# Improving Pedestrian Attribute Recognition With Weakly-Supervised Multi-Scale Attribute-Specific Localization

Chufeng Tang[1]   Lu Sheng[2]   Zhaoxiang Zhang[3]   Xiaolin Hu[1]*

[1]State Key Laboratory of Intelligent Technology and Systems,
Institute for Artificial Intelligence, Department of Computer Science and Technology,
Beijing National Research Center for Information Science and Technology, Tsinghua University
[2]College of Software, Beihang University  [3]Institute of Automation, Chinese Academy of Sciences

{tcf18@mails,xlhu@mail}.tsinghua.edu.cn   lsheng@buaa.edu.cn   zhaoxiang.zhang@ia.ac.cn

## Abstract

*Pedestrian attribute recognition has been an emerging research topic in the area of video surveillance. To predict the existence of a particular attribute, it is demanded to localize the regions related to the attribute. However, in this task, the region annotations are not available. How to carve out these attribute-related regions remains challenging. Existing methods applied attribute-agnostic visual attention or heuristic body-part localization mechanisms to enhance the local feature representations, while neglecting to employ attributes to define local feature areas. We propose a flexible Attribute Localization Module (ALM) to adaptively discover the most discriminative regions and learns the regional features for each attribute at multiple levels. Moreover, a feature pyramid architecture is also introduced to enhance the attribute-specific localization at low-levels with high-level semantic guidance. The proposed framework does not require additional region annotations and can be trained end-to-end with multi-level deep supervision. Extensive experiments show that the proposed method achieves state-of-the-art results on three pedestrian attribute datasets, including PETA, RAP, and PA-100K.*

## 1. Introduction

Recognition of pedestrian attributes, *e.g.* gender, age, and clothing style, has drawn extensive attention because of its great potential in video surveillance applications, such as face verification [10], person retrieval [2, 27], and person re-identification [11, 22, 30]. Recently, methods based on the Convolutional Neural Networks (CNN) [6, 8] achieve great success in pedestrian attribute recognition by learning powerful features from images. Some existing works [13, 28]
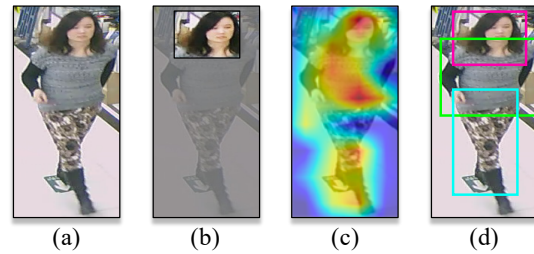
---

*Corresponding author.



Figure 1. Attentive regions generated by different methods when recognizing the attribute *Longhair*. (a) The original input image. (b) Attribute-specific region generated by our proposed method, which is indeed localized into a head-related region. (c) Attention mask generated by attribute-agnostic attention methods [20, 24, 38], which covers a broad region but not specific to *Longhair*. (d) Body parts generated by part-based methods [15, 19, 34, 35], which extract features from these body parts.

treat pedestrian attribute recognition as a multi-label classification problem and extract feature representations only from the whole input images. These holistic methods usually rely on global features, but regional features are more significant for fine-grained attribute classification.

Intuitively, attributes can be localized into some relevant regions in a pedestrian image. As illustrated in Figure 1 (b), when recognizing *Longhair*, it is reasonable to focus on the head-related regions. Recent methods attempt to leverage the attention localization to promote learning discriminative features for attribute recognition. A popular solution [20, 24, 38] is to employ the visual attention mechanism to capture the most relevant features. These methods usually generate attention masks from certain layers and then multiply them to corresponded feature maps so as to extract the attentive features. However, it is ambiguous which mask encodes a given attribute's location, and there is no specific mechanism that guarantees the correspondences between attributes and attention masks. As shown in Figure 1 (c), the learned attention mask attends to a broad region which

is not specific to the required attribute *Longhair*. An alternative way is to leverage predefined rigid parts [40] or external part localization modules [15, 19, 34, 35]. Some works apply body-parts detection [35], pose estimation [15, 34] and region proposals [19] to learn part-based local features. As shown in Figure 1 (d), these methods extract local features from the localized body parts (*e.g.* head, torso, and legs). However, most of them just fuse the part-based features with global features, which still fail to indicate the attribute-region correspondence but require extra computational resources for sophisticated part localization.

Different from these methods, we propose a flexible *Attribute Localization Module* (ALM) that can automatically discover the discriminative regions and extract region-based feature representations in an attribute-specific manner. Specifically, the ALM consists of a tiny channel-attention sub-network to fully exploit the inter-channel dependencies of the input features, followed by a spatial transformer [9] to localize the attribute-specific regions adaptively. Moreover, we embed multiple ALMs at different feature levels and introduce a feature pyramid architecture by integrating high-level semantics to reinforce the attribute localization at low-levels. In addition, ALMs at different feature levels are trained by the same set of attribute supervisions, called *deep supervision* [12, 32], where the final predictions are obtained through a voting scheme to output the maximum responses across different feature levels. This voting scheme will suggest a best prediction occurs in one feature level that has the most accurate attribute region, without interference of negative features from inappropriate regions. The proposed framework is end-to-end trainable and requires only image-level annotations. The contributions of this work can be summarized as follows:

- We propose an end-to-end trainable framework which performs attribute-specific localization at multiple scales to discover the most discriminative attribute regions in a weakly-supervised manner.
- We propose a feature pyramid architecture by leveraging both low-level details and high-level semantics to enhance the multi-scale attribute localization and region-based feature learning in a mutually reinforcing manner. The multi-scale attribute predictions are further fused by an effective voting scheme.
- We conduct extensive experiments on three publicly available pedestrian attribute datasets (PETA [1], RAP [16], and PA-100K [20]) and achieve significant improvement over the previous state-of-the-art methods.

## 2. Related Works

**Pedestrian Attribute Recognition**. Earlier pedestrian attribute recognition methods [1, 11, 39] rely on hand-crafted features such as color and texture histograms, and trained separately. However, the performance of these traditional methods is far from satisfactory. More recently, methods based on the Convolutional Neural Networks achieved great success in pedestrian attribute recognition. Wang *et al.* [31] give a brief review of these methods. Sudowe *et al.* [28] propose a holistic CNN model to jointly learn different attributes. Li *et al.* [13] formulate pedestrian attribute recognition as a multi-label classification problem and propose an improved cross-entropy loss function. However, the performance of these holistic methods is limited due to the lack of consideration of the prior information in attributes. Some recent approaches attempt to exploit the spatial relations and semantic relations among attributes to further improve the recognition performance. These methods can be classified into three basic categories: (1) **Relation-based**: Some works [29, 37] exploit semantic relations to assist attribute recognition. Wang *et al.* [29] propose a CNN-RNN based framework to exploit the interdependency and correlation among attributes. Zhao *et al.* [37] divide the attributes into several groups and attempt to explore the intra-group and inter-group relationships. However, these methods require manually defined rules, *e.g.* prediction order, attribute group, which are hard to determine in real applications. (2) **Attention-based**: Some researchers [20, 24, 25, 38] introduce the visual attention mechanism in attribute recognition. Liu *et al.* [20] propose a multi-directional attention model to learn multi-scale attentive features for pedestrian analysis. Sarafianos *et al.* [24] extend the spatial regularization module [38] to learn effective attention maps at multiple scales. Although recognition accuracy has been improved, these methods are attribute-agnostic and fail to take the attribute-specific information into consideration. (3) **Part-based**: The part-based methods usually extract features from some localized body-parts. Zhu *et al.* [40] divide the whole image into 15 rigid patches and fuse features from different patches. Yang *et al.* [34] and Li *et al.* [15] leverage external pose estimation module to localize body-parts. Liu *et al.* [19] also explore attribute regions in a weakly supervised manner while they assign attribute regions to some fixed proposals generated by EdgeBoxes [42] in advance, which is not fully-adaptive and end-to-end trainable. These methods rely either on predefined rigid parts or on sophisticated part localization mechanisms, which are less robust to pose variances and require extra computational resources. By contrast, the proposed method localizes the most discriminative regions in an attribute-specific manner, which is not considered in most of the existing works.

