# HydraPlus-Net: Attentive Deep Features for Pedestrian Analysis

Xihui Liu[1,2*],   Haiyu Zhao[2*],   Maoqing Tian[2],   Lu Sheng[1]

Jing Shao[2†],   Shuai Yi[2],   Junjie Yan[2],   Xiaogang Wang[1]

[1]The Chinese University of Hong Kong    [2]SenseTime Group Limited

shaojing@sensetime.com

## Abstract

*Pedestrian analysis plays a vital role in intelligent video surveillance and is a key component for security-centric computer vision systems. Despite that the convolutional neural networks are remarkable in learning discriminative features from images, the learning of comprehensive features of pedestrians for fine-grained tasks remains an open problem. In this study, we propose a new attention-based deep neural network, named as HydraPlus-Net (HP-net), that multi-directionally feeds the multi-level attention maps to different feature layers. The attentive deep features learned from the proposed HP-net bring unique advantages: (1) the model is capable of capturing multiple attentions from low-level to semantic-level, and (2) it explores the multi-scale selectiveness of attentive features to enrich the final feature representations for a pedestrian image. We demonstrate the effectiveness and generality of the proposed HP-net for pedestrian analysis on two tasks, i.e. pedestrian attribute recognition and person re-identification. Intensive experimental results have been provided to prove that the HP-net outperforms the state-of-the-art methods on various datasets.[1]*

## 1. Introduction

Pedestrian analysis is a long-lasting research topic because of the continuing demands for intelligent video surveillance and psychological social behavior researches. Particularly, with the explosion of researches about the deep convolutional neural networks in recent computer vision community, a variety of applications categorized as the pedestrian analysis, *e.g.* pedestrian attribute recognition, person re-identification and *etc.*, have received remarkable improvements and presented potentialities for practical us-

---

[*]X. Liu and H. Zhao share equal contribution.

[†]J. Shao is the corresponding author.
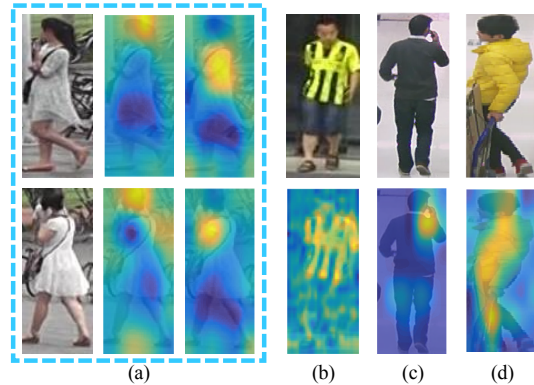
[1]https://github.com/xh-liu/HydraPlus-Net



Figure 1. Pedestrian analysis needs a comprehensive feature representation from multi-levels and scales. (a) *Semantic-level*: attending features around local regions facilitates distinguishing persons that own similar appearances at a glance, such as "long hair" vs. "short hair" and "long-sleeves" vs. "short-sleeves". (b) *Low-level*: some patterns like "clothing stride" can be well captured by low-level features rather than those in high-level. (c-d) *Scales*: multi-scale attentive features benefit describing person's characteristics, where the small-scale attention map in (c) corresponds to the "phone" and a large-scale one in (d) for a global understanding.

age in modern surveillance system. However, the learning of feature representation for pedestrian images, as the backbone for all those applications, still confronts critical challenges and needs profound studies.

At first, most traditional deep architectures have not extracted the detailed and localized features complementary to the high-level global features, which are especially effective for fine-grained tasks in pedestrian analysis. For example, it is difficult to distinguish two instances if no semantic features are extracted around hair and shoulders, as shown in Fig. 1(a). Also in Fig. 1(c), the effective features should be located within a small-scale head-shoulder region if we want to detect the attribute "calling". However, existing arts merely extract global features [13, 24, 30] and are hardly effective to location-aware semantic pattern extraction. Furthermore, it is well-known that multi-level features

aid diverse vision tasks [21, 6]. Similar phenomenon has also happened in the pedestrian analysis, such as the pattern "clothing stride" shown in Fig. 1(b) should be inferred from low-level features, while the attribute "gender" in Fig. 1(d) is judged by semantic understanding of the whole pedestrian image. Unlike previous approaches that mainly generate the global feature representations, the proposed feature representation encodes multiple levels of feature patterns as well as a mixture of global and local information, and thus it owns a potential capability for multi-level pedestrian attribute recognition and person re-identification.

Facing the drawbacks of recent methods for pedestrian analysis, we try to tackle the general feature learning paradigm for pedestrian analysis by a multi-directional network, called HydraPlus-Net, which is proposed to better exploit the global and local contents with multi-level feature fusion of a single pedestrian image. Specifically, we propose a multi-directional attention (MDA) module that aggregates multiple feature layers within the attentive regions extracted from multiple layers in the network. Since the attention maps are extracted from different semantic layers, they naturally abstract different levels of visual patterns of the same pedestrian image. Moreover, filtering multiple levels of features by the same attention map results in an effective fusion of multi-level features from a certain local attention distribution. After applying the MDA to different layers of the network, the multi-level attentive features are fused together to form the final feature representation.

