independent attribute dataset, then, fine-tuned on another dataset only annotated
806 with person IDs. After that, they estimate attribute labels for target dataset us-
807 ing the updated deep CNN model. (Khamis et al., 2014) propose to integrate a
808 semantic aspect into regular appearance-based methods. They jointly learn a dis-
809 criminative projection to a joint appearance-attribute subspace, which could ef-
810 fectively leverage the interaction between attributes and appearance for matching.
811 (Li et al., 2015) also present a comprehensive study on clothing attributes to assist
812 person re-ID. They first extract the body parts and their local features to allevi-
813 ate the pose-misalignment issues. Then, they propose a latent SVM based person

814 re-ID approach to model the relations between low-level part features, middle-
815 level clothing attributes and high-level re-ID labels of person pairs. They treat the
816 clothing attributes as real-value variables instead of using them as discrete vari-
817 ables to obtain better person re-ID performance. (Layne et al., 2012) and (Layne
818 et al., 2014) are all learn an attribute-center representation to describe people and
819 a metric to compare attribute profiles. (Layne et al., 2012) also achieve better
820 re-ID performance by learning a selection and weighting of mid-level semantic
821 attributes for the description of people. (Schumann and Stiefelhagen, 2017) first
822 train an attribute classifier and take its responses into the learning of person re-
823 ID model based on CNNs. (Li et al., 2019) find that attributes are related to
824 specific local regions and utilize the attribute detection to generate correspond-
825 ing attribute-part detectors. This will handle the body part misalignment problem
826 significantly for the re-ID task. (Ling et al., 2019) propose a multi-task learning
827 network with multiple classification and verification losses for person re-ID which
828 closely combine person identity and pedestrian attribute task. In (Su et al., 2017),
829 the authors use the idea of multi-shot re-identification for person re-ID instead of a
830 single prob image. Specifically, they utilize low-level features, attributes and inter-
831 attribute correlations to make their model robust under the multi-camera setting.
832 (Chen et al., 2018) also develop a CNN-based pedestrian attribute-assisted person
833 re-identification framework. They first learn the attribute with a part-specific
834 CNN and fuse them with low-level robust LOMO features. Then, they merge the
835 learned attribute CNN embedding with identification CNN embedding under a
836 triplet structure for person re-ID.

837 There are also some works integrating pedestrian attributes for person retrieval
838 and human active recognition. For the person retrieval, (Wang et al., 2013) lever-
839 age low-level features (e.g., color) and high-level features (i.e. the person at-
840 tributes) of clothing to tackle the issues caused by geometric deformation, oc-
841 clusion and clutter background. Their content-based image retrieval algorithm is
842 developed based on the bag-of-visual-words model. More importantly, they pro-
843 pose a re-ranking approach to improve the search result by exploiting attributes,
844 such as the type of clothing, sleeves and patterns. (Chen et al., 2015) approach the
845 problem of describing people by first mining clothing attributes with fine-grained
846 attribute labels from online shopping stores. Then, they use a double-path deep
847 domain adaptation network to bridge the gap between the collected images and
848 practical testing data. Their work validate the effectiveness and importance of
849 person attributes for people describe. For the human active recognition, there is a
850 literature review summarized by (Ziaeeafard and Bergevin, 2015) which also men-
851 tion that the attributes are an element of semantic space and are effective features

852 describing a basic or an intrinsic characteristic of an activity. In addition, (Liu
853 et al., 2011) validate that attributes enable the construction of more descriptive
854 models for human action recognition. They select attributes in a discriminative
855 fashion or coherently integrate with data-driven attributes to make the attribute set
856 more descriptive.

857 Due to the pedestrian attribute recognition is mainly focus on the clothing fea-
858 ture studied in many other research topics, such as part-detection, pose estimation
859 (Murphy-Chutorian and Trivedi, 2008) and human parsing (Huang et al., 2018).
860 But these tasks have their own emphasized point, for example: part-detection aims
861 at locating the local parts of object using a bounding box; pose estimation focuses
862 on locating the key points of people which will be useful for human activity recog-
863 nition; And human parsing is a more fine-grained pixel-wise segmentation of hu-
864 man body which is more difficult than pedestrian attribute recognition. However,
865 these tasks can be learned in a joint manner due to these tasks are all focus on
866 human body and also have their own emphasized point. Actually, the multi-task
867 learning has been studied for a long time in machine learning, pattern recognition
868 and computer vision community. The joint learning of pedestrian attribute recog-
869 nition and other tasks also validate the effectiveness of such multi-task setting,
870 such as joint PAR and person re-ID algorithms described above.

871 **5. Benchmarks**

872 *5.1. Datasets*

873 Unlike other tasks in computer vision, for pedestrian attribute recognition, the
874 annotation of dataset contains many labels at different levels. For example, hair
875 style, color, hat and glass, are seen as specific low-level attributes and correspond
876 to different areas of the images; while some attributes are abstract concepts, such
877 as gender, orientation and age, which do not correspond to certain regions, we
878 consider these attributes as high-level attributes. Furthermore, human attribute
879 recognition is generally severely affected by environmental or contextual factors,
880 such as viewpoints, occlusions and body parts. In order to facilitate the study,
881 some datasets provide annotations of perspective, parts bounding box, occlusion.

882 By reviewing related work in recent years, we have found and summarized
883 several datasets which are used to research pedestrian attribute recognition. As
884 shown in Table 1, we only show some important parameters of these benchmark
885 datasets, such as image numbers, attribute numbers, image source and correspond-
886 ing project pages due to the limited space in this paper. For more detailed infor-

887 mation of these datasets, please visit our project page for the arXiv version (Xiao
 888 et al., 2019).