**Weakly Supervised Attention Localization**. In addition to pedestrian attribute recognition, the idea of performing attention localization without region annotations is also extensively investigated in other visual tasks. Jaderberg *et al.* [9] propose the well-known *Spatial Transformer Network* (STN) which can extract attentional regions with any spatial transformation in an end-to-end trainable manner.
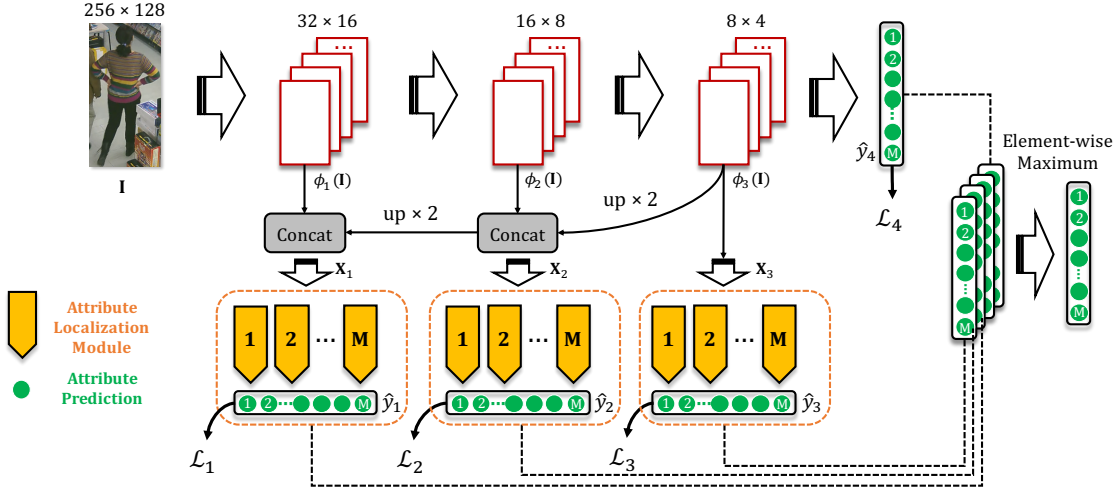
Figure 2. Overview of the proposed framework. The input pedestrian image is fed into the main network with both bottom-up and top-down pathways. Features combined from different levels are fed into multiple *Attribute Localization Modules* (Figure 3), which perform attribute-specific localization and region-based feature learning. Outputs from different branches are trained with *deep supervision* and aggregated through an element-wise maximum operation for inference. $M$ is the total number of attributes. Best viewed in color.
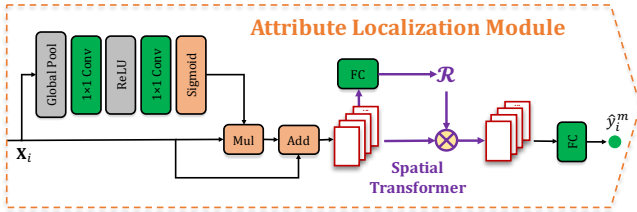


Figure 3. Details of the proposed *Attribute Localization Module* (ALM), which consists of a tiny channel-attention sub-network and a simplified spatial transformer. The ALM takes the combined features $\mathbf{X}_i$ as input and produces an attribute-specific prediction. Each ALM only serves one attribute at a singe level.

Some recent works [14, 17] adopt STN to localize body-parts for person re-identification. Fu *et al*. [3] attempt to recursively learn discriminative region for fine-grained image recognition. Wang *et al*. [33] search the discriminative regions with STN and LSTM for multi-label classification, while not in a label-specific manner. The proposed method is inspired by these works but can adaptively localize the individual informative regions for each attribute.

**Feature Pyramid Architecture**. There are several works exploiting top-down or skip connections that incorporate features across levels, *e.g.* U-Net [23], Stacked hourglass network [21]. The proposed feature pyramid architecture is similar to *Feature Pyramid Networks* (FPN) [18], which have been studied in various object detection and segmentation models [26, 41]. To the best of our knowledge, this work is the first attempt of employing these ideas to localize attentive regions for pedestrian attribute recognition.

## 3. Proposed Method

The overview of the proposed framework is illustrated in Figure 2. As shown, the proposed framework consists of a main network with feature pyramid structures, and a group of *Attribute Localization Modules* (ALM) applied to different feature levels. The input pedestrian image is first fed into the main network without additional region annotations, and a prediction vector is obtained at the end of the bottom-up pathway. The details of ALM are shown in Figure 3. Each ALM only perform attribute localization and region-based feature learning for one attribute at a single feature level. The ALMs at different feature levels are trained in a *deep supervision* manner. Formally, given an input pedestrian image $\mathbf{I}$ along with its corresponding attribute labels $y = \left[y^1, y^2, \ldots, y^M\right]^T$ where $M$ is ths total number of attributes in the dataset and $y^m, m \in 1, \ldots, M$ is a binary label that indicates the presence of the $m$-th attribute if $y^m = 1$, and $y^m = 0$ otherwise. We adopt the BN-Inception [8] architecture as the backbone network in our framework. In principle, the backbone can be replaced with any other CNN architecture. Implementation details are shown in Appendix A.

### 3.1. Network Architecture

The key idea of this work is to perform attribute-specific localization for improving attribute recognition. It is well known that features in deeper CNN layers have coarser resolutions. Even though we can precisely localize the attribute regions based on semantically stronger features, it is still difficult to extract region-based discriminative features since some finer details may disappear. In contrast,

features in lower layers always capture rich details but poor contextual information, resulting in unreliable attribute localization. Obviously, low-level details and high-level semantics are complementary to each other. Therefore, we propose a feature pyramid architecture, inspired by the FPN alike models [18, 41], to enhance the attribute localization and region-based feature learning in a mutually reinforcing manner. As illustrated in Figure 2, the proposed feature pyramid architecture consists of a bottom-up pathway and a top-down pathway.