Table 1: An overview of pedestrian attribute datasets (the # denotes the number of).

| Dataset | # Pedestrians | #Attributes (Binary/Multi-class) | Source |
|-----------------------------|---------------|----------------------------------|---|
| PETA | 19000 | 61/4 | outdoor & indoor |
| RAP | 41585 | 69/3 | indoor |
| RAP-2.0 | 84928 | 69/3 | indoor |
| PA-100K | 100000 | 26/0 | outdoor |
| WIDER | 13789 | 14/0 | WIDER images (Xiong et al., 2015) |
| Market-1501 | 32668 | 26/1 | outdoor |
| DukeMTMC | 34183 | 23/0 | outdoor |
| PARSE-27K | 27000 | 8/2 | outdoor |
| APIS | 3661 | 11/2 | KITTI (Geiger et al., 2012) , CBCL Street Scenes (Bileschi, 2006), INRIA (Dalal and Triggs, 2005) and SVS |
| HAT | 9344 | 27/0 | image site Flickr |
| CRP | 27454 | 1/13 | outdoor |
| CAD | 1856 | 23/3 | image site Sartorialist ² and Flickr |
| BAP | 8035 | 9/0 | H3D (Bourdev and Malik, 2009) dataset PASCAL VOC 2010 |
| UAV-Human (Li et al., 2021) | 22,263 | 7/0 | outdoor (UAV) |

889 5.2. Evaluation Criteria

890 The performance of attribute classification can be evaluated with the Receiver
 891 Operating Characteristic (ROC) and the Area Under the average ROC Curve (AUC)
 892 which are calculated by two indicators, the recall rate and false positive rate. The
 893 recall rate is the fraction of the correctly detected positives over the total amount of
 894 positive samples, and the false positive rate means the fraction of the misclassified
 895 negatives out of the whole negative samples. At various threshold settings, a ROC
 896 curve can be drawn by plotting the recall rate vs. the false positive rate. However,
 897 seldom of PAR algorithms adopt these two metrics except for (Zhu et al., 2013).
 898 The Geometric Mean (G-mean) is used by (Chen et al., 2012) for the evaluation,
 899 which is a popular evaluation metric for unbalanced data classification.

900 In addition to aforementioned metrics, the mean accuracy (mA) is also used
 901 to evaluate the attribute recognition algorithms. For each attribute, mA calcu-
 902 lates the classification accuracy of positive and negative samples respectively, and
 903 then gets their average values as the recognition result for the attribute. Finally, a
 904 recognition rate is obtained by taking an average over all attributes. The evaluation
 905 criterion can be calculated through the following formula:

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

906 where L is the number of attributes. TP_i and TN_i are the number of correctly

907 predicted positive and negative examples respectively, P_i and N_i are the number
 908 of positive and negative examples respectively.

909 Aforementioned evaluation criteria treat each attribute independently and ig-
 910 nore the inter-attribute correlation which exists naturally in multi-attribute recog-
 911 nition problem. (Li et al., 2016) named these metrics as *label-based* criteria
 912 and propose to use the *example-based* evaluation criteria inspired by a fact that
 913 example-based evaluation captures better the consistence of prediction on a given
 914 pedestrian image. Four widely used metrics, i.e., accuracy, precision, recall rate
 915 and F1 value, can be defined as:

$$Acc = \frac{1}{N} \sum_{i=1}^N \frac{|Y_i \cap f(x_i)|}{|Y_i \cup f(x_i)|}, \quad Prec = \frac{1}{N} \sum_{i=1}^N \frac{|Y_i \cap f(x_i)|}{|f(x_i)|}, \quad Rec = \frac{1}{N} \sum_{i=1}^N \frac{|Y_i \cap f(x_i)|}{|Y_i|}, \quad F1 = \frac{2 * Prec * Rec}{Prec + Rec} \quad (2)$$

916 where N is the number of examples, Y_i is the ground truth positive labels of the
 917 i -th example, $f(x)$ returns the predicted positive labels for i -th example. And $|\cdot|$
 918 means the set cardinality. Due to the ROC, AUC and G-mean are only used in
 919 a few PAR works, thus, we only report the main experimental results based on
 920 mAP, accuracy, precision, recall and F1 value in Table 2 and Table 3.

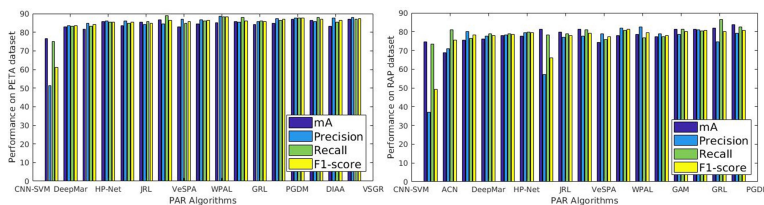


Figure 2: Comparison of selected 17 PAR algorithms (from 2014 to 2020) on the PETA and RAP dataset.

921 5.3. Performance Evaluation

922 In this section, we give a brief introduction to the performance of selected 17
 923 PAR algorithms proposed from 2014 to 2020. As shown in Fig. 2, we can find that
 924 the baseline method CNN-SVM is outperformed by recent deep learning based
 925 PAR approaches significantly on both large scale benchmark datasets RAP and
 926 PETA. Specifically, recent deep learning approaches improve the baseline from
 927 about 50+% to 80+% on multiple evaluation metrics. These experimental results
 928 fully demonstrate the effectiveness and advantages of deep learning based PAR
 929 algorithms. Interestingly, we also find that the accuracy of current deep learning
 930 based methods are comparable, and there is no significant improvement of current

931 methods (in 2020) compared with deep PAR algorithms proposed in several years
932 ago. More detailed experimental results of these methods can be found in Table 2
933 and Table 3. Therefore, how to design new modules for the further improvement
934 of PAR results in future works? In the following section, we propose several
935 possible research directions for PAR.

936 6. Future Research Directions

937 **More Accurate and Efficient Part Localization Algorithm** Human beings
938 could recognize the detailed attributes information in an very efficient way, be-
939 cause we can focus on specific regions in a glimpse and reason the attribute based
940 on the local and global information. Therefore, it is an intuitive idea to design
941 algorithms which can detect the local parts for accurate attribute recognition. Ac-
942 cording to section 3.2, it is easy to find that researchers are indeed more interested
943 in mining local parts of human body. They use manual annotated or detected
944 human body or pose information for the part localization. There are also some al-
945 gorithms attempting to propose unified framework in a weakly supervised manner
946 to jointly handle the attribute recognition and localization. We think this will also
947 be a good and useful research direction for pedestrian attribute recognition.

948 **Deep Generative Models for Data Augmentation** In recent years, the deep
949 generative models have made great progress and many algorithms are proposed.
950 One intuitive research direction is how can we use deep generative models to
951 handle the issues of low-quality person images or unbalanced data distribution?
952 There are already many researches who focus on image generation with the guid-
953 ance of text, attribute or pose information. The generated images can be used
954 in many other tasks for data augmentation, for example, object detection, person
955 re-identification and visual tracking (Wang et al., 2018). It is also worthy to de-
956 sign new algorithms to generate pedestrian images according to given attributes
957 to augment the training data.