The bottom-up pathway, implemented by BN-Inception network, consists of multiple `inception` blocks with different feature levels. In this paper, we conduct attribute localization with bottom-up features generated from three different levels: the `incep_3b`, `incep_4d`, and `incep_5b` block respectively, where they have strides of $\{8, 16, 32\}$ pixels with respect to the input image. The selected `inception` blocks are both at the end of their corresponded stages, where blocks of the same stage keep the same feature maps resolution, since we believe the last block should have strongest features. Given an input image $\mathbf{I}$, we denote the bottom-up features generated from the above blocks as $\phi_i(\mathbf{I}) \in \mathbb{R}^{H_i \times W_i \times C_i}, i \in \{1, 2, 3\}$. For $256 \times 128$ RGB input images, the spatial size $H_i \times W_i$ equal to $32 \times 16$, $16 \times 8$, and $8 \times 4$ respectively.

In addition, the top-down pathway contains three lateral connections and two top-down connections, as shown in Figure 2. The lateral connections are simply used to reduce the dimensionalities of bottom-up features to $d$, where $d = 256$ in our implementation. The higher level features are transmitted through the top-down connections and meanwhile go through an upsampling operation. Afterward, features from adjacent levels are concatenated as follows:

$$\mathbf{X}_i = \{f(\phi_i(\mathbf{I})), g(\mathbf{X}_{i+1})\}, i \in \{1, 2\}, \qquad (1)$$

where $f$ is a $1 \times 1$ convolutional layer for dimensionality reduction, $g$ refers to upsampling with nearest neighbor interpolation. Since the highest level features have no top-down connection, we only conduct dimensionality reduction for $\phi_3(\mathbf{I})$:

$$\mathbf{X}_3 = f(\phi_3(\mathbf{I})). \qquad (2)$$

The channel size of $\mathbf{X}_i$ equal to $d, 2d, 3d$ for $i \in \{1, 2, 3\}$. The combined features $\mathbf{X}_i$ are used for attribute-specific localization.

## 3.2. Attribute Localization Module

As mentioned in Section 1, several existing methods attempt to extract local features through attribute-agnostic visual attention, predefined rigid parts or external part localization modules. However, these methods are not the optimal solution since they overlook the significance of attribute-specific localization. As shown in Figure 1 (c,d), attentive regions belong to different attributes are mixed together, which is inconsistent with the original intention that narrowing the attentive region for improving attribute recognition. We believe that attribute-specific localization is a better choice since it can disentangle the confused attention masks into several individual regions, where each region for a specific attribute. Moreover, the learned attribute-specific regions are more interpretable since we can observe the attribute-region correspondence intuitively. What we need is a mechanism that can learn an individual bounding box, representing the discriminative region, in feature maps for a given attribute. The well-known RoI pooling technique [4] is inappropriate since it requires region annotations, which are not available in pedestrian attribute datasets. Inspired by the recent success of *Spatial Transformer Network* (STN) [9], we propose a flexible *Attribute Localization Module* (ALM) to automatically discover the discriminative regions for each attribute in a weakly-supervised manner. The overview of the proposed ALM is illustrated in Figure 3.

As shown, each ALM contains a spatial transformer layer originates from STN. STN is a differentiable module which is capable of applying a spatial transformation to a feature map, *e.g.* cropping, translation, and scaling. In this paper, we adopt a simplified version of STN since we treat the attribute region as a simple bounding box, which can be realized through the following transformation:

$$\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = \begin{bmatrix} s_x & 0 & t_x \\ 0 & s_y & t_y \end{bmatrix} \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix}, \qquad (3)$$

where $s_x$, $s_y$ are scaling parameters, and $t_x$, $t_y$ are translation parameters, the expected bounding box can be obtained through these four parameters. $(x_i^s, y_i^s)$ and $(x_i^t, y_i^t)$ are the source coordinates and target coordinates of the $i$-th pixel. To some extent, this simplified spatial transformer can be viewed as a differentiable RoI pooling, which is end-to-end trainable without region annotations. To accelerate the convergence, we simply constrain $s_x, s_y$ to $(0, 1)$ and $t_x, t_y$ to $(-1, 1)$ by a *sigmoid* and *tanh* activation, respectively.

In addition, we also introduce a tiny channel-attention sub-network, as shown in Figure 3. As mentioned above, the ALM takes the features combined from adjacent levels as input, where both finer details and strong semantics take the same proportion (both have $d$ channels), which means they equally contribute to attribute localization. However, the expected proportion should vary from attribute to attribute. For example, more details should be paid when recognizing finer attributes. Therefore, we introduce this channel-attention sub-network, similar to SE-Net [7], to modulate the inter-channel dependencies.

Specifically, the input features $\mathbf{X}_i$ pass through a series

of linear and nonlinear layers, producing a weight vector for feature recalibration across channels. The reweighted features are obtained by channel-wise multiplying the weight vector with $\mathbf{X}_i$, and an extra residual link is applied to preserve the complementary information. Subsequently, a fully-connected layer is applied to estimate the transformation matrix, denoted as $\mathcal{R}$, and then the region-based features sampled by bilinear interpolation are used for attribute classification. We simply formulate the prediction belong to $m$-th attribute at $i$-th level as:

$$\hat{y}_i^m = ALM_i^m(\mathbf{X}_i). \tag{4}$$

### 3.3. Deep Supervision

As illustrated in Figure 2, four individual prediction vectors are obtained from three ALM groups and one global branch. We apply the *deep supervision* [12, 32] mechanism for training where the four individual predictions are directly supervised by ground-truth labels. During inference, multiple prediction vectors are aggregated through an effective voting scheme that producing the maximum responses across different feature levels. The intuition behind this design is that each ALM should directly take the feedback about whether the localized region is accurate. If we only preserve the supervision of the fused predictions (maximum or averaging), the gradients are not informative enough of how each level performs, such that some branches are trained insufficiently. The maximum voting scheme is applied to choose the best predictions from different levels with the most accurate attribute region.

Specifically, we adopt the weighted binary cross-entropy loss function [13] at each stage, formulated as follow:

$$
\begin{aligned}
\mathcal{L}_i(\hat{y}_i, y) = -\frac{1}{M} \sum_{m=1}^{M} \gamma^m ( y^m \log(\sigma(\hat{y}_i^m)) \\
+ (1 - y^m) \log(1 - \sigma(\hat{y}_i^m)) ),
\end{aligned} \tag{5}
$$

where $\gamma^m = e^{-a_m}$ is the loss weight for $m$-th attribute and $a_m$ is the prior class distribution of $m$-th attribute, $M$ is the number of attributes, $i$ represents the $i$-th branch, where $i \in \{1, 2, 3, 4\}$, and $\sigma$ refers to the *sigmoid* activation. The total training loss is calculated by summing over the four individual loss: $\mathcal{L} = \sum_{i=1}^{4} \mathcal{L}_i$.

## 4. Experiments

### 4.1. Datasets and Evaluation Metrics

The proposed method is evaluated on three publicly available pedestrian attribute datasets: (1) The **PETA** dataset [1] consists of 19,000 images with 61 binary attributes and 4 multi-class attributes. Following the previous works [1, 25], the whole dataset is randomly partitioned

into three subsets: 9,500 for training, 1,900 for verification and 7,600 for testing. We choose 35 attributes which the positive ratio is higher than $5\%$ for evaluation. (2) The **RAP** dataset [16] contains 41,585 images which are collected from 26 indoor surveillance cameras, where each image is annotated with 72 fine-grained attributes. Following the official protocol [16], we split the whole dataset into 33,268 training images and 8,317 test images. Only 51 binary attributes with the positive ratio higher than $1\%$ are selected for evaluation. (3) The **PA-100K** dataset [20] is to-date the largest dataset for pedestrian attribute recognition, which contains 100,000 pedestrian images in total collected from outdoor surveillance cameras. Each image is annotated with 26 commonly used attributes. According to the official setting [20], the whole dataset is randomly split into 80,000 training images, 10,000 validation images and 10,000 test images.

We adopt two types of metrics for evaluation [16]: (1) Label-based: we calculate the mean accuracy (**mA**) as the mean of positive accuracy and negative accuracy for each attribute. The mA criterion can be formulated as:

$$mA = \frac{1}{2N} \sum_{i=1}^{M} \left( \frac{TP_i}{P_i} + \frac{TN_i}{N_i} \right), \tag{6}$$

where $N$ is the number of examples and $M$ is the number of attributes; $P_i$ and $TP_i$ are the number of positive examples and correctly predicted positive examples of the $i$-th attribute respectively; $N_i$ and $TN_i$ are defined similarly. (2) Instance-based: we adopt four well-known criteria: **accuracy**, **precision**, **recall** and **F1 score**, details are omitted.

### 4.2. Effectiveness of Critical Components

As shown in Table 1, starting with the BN-Inception baseline, we gradually append each component and meanwhile compare it with several variants. (1) **Attribute Localization Module**: We first evaluate the contribution of the simplified ALM (without channel-attention sub-network) by embedding ALMs at the final layer (`incep_5b`). The increased mA and F1 scores demonstrate the effectiveness of attribute-specific localization. Based on this fact, we further embed multiple ALMs at different feature levels (`incep_3b,4d,5b`), and a greater improvement is achieved ($3.1\%$ and $1.3\%$ in mA and F1, respectively). Considering the model complexity, we limit the number of levels to three in our framework. (2) **Top-down Guidance**: Secondly, we evaluate the impact of the proposed feature pyramid architecture by comparing with three variants, which are different in how to combine features from different levels. The first one is implemented by element-wise adding the features from different levels, like the original FPN [18], but the performance decreases. The poor results suggest that some essential information may disappear if we disregard the feature mismatching problem. The

| Metric / Component | mA | F1 |
|---|---|---|
| Baseline | 75.76 | 78.20 |
| ALM at Single Level (5b) | 77.45 | 79.14 |
| **ALM** at Multiple Levels (3b,4d,5b) | 78.89 | 79.50 |
| Top-down (Addition) | 78.51 | 79.42 |
| Top-down (Concatenation) | 79.93 | 79.91 |
| **Top-down** (Channel Attention) | 80.61 | 79.98 |
| Deep Supervision (Averaging) | 80.70 | 80.04 |
| **Deep Supervision** (Maximum) (**Ours**) | **81.87** | **80.16** |
| **Ours** w/o ALMs | 78.91 | 79.55 |

Table 1. Performance comparisons on RAP dataset when gradually adding each proposed component to the baseline model (except the last row). Variants of the same component lie in the same group. **Bold** means the setting adopted in our final framework.
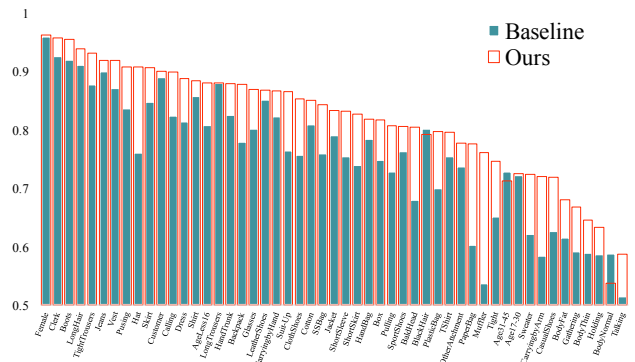


Figure 4. Attribute-wise mA comparison on RAP dataset between our proposed method and the baseline model. The bars are sorted in descending order according to the larger mA between the two models. We can observe significant improvements on some fine-grained attributes, *e.g. BaldHead*, *Hat* and *Muffler*.



Figure 5. Visualization of attribute localization results at different feature levels. Best viewed in color.

improved concatenation version achieves better results (improves $1.0\%$ in mA), which shows the success of high-level top-down guidance. Moreover, the introduced channel-attention sub-network further improves mA a lot to $80.61\%$ by modulating the inter-channel dependencies. (3) **Deep Supervision**: As mentioned in Section 3.3, the obtained gradients with only the supervision of fused predictions are not informative enough of how each level performs, while some branches are trained insufficiently. To address this problem, ALMs at different levels are trained with *deep supervision* mechanism. For inference, the experimental results suggest that element-wise maximum is a superior ensemble method than averaging since some weaker existences are ignored in averaging.

Removing all ALMs while keeping others unchanged results in a significant drop (last row in Table 1), which further confirmed the effectiveness of ALMs. Compared with the baseline, the final model achieves a remarkable performance, improving $6.1\%$ and $1.9\%$ in mA and F1 metrics, respectively. Figure 4 shows the attribute-wise mA comparison between the proposed method and baseline model on RAP dataset. As shown, the proposed method achieves significant improvement on a number of attributes, especially some fine-grained attributes, *e.g. BaldHead*$(23.1\%)$, *Hat*$(12.4\%)$ and *Muffler*$(13.5\%)$. The accurate recognition of these attributes shows the effectiveness of the proposed attribute-specific localization module.

### 4.3. Visualization of Attribute Localization

Through the above quantitative evaluation, we can observe significant improvements on some fine-grained attributes. In this subsection, we visualize the localized attribute regions from different feature levels for qualitative analysis. In our implementation, the attribute regions are located within the feature maps, while the correspondence between a feature map pixel and an image pixel is not unique. For a relatively coarse visualization, we simply map a feature-level pixel to the center of the receptive field on the input image, like SPPNet [5]. As shown in Figure 5, we display several examples belong to six different attributes, covering both abstract and concrete attributes. As we can see, the proposed ALMs can successfully localize these concrete attributes, *e.g. Backpack*, *PlasticBag*, and *Hat*, into the corresponded informative regions, despite the extreme occlusions (a, c) or pose variances (e). While recognizing the more abstract attributes *Clerk* and *BodyFat*, the ALMs tend to explore the larger regions, since they often require high-level semantics from the whole image. In addition, a failure case is also provided, as shown in Figure 5(d). The ALMs fail to localize the expected regions at two lower levels when recognizing *BaldHead*. We believe that this prob-

lem originates from the highly imbalanced data distribution, where only $0.4$ percent of images are annotated with *Bald-Head* in the RAP dataset. Although these localized attribute regions are relatively coarse, it is still acceptable for recognizing attributes because they indeed capture these most discriminative regions with large overlap.