Table 2: An overview of PAR algorithms reviewed in this paper (Part-I).

| Algorithm | Part | Attention | Seq. | C. L. | Graphic | Groups | Loss | Accuracy |
|---|------|-----------|------|-------|---------|--------|--------------------------------|--|
| Poselets (Bourdev et al., 2011) (ICCV-2011) | ✓ | | | | | | - | mAP BAP/Attributes25K: 65.18/51.06 |
| DCSA (Chen et al., 2012) (ECCV-2012) | ✓ | | | | ✓ | | SVM | - |
| RAD (Joo et al., 2013) (ICCV-2013) | ✓ | | | | | | | mAP HAT: 59.3 |
| PANDA (Zhang et al., 2014) (CVPR-2014) | ✓ | | | | | | SVM | mAP BAP/Attributes25K: 78.98/70.74 |
| ACN (Sudowe et al., 2015) (ICCVW-2015) | | | | | | | KL-loss | PARSE-27K: 63.6, mAP HATDB: 66.2, BAP: 80.02 |
| DeepSAR (Li et al., 2015) (ACPR-2015) | | | | | | | Softmax Loss | Accuracy PETA: 81.3 |
| DeepMAR (Li et al., 2015) (ACPR-2015) | | | | | | | Weighted Cross-entropy Loss | Accuracy PETA: 82.6 |
| MTCNN (Abdulnabi et al., 2015) (TMM-2015) | | | | | | ✓ | Softmax Loss | Accuracy AwA: 81.19 |
| MLCNN (Zhu et al., 2015) (ICB-2015) | ✓ | | | | | | Softmax Loss | VPeR: 74.1, Accuracy GRID: 73.2 |
| AAWP (Ghoxari et al., 2015) (ICCV-2015) | ✓ | | | | | | SVM | mAP BAP: 83.6 |
| ARAP (Luwei Yang and Tan, 2016) (BMVC-2016) | ✓ | | | | | | Softmax Loss | MPI-AlexNet: 78.00/73.2/77.74, Accuracy Garment-AlexNet: 76.24/67.70/77.48 |
| DeepCAMP (Diba et al., 2016) (CVPR-2016) | ✓ | | | | | | Softmax Loss | mAP BAP: 86.6 |
| DHC (Li et al., 2016) (ECCV-2016) | ✓ | | | | | | Cross-entropy Loss | BAP: 92.2, mAP HAT: 78.0 |
| PatchI (Sudowe and Leibe, 2016) (BMVC-2016) | | | | | | | Cross-entropy Loss | WIDER: 81.3 |
| HydraPlus-Net (Liu et al., 2017) (ICCV-2017) | | ✓ | | | | | Softmax Loss | mAP PARSE-27K: 72.76 mAP Acc/Prec/Recall/F1, PA-100K: 74.21/72.19/82.97/82.09/82.53, PETA: 81.77/76.13/84.92/83.24/84.07, RAP: 76.12/65.39/77.53/78.79/78.05 |
| CAM (Guo et al., 2017) (PRL-2017) | | ✓ | | | | | Exponential Loss | BAP: 89.9 mAP WIDER: 82.9 |
| JRL (Wang et al., 2017) (ICCV-2017) | ✓ | | ✓ | | | | Cross-entropy Loss | mAP/Prec/Recall/F1, PETA: 85.67/86.03/85.34/85.42, RAP: 77.81/78.11/78.98/78.58 |
| WPAL (Zhou et al., 2017) (BMVC-2017) | | | | | | | Weighted Cross-entropy Loss | mAP Acc/Prec/Recall/F1, PETA: 85.50/76.98/84.07/85.78/84.90, RAP: 81.25/50.30/57.17/78.39/66.12 |
| AWMT (He et al., 2017) (MM-2017) | | ✓ | | ✓ | | | Cross-entropy Loss | CelebA: 91.80, Duke: 87.53 |
| MTCT (Dong et al., 2017) (WACV-2017) | | | | ✓ | | ✓ | t-STE Loss | Accuracy/Precision/Recall Street data-c: 64.35/64.97/75.66 |
| CILICIA (Sarafianos et al., 2017) (ICCV-2017) | | | | ✓ | | ✓ | Categorical Cross-entropy Loss | SoBIR: 73.1 Accuracy: VIPeR: 80.5 |
| FaFS (Lu et al., 2017) (CVPR-2017) | | | | | | ✓ | Cross-entropy Loss | Accuracy/Top-10 Recall/CelebA: 91.02/71.38 |
| GAM (Fabbri et al., 2017) (AVSS-2017) | ✓ | | | | | | Cross-entropy Loss | mAP Acc/Prec/Recall/F1, RAP: 79.73/83.97/76.96/78.72/77.83 |
| MTA-Net (Ji et al., 2020) (PRL-2020) | | ✓ | ✓ | | | | Focal Balance Loss | mAP Acc/Prec/Recall/F1, RAP: 77.62/67.17/79.72/78.44/79.07, PETA: 84.62/78.80/85.67/86.42/86.04 |

Table 3: An overview of PAR algorithms reviewed in this paper (Part-II).

| Algorithm | Part | Attention | Seq. | C. L. | Graphic | Groups | Loss | Accuracy |
|--|------|-----------|------|-------|---------|--------|------------------------------------|---|
| A-ADG (Park et al., 2018) (TPAMI-2018) | ✓ | | | | ✓ | | - | maP/mAC BAP: 91.684.3 |
| GRL (Zhao et al., 2018) (IJCAI-2018) | ✓ | | ✓ | | | ✓ | Cross-entropy Loss | ma/Prec/Rec/F ₁ , PETA: 86.70/84.34/88.82/86.51 RAP: 81.20/77.70/80.90/79.29 |
| LGNNet (Liu et al., 2018) (BMVC-2018) | ✓ | | | | | | Softmax Loss | ma/Acc/Prec/Rec/F ₁ , RAP: 78.68/68.00/80.36/79.82/80.09, PA-100K: 76.96/73.55/86.99/83.17/85.04 |
| FGDM (Li et al., 2018) (ICME-2018) | ✓ | | | | | | Weighted Cross-entropy Loss | ma/Acc/Prec/Rec/F ₁ , PETA: 82.97/78.08/86.86/84.68/85.76, RAP: 74.31/64.57/78.86/75.90/77.35, PA-100K: 74.95/73.08/84.36/82.24/83.29 |
| DJAA (Sarafianos et al., 2018) (ECCV-2018) | | ✓ | | | | | Weighted Focal Loss | ma/Acc/Prec/Rec/F ₁ , PETA: 84.59/78.56/86.79/86.12/86.46 |
| VSGR (HUANG, 2019) (AAAI-2019) | ✓ | | ✓ | | ✓ | | Cross-entropy Loss | ma/Acc/Prec/Rec/F ₁ , RAP: 77.91/70.04/82.05/80.64/81.34, PA-100K: 79.52/80.58/89.04/87.15/88.26, PETA: 85.21/81.25/88.43/88.42/88.42 |
| RCRA (Xin Zhao and Yan, 2019) (AAAI-2019) | | ✓ | ✓ | | | ✓ | Weighted Cross-entropy Loss | ma/Prec/Rec/F ₁ , RAP: 78.47/82.67/76.65/79.54, PETA: 85.78/85.42/88.02/86.07 |
| $I A^2$ -Net (Ji et al., 2019) (PRL-2019) | | ✓ | ✓ | | | ✓ | Focal Cross-entropy Loss | ma/Acc/Prec/Rec/F ₁ , RAP: 77.44/67.75/79.01/77.45/78.03, PETA: 84.13/78.62/85.73/86.07/85.88 |
| JLPLS-PAA (Tan et al., 2019) (TIP-2019) | | ✓ | | | | | Cross-entropy Loss | ma/Acc/Prec/Rec/F ₁ , RAP: 81.25/67.91/78.56/81.45/79.98, PETA: 84.88/79.46/87.42/86.33/86.87, PA-100K: 81.61/78.89/86.83/87.73/87.27 |
| CoCNN (Kai Han, 2019) (IJCAI-2019) | ✓ | | | | | ✓ | Cross-entropy Loss | ma/Acc/Prec/Rec/F ₁ , RAP: 81.42/68.37/81.04/80.27/80.65, PETA: 86.97/79.95/87.58/87.73/87.65, PA-100K: 80.56/78.30/89.49/84.36/86.85 |
| DCL (Wang et al., 2019) (ICCV-2019) | | | | ✓ | | | Cross-entropy + Triplet Loss | ma: RAP/CelebA: 83.7/89.05 |
| ALM (Tang et al., 2019) (ICCV-2019) | ✓ | ✓ | | | | | Weighted Binary Cross-entropy Loss | ma/Acc/Prec/Rec/F ₁ , RAP: 81.87/68.17/74.71/86.48/80.16, PETA: 86.30/79.52/85.65/88.09/86.85, PA-100K: 80.68/77.08/84.21/88.84/86.46 |
| HAR (Wu et al., 2019) (AAAI-2020) | ✓ | ✓ | | | | ✓ | Cross-entropy Loss | ma/Acc/Prec/Rec/F ₁ , WIDER: maP: 87.3 |
| HFE (Yang et al., 2020) (CVPR-2020) | | | | | | ✓ | Cross-entropy Loss & HFE Loss | Duke: 91.77. Market 501: 92.90 |
| CAS (Zeng et al., 2020) (ICME-2020) | | ✓ | | | | | Cross-entropy Loss | ma/Acc/Prec/Rec/F ₁ , PA-100K: 77.20/78.09/88.46/84.86/86.62 |
| CRM (Tan et al., 2020) (AAAI-2020) | | ✓ | | | | | Cross-entropy Loss | ma/Acc/Prec/Rec/F ₁ , PETA: 86.96/80.38/87.81/87.09/87.45, RAP: 83.69/69.15/79.31/82.40/80.82, PA-100K: 82.31/79.47/87.45/87.77/87.61 |

958 **Further Exploring the Visual Attention Mechanism** Visual attention has
959 drawn more and more researcher’s attention in recent years. It is still one of the
960 most popular techniques used in nowadays and integrated with every kind of deep
961 neural networks in many tasks. Just as noted in (Mnih et al., 2014), one important
962 property of human perception is that one does not tend to process a whole scene
963 in its entirety at once. Instead, humans focus attention selectively on parts of the
964 visual space to acquire information when and where it is needed, and combine in-
965 formation from different fixations over time to build up an internal representation
966 of the scene, guiding future eye movements and decision making. It also substan-
967 tially reduces the task complexity as the object of interest can be placed in the
968 center of the fixation and irrelevant features of the visual environment (“clutter”)
969 outside the fixated region are naturally ignored. Designing novel attention mech-
970 anism or borrowing from other research domains for pedestrian attribute recogni-
971 tion maybe be an important research direction in the future.

972 **Newly Designed Loss Functions** In recent years, there are many loss func-
973 tions proposed for deep neural network optimization, such as (Weighted) Cross
974 Entropy Loss, Contrastive Loss, Center Loss, Triplet Loss, Focal Loss. Re-
975 searchers also design new loss functions for the PAR, such as WPAL and AWMT,
976 to further improving their recognition performance. It is a very important direc-
977 tion to study the influence of different loss functions for PAR.

978 **Exploring More Advanced Network Architecture** Existing PAR models
979 adopts off the shelf pre-trained network on large scale dataset, as their backbone
980 network architecture. Seldom of them consider the unique characteristics of PAR
981 and design novel networks. Some novel networks are proposed in recent years,
982 such as capsule network, however, there are still no attempts to use such networks
983 for PAR. There are also works demonstrating that the deeper network architec-
984 ture the better recognition performance we can obtain. Nowadays, Automatic
985 Machine Learning solutions (AutoML) draw more and more attentions and many
986 development tools are also released for the development, such as: AutoWEKA
987 and Auto-sklearn. Therefore, it will be a good choice to design specific networks
988 for person attribute recognition in future works with aforementioned approaches.

989 **Prior Knowledge guided Learning** Different from regular classification task,
990 pedestrian attribute recognition always have its own characteristics due to the pref-
991 erence of human beings or natural constraints. It is an important research direction
992 to mining the prior or common knowledge for the PAR. For example, we wear dif-
993 ferent clothes in various seasons, temperatures or occasions. On the other hand,

994 some researchers attempt to use the history knowledge (such as: Wikipedia³) to
995 help improve their overall performance. Therefore, how to use this information
996 to explore the relations between person attributes or help the machine learning
997 model to further understanding the attributes is still an unstudied problem.