### 4.4. Different Attribute-Specific Methods

The most significant contribution of this work is the idea of localizing an individual informative region for each attribute, which we called **attribute-specific** and was not well investigated in previous works. In this subsection, we conduct experiments to demonstrate the advantages of our proposed method by comparing with other attribute-specific localization methods, such as visual attention and predefined parts. Different from the attribute-agnostic attention masks and body-parts illustrated in Figure 1, we extend them to an attribute-specific version for comparison. Firstly, we replace the proposed ALM with a spatial attention module while keeping others unchanged for a fair comparison. In detail, we generate individual attention masks for each attribute through a global cross-channel averaging layer and a $3 \times 3$ convolutional layer, like HA-CNN [17]. For another comparison model, we divide the whole image into three rigid parts (head, torso, and legs) and extract part-based features with an RoI pooling layer, then manually define the attribute-part relations, *e.g.* recognizing *hat* only from the head part. More details about the compared methods are shown in Appendix B. Experimental results are listed in Table 2. As expected, the proposed method largely outperforms the other two methods (improving $5.3\%$ and $3.5\%$ in mA, respectively).

To better understanding the differences, we visualize these localization results in Figure 6. As we can see, the attribute regions generated by ALMs are the most accurate and discriminative one. Although the attention-based model achieves a not-bad result, the generated attention masks may attend to the irrelevant or biased regions. While recognizing *Box*, the attention masks fail to cover the expected regions, and we also observed that they tend to localize almost the same regions wherever the boxes are. By contrast, the proposed method can successfully handle the location uncertainties and pose variances. We provide more visualization results in Figure S4.

To some extent, the methods relying on attention masks and rigid parts are at two extremes. The former attempts to completely cover the informative pixels in a highly adaptive way, but mostly fails since we have only image-level annotations. The latter one just totally discards the adaptive factors, which are less robust to pose variances. Therefore, the proposed method attempts to achieve a balance between these two extremes, by constraining the attentional regions to several bounding boxes, which relatively coarse but more

| Metric / Method | mA | F1 |
|---|---|---|
| Rigid Part | 76.56 | 78.84 |
| Attention Mask | 78.35 | 79.51 |
| **Attribute Region** | **81.87** | **80.16** |

Table 2. Experimental results of different attribute-specific localization methods evaluated on RAP dataset.
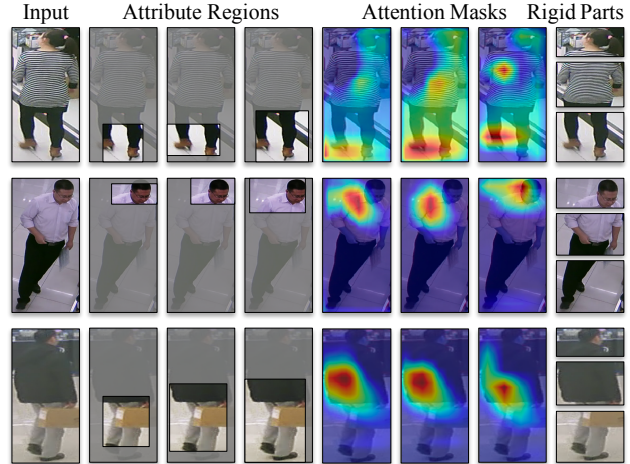


Figure 6. Case studies of different attribute-specific localization methods on three different attributes: ***Boots*** (Top), ***Glasses*** (Middle), and ***Box*** (Bottom). Different from Figure 1, the attention masks and body-parts are applied in an attribute-specific manner.

interpretable and controllable.

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

In this subsection, we compare the performance of our proposed method against several state-of-the-art methods. As mentioned in Section 2, we divide these methods into four categories: (1) Holistic methods including ACN [28] and DeepMar [13], which first take CNN to jointly learn multiple attributes. (2) Relation-based methods including JRL [29] and GRL [37], which both exploit the semantic relations by a CNN-RNN based model. (3) Attention-based methods including HP-Net [20] and DIAA [19] relying on multi-scale attention mechanism, and VeSPA [25] which perform view-specific attribute prediction through a coarse view predictor. (4) Part-based methods including recently proposed PGDM [15] and LG-Net [19], which relying on external pose estimation or region proposal module.

Table 3 and Table 4 show the comparison results on three different datasets. The results suggest that our proposed method achieves superior performances compared with existing works under both label-based and instance-based metrics on all three datasets. Compared with the previous methods relying on attribute-agnostic attention or extra part localization mechanism, the proposed method can achieve a significant improvement across all datasets, which

| Dataset | PETA | | | | | RAP | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Method \\ Metric | mA | Accu | Prec | Recall | F1 | mA | Accu | Prec | Recall | F1 | #P | GFLOPs |
| ACN [28] | 81.15 | 73.66 | 84.06 | 81.26 | 82.64 | 69.66 | 62.61 | 80.12 | 72.26 | 75.98 | - | - |
| DeepMar [13] | 82.89 | 75.07 | 83.68 | 83.14 | 83.41 | 73.79 | 62.02 | 74.92 | 76.21 | 75.56 | 58.5M | **0.72** |
| JRL [29] | 85.67 | - | 86.03 | 85.34 | 85.42 | 77.81 | - | 78.11 | 78.98 | 78.58 | - | - |
| JRL* [29] | 82.13 | - | 82.55 | 82.12 | 82.02 | 74.74 | - | 75.08 | 74.96 | 74.62 | - | - |
| GRL [37] | **86.70** | - | 84.34 | **88.82** | 86.51 | 81.20 | - | 77.70 | 80.90 | 79.29 | >50M | >10 |
| HP-Net [20] | 81.77 | 76.13 | 84.92 | 83.24 | 84.07 | 76.12 | 65.39 | 77.33 | 78.79 | 78.05 | - | - |
| VeSPA [25] | 83.45 | 77.73 | 86.18 | 84.81 | 85.49 | 77.70 | 67.35 | 79.51 | 79.67 | 79.59 | 17.0M | > 3 |
| DIAA [24] | 84.59 | 78.56 | 86.79 | 86.12 | 86.46 | - | - | - | - | - | - | - |
| PGDM [15] | 82.97 | 78.08 | **86.86** | 84.68 | 85.76 | 74.31 | 64.57 | 78.86 | 75.90 | 77.35 | 87.2M | ≈1 |
| LG-Net [19] | - | - | - | - | - | 78.68 | 68.00 | **80.36** | 79.82 | 80.09 | >20M | > 4 |
| BN-Inception | 82.66 | 77.73 | 86.68 | 84.20 | 85.57 | 75.76 | 65.57 | 78.92 | 77.49 | 78.20 | **10.3M** | 1.78 |
| **Ours** | 86.30 | **79.52** | 85.65 | 88.09 | **86.85** | 81.87 | 68.17 | 74.71 | **86.48** | **80.16** | 17.1M | 1.95 |