998 **Multi-modal Pedestrian Attribute Recognition** Although existing single-
999 modal algorithms already achieve good performance on some benchmark dataset
1000 as mentioned above. However, as is known to all, the RGB image is sensitive
1001 to illumination, bad weather (such as: rain, snow, fog), night time, *etc.* It seems
1002 impossible for us to achieve accurate pedestrian attribute recognition in all day
1003 and all weather. But the actual requirement of intelligent surveillance needs far
1004 more than this target. How can we bridge this gap? One intuitive idea is to mine
1005 useful information from other modalities, such as thermal or depth sensors, to
1006 integrate with RGB sensor. There are already many works attempt to fuse these
1007 multi-modal data and improve their final performance significantly. We think the
1008 idea of multi-modal fusion could also help improve the robustness of pedestrian
1009 attribute recognition. The thermal images can highlight the contour of human and
1010 some other wearing or carrying objects.

1011 **Video based Pedestrian Attribute Recognition** Existing pedestrian attribute
1012 recognition is based on single image, however, we often obtain the video sequence
1013 captured by cameras in practical scenario. Although running existing algorithm on
1014 each video frame can be an intuitive and easy strategy, but the efficiency maybe
1015 the bottleneck for practical applications. Generally speaking, image based at-
1016 tribute recognition can only make use of the spatial information from the given
1017 image, which increases the difficulty of PAR due to the limited information. In
1018 contrast, given the video based PAR, we can jointly utilize the spatial and temporal
1019 information. The benefits can be listed as follows: 1). we can extend the attribute
1020 recognition into a more general case by defining more dynamic person attributes,
1021 such as “running man”; 2). the motion information can be used to reason the at-
1022 tributes which maybe hard to recognize in single image; 3). the general person
1023 attributes learned in videos can provide more helpful information for other video
1024 based tasks, such as video caption, video object detection. Therefore, how to rec-
1025 ognize human attributes in practical video sequence efficiently and accurately is a
1026 problem worth studying.

1027 **Joint Learning of Attribute and Other Tasks** Integrating the person attribute
1028 learning into the pipeline of other person related tasks is also an interesting and

³en.wikipedia.org

1029 important research direction. There are already many algorithms proposed by
1030 considering the person attributes into corresponding tasks, such as: attribute based
1031 pedestrian detection, visual tracking, person re-identification and social activity
1032 analysis. In the future, how to better explore the fine-grained person attributes for
1033 other tasks and also use other tasks for better human attribute recognition is an
1034 important research directions.

1035 **7. Conclusion**

1036 In this paper, we give a review of PAR from traditional approaches to deep
1037 learning based algorithms in recent years. Specifically, we first introduce the back-
1038 ground (problem formulation and challenging factors) of PAR. Then, we give a
1039 review of PAR algorithms from different perspectives, including: global based,
1040 part based, visual attention based, sequential prediction based, newly designed
1041 loss function based, curriculum learning based, graphic model based and other
1042 algorithms. After that, we discuss the specific attribute recognition, then, give
1043 a comparison between deep learning and traditional algorithm based PAR meth-
1044 ods. After that, we show the connections between PAR and other computer vision
1045 tasks. We summarize existing benchmarks proposed for PAR, including popular
1046 datasets and evaluation criteria, and also give a brief comparison of selected 17
1047 PAR algorithms on RAP and PETA dataset. Finally, we summarize this paper
1048 and give several possible research directions for PAR. However, due to the limited
1049 space in this paper, there are still many other works that may be related to PAR
1050 but not covered in this survey. For example, the history of the backbone deep
1051 networks used in deep PAR algorithms, the various machine learning techniques
1052 such as transfer learning, self-supervised learning, meta-learning, or active learn-
1053 ing which may inspire the researchers to design more advanced PAR algorithms.
1054 In our future works, we will summarize these techniques which may be useful for
1055 pedestrian attribute recognition.

1056 **Acknowledgements:** This work is jointly supported by Postdoctoral Innovative Tal-
1057 ent Support Program BX20200174, China Postdoctoral Science Foundation Funded Project
1058 2020M682828. National Nature Science Foundation of China (61976002, 62076003,
1059 61860206004), Australian Research Council Projects FL-170100117. We also thanks
1060 all the reviewers, AE and EiC for their valuable comments and suggestions.

1061 **References**

1062 Z. Chen, W. Ouyang, T. Liu, D. Tao, A shape transformation-based dataset augmentation
1063 framework for pedestrian detection, IJCV 129 (2021) 1121–1138.

- 1064 Y. Deng, P. Luo, C. C. Loy, X. Tang, Pedestrian attribute recognition at far distance, in:
1065 Proceedings of the 22nd ACM MM, 2014, pp. 789–792.
- 1066 P. Sudowe, H. Spitzer, B. Leibe, Person attribute recognition with a jointly-trained holistic
1067 cnn model, in: IEEE ICCV Workshops, 2015, pp. 87–95.
- 1068 D. Li, X. Chen, K. Huang, Multi-attribute learning for pedestrian attribute recognition in
1069 surveillance scenarios, in: ACPR, IEEE, 2015, pp. 111–115.
- 1070 A. H. Abdalnabi, G. Wang, J. Lu, K. Jia, Multi-task cnn model for attribute prediction,
1071 IEEE TMM 17 (2015) 1949–1959.
- 1072 H. S. C. L. C. X. C. X. Kai Han, Yunhe Wang, Attribute aware pooling for pedestrian
1073 attribute recognition, in: IJCAI, 2019.
- 1074 L. Bourdev, S. Maji, J. Malik, Describing people: A poselet-based approach to attribute
1075 classification, in: IEEE ICCV, 2011, pp. 1543–1550.
- 1076 J. Joo, S. Wang, S.-C. Zhu, Human attribute recognition by rich appearance dictionary,
1077 in: IEEE ICCV, 2013, pp. 721–728.
- 1078 S. Lazebnik, C. Schmid, J. Ponce, Beyond bags of features: Spatial pyramid matching for
1079 recognizing natural scene categories, in: IEEE CVPR, volume 2, 2006, pp. 2169–2178.
- 1080 N. Zhang, M. Paluri, M. Ranzato, T. Darrell, L. Bourdev, Panda: Pose aligned networks
1081 for deep attribute modeling, in: IEEE CVPR, 2014, pp. 1637–1644.
- 1082 G. Gkioxari, R. Girshick, J. Malik, Actions and attributes from wholes and parts, in:
1083 IEEE ICCV, 2015.
- 1084 R. Girshick, F. Iandola, T. Darrell, J. Malik, Deformable part models are convolutional
1085 neural networks, in: IEEE CVPR, 2015, pp. 437–446.
- 1086 J. Zhu, S. Liao, D. Yi, Z. Lei, S. Z. Li, Multi-label cnn based pedestrian attribute learning
1087 for soft biometrics, in: IEEE ICB, 2015, pp. 535–540.
- 1088 C. Tang, L. Sheng, Z. Zhang, X. Hu, Improving pedestrian attribute recognition with
1089 weakly-supervised multi-scale attribute-specific localization, in: IEEE ICCV, 2019,
1090 pp. 4997–5006.
- 1091 Y. W. S. L. Luwei Yang, Ligeng Zhu, P. Tan, Attribute recognition from adaptive parts,
1092 in: BMVC, 2016, pp. 81.1–81.11. doi:[10.5244/C.30.81](https://doi.org/10.5244/C.30.81).