Table 3. Quantitative comparisons against previous methods on PETA and RAP datasets. We divide these methods into four groups: holistic methods, relation-based methods, attention-based methods, and part-based methods, from top to bottom. JRL* is the single model version of JRL. The precision and recall metrics are not so reliable in class-imbalanced datasets while the mA and F1 score are more convictive. Best results are in **bold**. For RAP dataset, we further provide comparisons on the number of parameters (#P) and complexity (GFLOPs).

| Dataset | PA-100K | | | | |
|---|---|---|---|---|---|
| Method | mA | Accu | Prec | Recall | F1 |
| DeepMar [13] | 72.70 | 70.39 | 82.24 | 80.42 | 81.32 |
| HP-Net [20] | 74.21 | 72.19 | 82.97 | 82.09 | 82.53 |
| PGDM [15] | 74.95 | 73.08 | 84.36 | 82.24 | 83.29 |
| VeSPA [25] | 76.32 | 73.00 | 84.99 | 81.49 | 83.20 |
| LG-Net [19] | 76.96 | 75.55 | **86.99** | 83.17 | 85.04 |
| BN-Inception | 77.47 | 75.05 | 86.61 | 85.34 | 85.97 |
| **Ours** | **80.68** | **77.08** | 84.21 | **88.84** | **86.46** |

Table 4. Quantitative comparisons on PA-100K dataset.

demonstrates the effectiveness of attribute-specific localization. Although a slightly lower mA score is achieved than the relation-based method GRL on PETA dataset, due to their stronger Inception-v3 backbone network (with twice as many parameters as ours), we can still outperform them on other metrics and datasets. On the more challenging dataset PA-100K, the proposed method largely outperforms all previous works, improving 3.7% and 1.4% in mA and F1, respectively, over the second best results. Notably, the proposed method surpasses the baseline model with a significant margin, especially on the label-based metric mA (3.6%, 6.1%, and 3.2% on three datasets, respectively). Note that the proposed method often achieve a lower precision but higher recall, while these two metrics are not so reliable, especially in class-imbalanced datasets. Moreover, the two metrics are inversely correlated, *i.e.*, increase in one metric always leads to decrease in another (*e.g.*, by modulating the class weights in the loss function). The mA and F1 metrics are more appropriate in measuring the performance of an attribute recognition model. Our method consistently achieves the best results in these two metrics.

We provide a comparison of the computational cost for different methods (rightmost columns in Table 3) on RAP dataset. For the number of parameters, theoretically, there are totally ($\frac{C^2}{8} + 4C$) trainable parameters in each ALM: $4C$ from the STN module, $\frac{C^2}{8}$ from the channel-attention module, where $C$ is the number of input channels. As shown, the proposed model has much fewer trainable parameters than previous models. In terms of model complexity, even with 51 attributes, the proposed model is still light-weight as only 0.17 GFLOPs are added to the backbone network. The reason is that ALM contains only FC-layers (or $1 \times 1$ Conv), which involves much fewer FLOPs than $3 \times 3$ Conv-layers. In general, the entire model is much more efficient than previous models.

## 5. Conclusion

We propose an end-to-end framework for pedestrian attribute recognition, which can automatically localize the attribute-specific regions at multiple feature levels. Moreover, we apply a feature pyramid architecture to enhance the attribute localization and region-based feature learning in a mutually reinforcing manner. Experimental results on PETA, RAP, and PA-100K datasets show that the proposed method can significantly outperform most of the existing methods. The extensive analysis suggests that the proposed method can successfully localize the most informative region for each attribute in a weakly-supervised manner.

# References

[1] Yubin Deng, Ping Luo, Chen Change Loy, and Xiaoou Tang. Pedestrian attribute recognition at far distance. In *Proceedings of the 22nd ACM International Conference on Multimedia*, pages 789–792, 2014.

[2] Rogerio Feris, Russel Bobbitt, Lisa Brown, and Sharath Pankanti. Attribute-based people search: Lessons learnt from a practical surveillance system. In *Proceedings of International Conference on Multimedia Retrieval*, pages 153–160, 2014.

[3] Jianlong Fu, Heliang Zheng, and Tao Mei. Look closer to see better: Recurrent attention convolutional neural network for fine-grained image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4438–4446, 2017.

[4] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1440–1448, 2015.

[5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(9):1904–1916, 2015.

[6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.

[7] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7132–7141, 2018.

[8] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, pages 448–456, 2015.

[9] Max Jaderberg, Karen Simonyan, Andrew Zisserman, and Koray Kavukcuoglu. Spatial transformer networks. In *Advances in Neural Information Processing Systems*, pages 2017–2025, 2015.

[10] Neeraj Kumar, Alexander C Berg, Peter N Belhumeur, and Shree K Nayar. Attribute and simile classifiers for face verification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 365–372, 2009.

[11] Ryan Layne, Timothy M Hospedales, Shaogang Gong, and Q Mary. Person re-identification by attributes. In *Proceedings of the British Machine Vision Conference*, 2012.

[12] Chen-Yu Lee, Saining Xie, Patrick Gallagher, Zhengyou Zhang, and Zhuowen Tu. Deeply-supervised nets. In *Artificial Intelligence and Statistics*, pages 562–570, 2015.

[13] Dangwei Li, Xiaotang Chen, and Kaiqi Huang. Multi-attribute learning for pedestrian attribute recognition in surveillance scenarios. In *Proceedings of the IAPR Asian Conference on Pattern Recognition*, pages 111–115, 2015.

[14] Dangwei Li, Xiaotang Chen, Zhang Zhang, and Kaiqi Huang. Learning deep context-aware features over body and latent parts for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 384–393, 2017.

[15] Dangwei Li, Xiaotang Chen, Zhang Zhang, and Kaiqi Huang. Pose guided deep model for pedestrian attribute recognition in surveillance scenarios. In *Proceedings of the IEEE International Conference on Multimedia and Expo*, 2018.

[16] Dangwei Li, Zhang Zhang, Xiaotang Chen, Haibin Ling, and Kaiqi Huang. A richly annotated dataset for pedestrian attribute recognition. *arXiv preprint arXiv:1603.07054*, 2016.

[17] Wei Li, Xiatian Zhu, and Shaogang Gong. Harmonious attention network for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2285–2294, 2018.

[18] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2117–2125, 2017.

[19] Pengze Liu, Xihui Liu, Junjie Yan, and Jing Shao. Localization guided learning for pedestrian attribute recognition. In *Proceedings of the British Machine Vision Conference*, 2018.

[20] Xihui Liu, Haiyu Zhao, Maoqing Tian, Lu Sheng, Jing Shao, Shuai Yi, Junjie Yan, and Xiaogang Wang. Hydraplus-net: Attentive deep features for pedestrian analysis. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 350–359, 2017.

[21] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *Proceedings of the European Conference on Computer Vision*, pages 483–499, 2016.

[22] Peixi Peng, Yonghong Tian, Tao Xiang, Yaowei Wang, and Tiejun Huang. Joint learning of semantic and latent attributes. In *Proceedings of the European Conference on Computer Vision*, 2016.

[23] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 234–241, 2015.

[24] Nikolaos Sarafianos, Xiang Xu, and Ioannis A Kakadiaris. Deep imbalanced attribute classification using visual attention aggregation. In *Proceedings of the European Conference on Computer Vision*, pages 680–697, 2018.

[25] M Saquib Sarfraz, Arne Schumann, Yan Wang, and Rainer Stiefelhagen. Deep view-sensitive pedestrian attribute inference in an end-to-end model. In *Proceedings of the British Machine Vision Conference*, 2017.

[26] Abhinav Shrivastava, Rahul Sukthankar, Jitendra Malik, and Abhinav Gupta. Beyond skip connections: Top-down modulation for object detection. *arXiv preprint arXiv:1612.06851*, 2016.

[27] Behjat Siddiquie, Rogerio S Feris, and Larry S Davis. Image ranking and retrieval based on multi-attribute queries. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 801–808, 2011.

[28] Patrick Sudowe, Hannah Spitzer, and Bastian Leibe. Person attribute recognition with a jointly-trained holistic cnn model. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 87–95, 2015.

[29] Jingya Wang, Xiatian Zhu, Shaogang Gong, and Wei Li. Attribute recognition by joint recurrent learning of context and correlation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 531–540, 2017.

[30] Jingya Wang, Xiatian Zhu, Shaogang Gong, and Wei Li. Transferable joint attribute-identity deep learning for unsupervised person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2275–2284, 2018.

[31] Xiao Wang, Shaofei Zheng, Rui Yang, Bin Luo, and Jin Tang. Pedestrian attribute recognition: A survey. *arXiv preprint arXiv:1901.07474*, 2019.

[32] Yan Wang, Lequn Wang, Yurong You, Xu Zou, Vincent Chen, Serena Li, Gao Huang, Bharath Hariharan, and Kilian Q Weinberger. Resource aware person re-identification across multiple resolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8042–8051, 2018.

[33] Zhouxia Wang, Tianshui Chen, Guanbin Li, Ruijia Xu, and Liang Lin. Multi-label image recognition by recurrently discovering attentional regions. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 464–472, 2017.

[34] Luwei Yang, Ligen Zhu, Yichen Wei, Shuang Liang, and Ping Tan. Attribute recognition from adaptive parts. In *Proceedings of the British Machine Vision Conference*, 2016.

[35] Ning Zhang, Manohar Paluri, Marc'Aurelio Ranzato, Trevor Darrell, and Lubomir Bourdev. Panda: Pose aligned networks for deep attribute modeling. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1637–1644, 2014.

[36] Haiyu Zhao, Maoqing Tian, Shuyang Sun, Jing Shao, Junjie Yan, Shuai Yi, Xiaogang Wang, and Xiaoou Tang. Spindle net: Person re-identification with human body region guided feature decomposition and fusion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 907–915, 2017.

[37] Xin Zhao, Liufang Sang, Guiguang Ding, Yuchen Guo, and Xiaoming Jin. Grouping attribute recognition for pedestrian with joint recurrent learning. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence*, pages 3177–3183, 2018.

[38] Feng Zhu, Hongsheng Li, Wanli Ouyang, Nenghai Yu, and Xiaogang Wang. Learning spatial regularization with image-level supervisions for multi-label image classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5513–5522, 2017.

[39] Jianqing Zhu, Shengcai Liao, Zhen Lei, Dong Yi, and Stan Li. Pedestrian attribute classification in surveillance: Database and evaluation. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 331–338, 2013.

[40] Jianqing Zhu, Shengcai Liao, Dong Yi, Zhen Lei, and Stan Z Li. Multi-label cnn based pedestrian attribute learning for soft biometrics. In *Proceedings of the International Conference on Biometrics*, pages 535–540, 2015.

[41] Lei Zhu, Zijun Deng, Xiaowei Hu, Chi-Wing Fu, Xuemiao Xu, Jing Qin, and Pheng-Ann Heng. Bidirectional feature pyramid network with recurrent attention residual modules for shadow detection. In *Proceedings of the European Conference on Computer Vision*, pages 121–136, 2018.

[42] C. Lawrence Zitnick and Piotr Dollár. Edge boxes: Locating object proposals from edges. In *Proceedings of the European Conference on Computer Vision*, 2014.

# Appendix

## A. Implementation Details

We adopt the BN-Inception model pretrained from ImageNet as the backbone network. The proposed framework is implemented with PyTorch framework and trained end-to-end with only image-level annotations. We adopt Adam optimizer since it converges faster than SGD in our experiments with momentum set to $0.9$ and a weight decay equals to $0.0005$. The initial learning rate equals to $0.0001$ and the batch size is set to $32$. For RAP and PA-100K dataset, we train the model for 30 epochs and the learning rate decays by $0.1$ every 10 epochs. For the smaller PETA dataset, we double the training epochs. For data preprocessing, we resize the input pedestrian images to $256 \times 128$ and apply random horizontal mirroring and data shuffling for data augmentation.

## B. Different Attribute-Specific Methods

In Section 4.4, we compare the proposed method against the other two attribute-specific localization methods, including visual attention and rigid parts. Different from most existing attribute-agnostic attention-based and part-based methods, we build two attribute-specific models based on these ideas for comparison. Here we show the details of the compared models.

**Attention Masks Model**. We replace the proposed ALM with a spatial attention module while keeping others unchanged for fair comparison. The spatial attention module is implemented by a tiny 3-layers sub-network, as shown in Figure S2, which is inspired by HA-CNN [17]. The input features $\mathbf{X}_i \in \mathbb{R}^{H \times W \times C}$ at the $i$-th level (a certain layer in the backbone network, totally three levels) are first fed into a cross-channel averaging layer. A $3 \times 3$ Conv-BatchNorm-ReLU block is followed to generate the expected attention mask $\mathbf{S}_i^m \in \mathbb{R}^{H \times W \times 1}$, which is used for localizing the $m$-th attribute at the $i$-th level. All channels share the identical spatial attention mask. Subsequently, the attentive features are obtained by channel-wise multiplying the attention mask with the input features, and the corresponding prediction is calculated as follows:

$$\hat{y}_i^m = f(\mathbf{S}_i^m \cdot \mathbf{X}_i), \tag{S1}$$

where $f$ denotes a fully-connected layer. Each spatial attention module only serves one attribute at a singe level, the same as Figure 3.

| Region | Attributes |
|--------|-----------|
| Head | BaldHead, LongHair, BlackHair, Hat, Glasses, Muffler, Calling |
| Torso | Shirt, Sweater, Vest, TShirt, Cotton, Jacket, Suit-Up, Tight, ShortSleeve, LongTrousers, Skirt, ShortSkirt, Dress, Jeans, TightTrousers, CarryingbyArm, CarryingbyHand |
| Legs | LeatherShoes, SportShoes, Boots, ClothShoes, CasualShoes |
| Whole | Female, AgeLess16, Age17-30, Age31-45, BodyFat, BodyNormal, BodyThin, Customer, Clerk, Backpack, SSBag, HandBag, Box, PlasticBag, PaperBag, HandTrunk, OtherAttchment, Talking, Gathering, Holding, Pushing, Pulling |

Table S1. Attribute-region correspondence in RAP dataset.