- 1093 A. Diba, A. Mohammad Pazandeh, H. Pirsiavash, L. Van Gool, Deepcamp: Deep convo-
1094 lutional action & attribute mid-level patterns, in: IEEE CVPR, 2016, pp. 3557–3565.
- 1095 D. Li, X. Chen, Z. Zhang, K. Huang, Pose guided deep model for pedestrian attribute
1096 recognition in surveillance scenarios, in: IEEE ICME, 2018, pp. 1–6.
- 1097 Y. Li, C. Huang, C. C. Loy, X. Tang, Human attribute recognition by deep hierarchical
1098 contexts, in: ECCV, Springer, 2016, pp. 684–700.
- 1099 P. Liu, X. Liu, J. Yan, J. Shao, Localization guided learning for pedestrian attribute
1100 recognition, in: BMVC, 2018.
- 1101 X. Liu, H. Zhao, M. Tian, L. Sheng, J. Shao, S. Yi, J. Yan, X. Wang, Hydraplus-net:
1102 Attentive deep features for pedestrian analysis, in: IEEE ICCV, 2017, pp. 350–359.
- 1103 M. S. Sarfraz, A. Schumann, Y. Wang, R. Stiefelhagen, Deep view-sensitive pedestrian
1104 attribute inference in an end-to-end model, arXiv:1707.06089 (2017).
- 1105 N. Sarafianos, X. Xu, I. A. Kakadiaris, Deep imbalanced attribute classification using
1106 visual attention aggregation, in: ECCV, Springer, 2018, pp. 708–725.
- 1107 H. Guo, X. Fan, S. Wang, Human attribute recognition by refining attention heat map,
1108 Pattern Recognition Letters 94 (2017) 38–45.
- 1109 Z. Tan, Y. Yang, J. Wan, H. Wan, G. Guo, S. Z. Li, Attention based pedestrian attribute
1110 analysis, IEEE TIP (2019).
- 1111 Z. Ji, E. He, H. Wang, A. Yang, Image-attribute reciprocally guided attention network for
1112 pedestrian attribute recognition, Pattern Recognition Letters 120 (2019) 89–95.
- 1113 M. Wu, D. Huang, Y. Guo, Y. Wang, Distraction-aware feature learning for human at-
1114 tribute recognition via coarse-to-fine attention mechanism, arXiv:1911.11351 (2019).
- 1115 H. Zeng, H. Ai, Z. Zhuang, L. Chen, Multi-task learning via co-attentive sharing for
1116 pedestrian attribute recognition, in: ICME, IEEE, 2020, pp. 1–6.
- 1117 S. Zhang, Z. Song, X. Cao, H. Zhang, J. Zhou, Task-aware attention model for clothing
1118 attribute prediction, TCSVT 30 (2019) 1051–1064.
- 1119 J. Wang, Y. Yang, J. Mao, Z. Huang, C. Huang, W. Xu, Cnn-rnn: A unified framework
1120 for multi-label image classification, in: IEEE CVPR, 2016, pp. 2285–2294.
- 1121 J. Wang, X. Zhu, S. Gong, W. Li, Attribute recognition by joint recurrent learning of
1122 context and correlation, in: IEEE ICCV, 2017, pp. 531–540.

- 1123 X. Zhao, L. Sang, G. Ding, Y. Guo, X. Jin, Grouping attribute recognition for pedestrian
1124 with joint recurrent learning., in: IJCAI, 2018, pp. 3177–3183.
- 1125 H. Liu, J. Wu, J. Jiang, M. Qi, R. Bo, Sequence-based person attribute recognition with
1126 joint ctc-attention model, arXiv preprint arXiv:1811.08115 (2018).
- 1127 G. D. J. H. N. D. Xin Zhao, Liufang Sang, C. Yan, Recurrent attention model for pedes-
1128 trian attribute recognition, in: AAAI, 2019.
- 1129 Y. Zhou, K. Yu, B. Leng, Z. Zhang, D. Li, K. Huang, B. Feng, C. Yao, et al., Weakly-
1130 supervised learning of mid-level features for pedestrian attribute recognition and local-
1131 ization, in: BMVC, 2017.
- 1132 K. He, Z. Wang, Y. Fu, R. Feng, Y.-G. Jiang, X. Xue, Adaptively weighted multi-task
1133 deep network for person attribute classification, in: ACM MM, 2017, pp. 1636–1644.
- 1134 E. Yaghoubi, D. Borza, J. Neves, A. Kumar, H. Proença, An attention-based deep learning
1135 model for multiple pedestrian attributes recognition, arXiv preprint arXiv:2004.01110
1136 (2020).
- 1137 J. Yang, J. Fan, Y. Wang, Y. Wang, W. Gan, L. Liu, W. Wu, Hierarchical feature embed-
1138 ding for attribute recognition, in: IEEE CVPR, 2020, pp. 13055–13064.
- 1139 J. Jia, H. Huang, W. Yang, X. Chen, K. Huang, Rethinking of pedestrian attribute recogni-
1140 tion: Realistic datasets with efficient method, arXiv preprint arXiv:2005.11909 (2020).
- 1141 Z. Ji, Z. Hu, E. He, J. Han, Y. Pang, Pedestrian attribute recognition based on multiple
1142 time steps attention, Pattern Recognition Letters (2020).
- 1143 Q. Dong, S. Gong, X. Zhu, Multi-task curriculum transfer deep learning of clothing
1144 attributes, in: IEEE WACV, 2017, pp. 520–529.
- 1145 N. Sarafianos, T. Giannakopoulos, C. Nikou, I. A. Kakadiaris, Curriculum learning for
1146 multi-task classification of visual attributes, in: IEEE ICCV, 2017, pp. 2608–2615.
- 1147 J. Zhu, S. Liao, Z. Lei, S. Z. Li, Multi-label convolutional neural network based pedestrian
1148 attribute classification, IVC 58 (2017) 224–229.
- 1149 N. Sarafianos, T. Giannakopoulos, C. Nikou, I. A. Kakadiaris, Curriculum learning of
1150 visual attribute clusters for multi-task classification, Pattern Recognition 80 (2018)
1151 94–108.
- 1152 Y. Wang, W. Gan, W. Wu, J. Yan, Dynamic curriculum learning for imbalanced data
1153 classification, ICCV (2019).

- 1154 H. Chen, A. Gallagher, B. Girod, Describing clothing by semantic attributes, in: ECCV,
1155 Springer, 2012, pp. 609–623.
- 1156 S. Park, B. X. Nie, S.-C. Zhu, Attribute and-or grammar for joint parsing of human pose,
1157 parts and attributes, IEEE TPAMI 40 (2018) 1555–1569.
- 1158 Q. L. X. Z. R. H. K. HUANG, Visual-semantic graph reasoning for pedestrian attribute
1159 recognition, in: AAAI, 2019.
- 1160 Z. Tan, Y. Yang, J. Wan, G. Guo, S. Z. Li, Relation-aware pedestrian attribute recognition
1161 with graph convolutional networks., in: AAAI, 2020, pp. 12055–12062.
- 1162 W. Wang, Y. Xu, J. Shen, S.-C. Zhu, Attentive fashion grammar network for fashion
1163 landmark detection and clothing category classification, in: IEEE CVPR, 2018, pp.
1164 4271–4280.
- 1165 P. Sudowe, B. Leibe, Patchit: Self-supervised network weight initialization for fine-
1166 grained recognition., in: BMVC, 2016.
- 1167 Y. Lu, A. Kumar, S. Zhai, Y. Cheng, T. Javidi, R. Feris, Fully-adaptive feature sharing
1168 in multi-task networks with applications in person attribute classification, in: CVPR,
1169 volume 1, 2017, p. 6.
- 1170 M. Fabbri, S. Calderara, R. Cucchiara, Generative adversarial models for people attribute
1171 recognition in surveillance, in: IEEE AVSS, 2017, pp. 1–6.
- 1172 G. D. J. H. L. L. Liuyu Xiang, Xiaoming Jin, Incremental few-shot learning for pedestrian
1173 attribute recognition, in: IJCAI, 2019.
- 1174 M. Mirza, S. Osindero, Conditional generative adversarial networks, Manuscript:
1175 <https://arxiv.org/abs/1709.02023> (2014).
- 1176 Y. Zhang, P. Zhang, C. Yuan, Z. Wang, Texture and shape biased two-stream networks
1177 for clothing classification and attribute recognition, in: IEEE CVPR, 2020, pp. 13538–
1178 13547.
- 1179 J. Jia, H. Huang, X. Chen, K. Huang, Rethinking of pedestrian attribute recogni-
1180 tion: A reliable evaluation under zero-shot pedestrian identity setting, arXiv preprint
1181 arXiv:2107.03576 (2021).
- 1182 X. Zheng, Y. Guo, H. Huang, Y. Li, R. He, A survey to deep facial attribute analysis,
1183 arXiv preprint arXiv:1812.10265 (2018).

- 1184 B. Fasel, J. Luetttin, Automatic facial expression analysis: a survey, *Pattern recognition*
1185 36 (2003) 259–275.
- 1186 P. Rodríguez, G. Cucurull, J. M. Gonfaus, F. X. Roca, J. Gonzalez, Age and gender
1187 recognition in the wild with deep attention, *Pattern Recognition* 72 (2017) 563–571.
- 1188 K. Li, J. Xing, W. Hu, S. J. Maybank, D2c: Deep cumulatively and comparatively learning
1189 for human age estimation, *Pattern Recognition* 66 (2017) 95–105.
- 1190 J. Xing, K. Li, W. Hu, C. Yuan, H. Ling, Diagnosing deep learning models for high
1191 accuracy age estimation from a single image, *Pattern Recognition* 66 (2017) 106–116.
- 1192 G. Antipov, M. Baccouche, S.-A. Berrani, J.-L. Dugelay, Effective training of convolu-
1193 tional neural networks for face-based gender and age prediction, *Pattern Recognition*
1194 72 (2017) 15–26.
- 1195 H. Liu, J. Lu, J. Feng, J. Zhou, Group-aware deep feature learning for facial age estima-
1196 tion, *Pattern Recognition* 66 (2017) 82–94.
- 1197 K. Chen, K. Jia, Z. Zhang, J.-K. Kämäräinen, Spectral attribute learning for visual regres-
1198 sion, *Pattern Recognition* 66 (2017) 74–81.
- 1199 A. Hadid, M. Pietikäinen, Combining appearance and motion for face and gender recog-
1200 nition from videos, *Pattern Recognition* 42 (2009) 2818–2827.
- 1201 A. Branca, M. Leo, G. Attolico, A. Distanto, Detection of objects carried by people, in:
1202 *IEEE ICIP*, volume 3, 2002.
- 1203 F. Ghadiri, R. Bergevin, G.-A. Bilodeau, From superpixel to human shape modelling for
1204 carried object detection, *Pattern Recognition* 89 (2019) 134–150.
- 1205 D. Damen, D. Hogg, Detecting carried objects from sequences of walking pedestrians,
1206 *IEEE TPAMI* 34 (2011) 1056–1067.
- 1207 F. Ghadiri, R. Bergevin, G.-A. Bilodeau, Carried object detection based on an ensemble
1208 of contour exemplars, in: *ECCV*, Springer, 2016, pp. 852–866.
- 1209 Y. Tian, P. Luo, X. Wang, X. Tang, Pedestrian detection aided by deep learning seman-
1210 tic tasks, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern*
1211 *Recognition*, 2015, pp. 5079–5087.
- 1212 Y. Lin, L. Zheng, Z. Zheng, Y. Wu, Z. Hu, C. Yan, Y. Yang, Improving person re-
1213 identification by attribute and identity learning, *Pattern Recognition* (2019).