**Rigid Parts Model**. For attribute-specific part-based model, we replace ALM with a body-parts guided module, as shown in Figure S3. The key idea is to associate each attribute with a predefined body region, including *head, torso, legs*, and the whole image, *e.g.*, the *LongHair* attribute is associated with the head part. Since the body-part annotations are unavailable on most pedestrian attribute datasets, we adopt an external pose estimation model to localize the body parts, which is inspired by SpindleNet [36]. Specifically, we localize 14 human body keypoints for each pedestrian image using a pretrained pose estimation model [36]. The pedestrian image is then divided into three body-part regions based on these keypoints, as shown in Figure S1. In the body-parts guided module (Figure S3), the body-part-based local features are extracted from the input features $\mathbf{X}_i$ through an RoI pooling layer [4]. For attribute prediction, the most relevant features are selected according to the attribute-region correspondence, as listed in Table S1, *e.g.* recognizing *hat* using features only from the *head* part.
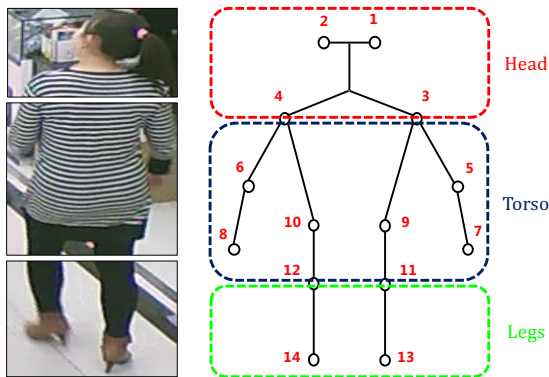


Figure S1. Illustration of body-parts generation. We divide a pedestrian image into three body-part regions (*head, torso, and legs*) based on 14 human body keypoints.

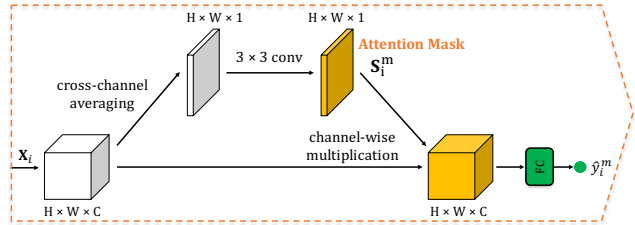We provide more localization results belong to different attributes, as shown in Figure S4.



Figure S2. Details of the spatial attention module for one attribute at a singe level. The expected attention mask follows a cross-channel averaging layer and a $3 \times 3$ Conv-BN-ReLU block.
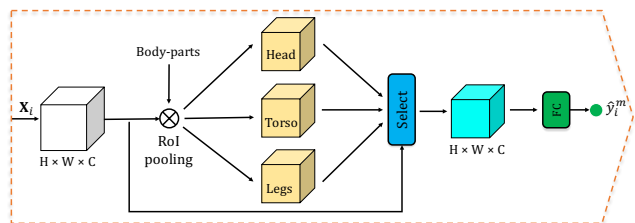


Figure S3. Details of the body-parts guided module for one attribute at a singe level. The three body-part regions are calculated based on several human body keypoints predicted by a pretrained pose estimation model. The local features belonging to different body-parts are extracted by an RoI pooling layer. The most relevant features are selected for attribute classification according to the predefined attribute-region correspondence (Table S1).
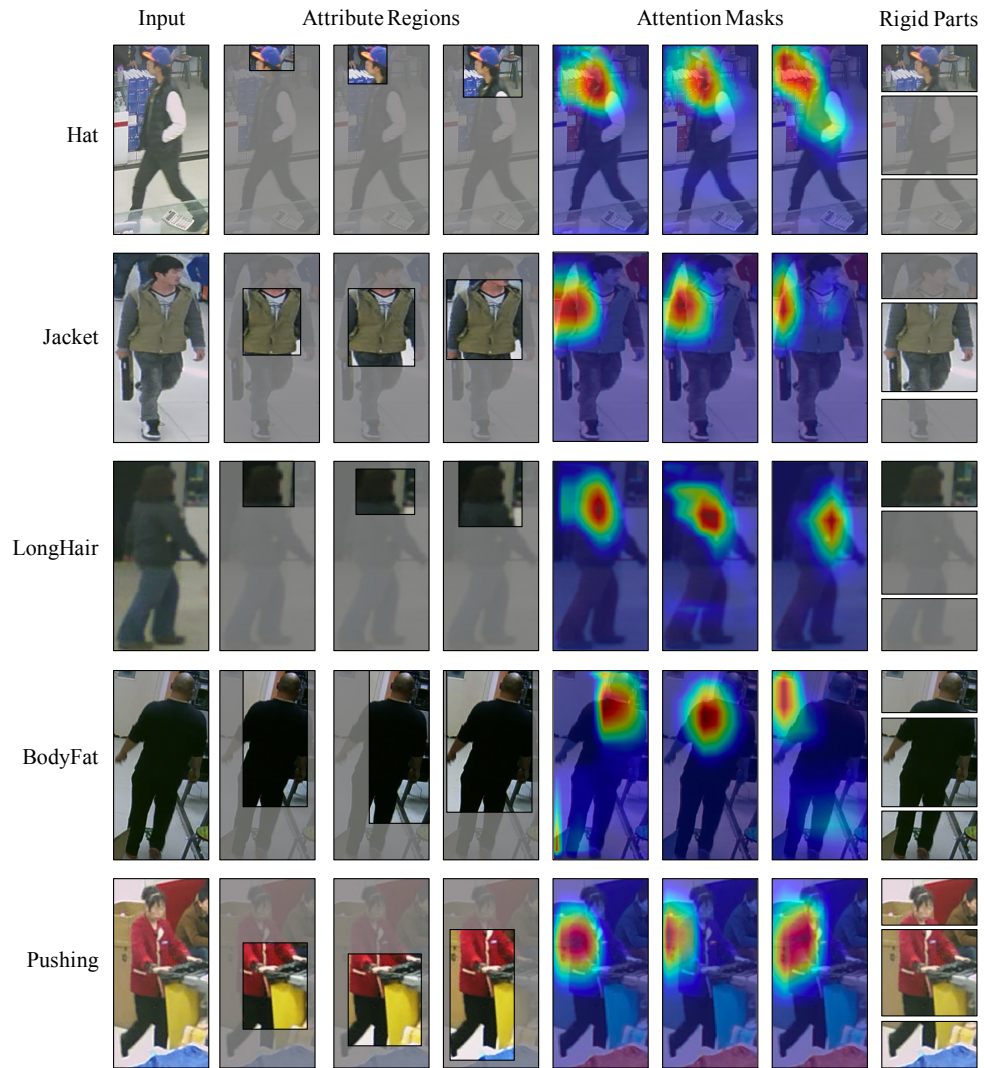
Figure S4. Case studies of different attribute-specific localization methods for five different attributes.