- 1214 K. Han, J. Guo, C. Zhang, M. Zhu, Attribute-aware attention model for fine-grained
1215 representation learning, in: ACM MM, 2018, pp. 2040–2048.
- 1216 C. Su, S. Zhang, J. Xing, W. Gao, Q. Tian, Deep attributes driven multi-camera person
1217 re-identification, in: ECCV, Springer, 2016, pp. 475–491.
- 1218 S. Khamis, C.-H. Kuo, V. K. Singh, V. D. Shet, L. S. Davis, Joint learning for attribute-
1219 consistent person re-identification, in: ECCV, Springer, 2014, pp. 134–146.
- 1220 A. Li, L. Liu, K. Wang, S. Liu, S. Yan, Clothing attributes assisted person reidentification,
1221 IEEE TCSVT 25 (2015) 869–878.
- 1222 R. Layne, T. M. Hospedales, S. Gong, Towards person identification and re-identification
1223 with attributes, in: ECCV, Springer, 2012, pp. 402–412.
- 1224 R. Layne, T. M. Hospedales, S. Gong, Attributes-based re-identification, in: Person
1225 Re-Identification, Springer, 2014, pp. 93–117.
- 1226 R. Layne, T. M. Hospedales, S. Gong, Q. Mary, Person re-identification by attributes., in:
1227 BMVC, volume 2, 2012, p. 8.
- 1228 A. Schumann, R. Stiefelhagen, Person re-identification by deep learning attribute-
1229 complementary information, in: IEEE CVPR Workshops, 2017, pp. 20–28.
- 1230 S. Li, H. Yu, W. Huang, J. Zhang, Attributes-aided part detection and refinement for
1231 person re-identification, arXiv preprint arXiv:1902.10528 (2019).
- 1232 H. Ling, Z. Wang, P. Li, Y. Shi, J. Chen, F. Zou, Improving person re-identification by
1233 multi-task learning, Neurocomputing 347 (2019) 109–118.
- 1234 C. Su, S. Zhang, F. Yang, G. Zhang, Q. Tian, W. Gao, L. S. Davis, Attributes driven
1235 tracklet-to-tracklet person re-identification using latent prototypes space mapping, Pat-
1236 tern Recognition 66 (2017) 4–15.
- 1237 Y. Chen, S. Duffner, A. Stoian, J.-Y. Dufour, A. Baskurt, Deep and low-level feature
1238 based attribute learning for person re-identification, IVC 79 (2018) 25–34.
- 1239 X. Wang, T. Zhang, D. R. Tretter, Q. Lin, Personal clothing retrieval on photo collections
1240 by color and attributes, IEEE TMM 15 (2013) 2035–2045.
- 1241 Q. Chen, J. Huang, R. Feris, L. M. Brown, J. Dong, S. Yan, Deep domain adaptation for
1242 describing people based on fine-grained clothing attributes, in: IEEE CVPR, 2015, pp.
1243 5315–5324.

- 1244 M. Ziaeefard, R. Bergevin, Semantic human activity recognition: A literature review,
1245 Pattern Recognition 48 (2015) 2329–2345.
- 1246 J. Liu, B. Kuipers, S. Savarese, Recognizing human actions by attributes, in: IEEE CVPR,
1247 2011, pp. 3337–3344.
- 1248 E. Murphy-Chutorian, M. M. Trivedi, Head pose estimation in computer vision: A survey,
1249 IEEE TPAMI 31 (2008) 607–626.
- 1250 L. Huang, J. Peng, R. Zhang, G. Li, L. Lin, Learning deep representations for semantic
1251 image parsing: a comprehensive overview, Frontiers of Computer Science 12 (2018)
1252 840–857.
- 1253 W. Xiao, Z. Shaofei, Y. Rui, L. Bin, T. Jin, Pedestrian attribute recognition: A survey,
1254 arXiv:1901.07474 (2019).
- 1255 Y. Xiong, K. Zhu, D. Lin, X. Tang, Recognize complex events from static images by
1256 fusing deep channels, in: CVPR, 2015, pp. 1600–1609.
- 1257 A. Geiger, P. Lenz, R. Urtasun, Are we ready for autonomous driving? the kitti vision
1258 benchmark suite, in: CVPR, 2012, pp. 3354–3361.
- 1259 S. M. Bileschi, StreetScenes: Towards scene understanding in still images, Technical
1260 Report, MASSACHUSETTS INST OF TECH CAMBRIDGE, 2006.
- 1261 N. Dalal, B. Triggs, Histograms of oriented gradients for human detection, in: CVPR,
1262 volume 1, 2005, pp. 886–893.
- 1263 L. Bourdev, J. Malik, Poselets: Body part detectors trained using 3d human pose annota-
1264 tions, in: ICCV, 2009, pp. 1365–1372.
- 1265 T. Li, J. Liu, W. Zhang, Y. Ni, W. Wang, Z. Li, Uav-human: A large benchmark for
1266 human behavior understanding with unmanned aerial vehicles, in: IEEE CVPR, 2021,
1267 pp. 16266–16275.
- 1268 J. Zhu, S. Liao, Z. Lei, D. Yi, S. Li, Pedestrian attribute classification in surveillance:
1269 Database and evaluation, in: IEEE ICCV Workshops, 2013, pp. 331–338.
- 1270 D. Li, Z. Zhang, X. Chen, H. Ling, K. Huang, A richly annotated dataset for pedestrian
1271 attribute recognition, arXiv:1603.07054 (2016).
- 1272 X. Wang, C. Li, B. Luo, J. Tang, Sint++: Robust visual tracking via adversarial positive
1273 instance generation, in: IEEE CVPR, 2018, pp. 4864–4873.
- 1274 V. Mnih, N. Heess, A. Graves, et al., Recurrent models of visual attention, in: NIPS,
1275 2014, pp. 2204–2212.