The proposed framework is evaluated on two representatives among the pedestrian analysis tasks, *i.e.* pedestrian attribute recognition and person re-identification (ReID), in which attribute recognition focuses on assigning a set of attribute labels to each pedestrian image while ReID aims to associate the images of one person across multiple cameras and/or temporal shots. Although pedestrian attribute recognition and ReID pay attention to different aspects of the input pedestrian image, these two tasks can be solved by learning a similar feature representation, since they are inherently correlated with similar semantic features and the success of one task will improve the performance of the other. Compared with existing approaches, our framework achieves the state-of-the-art performance on most datasets.

The contributions of this work are three-fold:

(1) A HydraPlus Network (HP-net) is proposed with the novel multi-directional attention modules to train multi-level and multi-scale attention-strengthened features for fine-grained tasks of pedestrian analysis.

(2) The HP-net is comprehensively evaluated on pedestrian attribute recognition and person re-identification. State-of-the-art performances have been achieved with significant improvements against the prior methods.

(3) A new large-scale pedestrian attribute dataset (PA-100K dataset) is collected with the most diverse scenes and
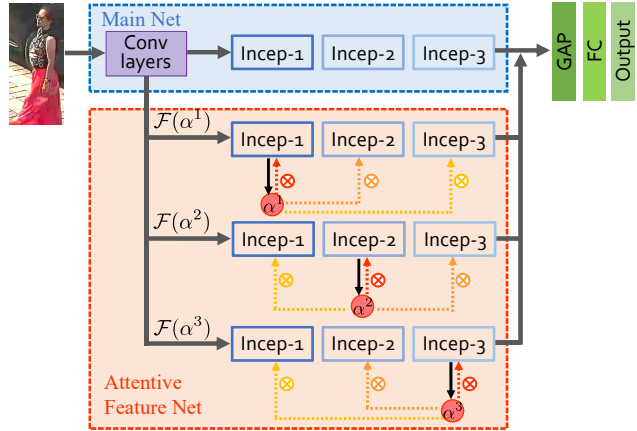


Figure 2. A deep HP-Net with a Main Net (M-net) and an Attentive Feature Net (AF-net). The AF-net comprises three multi-directional attention (MDA) modules (*i.e.* $\mathcal{F}(\alpha^i), i \in \Omega$). Each MDA module includes two components: (1) attention map generation with black solid lines, and (2) attentive features by masking the attention map to different levels of features in hot dash lines. A global average pooling and one fully-connected layer are applied to the concatenated features obtained from the M-net and AF-net.

the largest number of samples and instances up-to-date. The PA-100K dataset is more informative than the previous collections and helpful for various pedestrian analysis tasks.

## 2. Related Works

**Attention models** In computer vision, attention models have been used in tasks such as image caption generation [34], visual question answering [18, 33] and object detection [2]. Mnih *et al*. [20] and Xiao *et al*. [32] explored hard attention, in which the network attends to a certain region of the image or feature map. Compared to non-differentiable hard attention trained by reinforce algorithms [28], soft attention which weights the feature maps is differentiable and can be trained by back propagation. Chen *et al*. [4] introduced an attention to the multi-scale features, and Zagoruyko *et al*. [35] exploited attention in knowledge transfer. In this work, we design a multi-directional attention network for better pedestrian feature representation and apply it to both pedestrian attribute recognition and re-identification tasks. To the best of our knowledge, this is the first work to adopt attention idea in the aforementioned two tasks.

**Pedestrian attribute recognition** Pedestrian attribute has been an important research topic recently, due to its prospective application in video surveillance systems. Convolutional neural networks have achieved great success in pedestrian attribute recognition. Sudowe *et al*. [24] and Li *et al*. [13] proposed that jointly training multiple attributes can improve the performance of attribute recognition. Previous work also investigated the effectiveness of
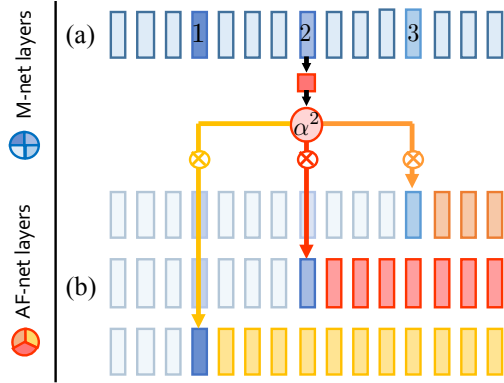
Figure 3. An example of the multi-directional attention (MDA) module $\mathcal{F}(\alpha^2)$. The M-net in (a) is presented with simplified layers (indexed by $i \in \Omega$) representing the output layers of three `inception` blocks. The attention map $\alpha^2$ is generated from block 2 and multi-directionally masks three adjacent blocks.
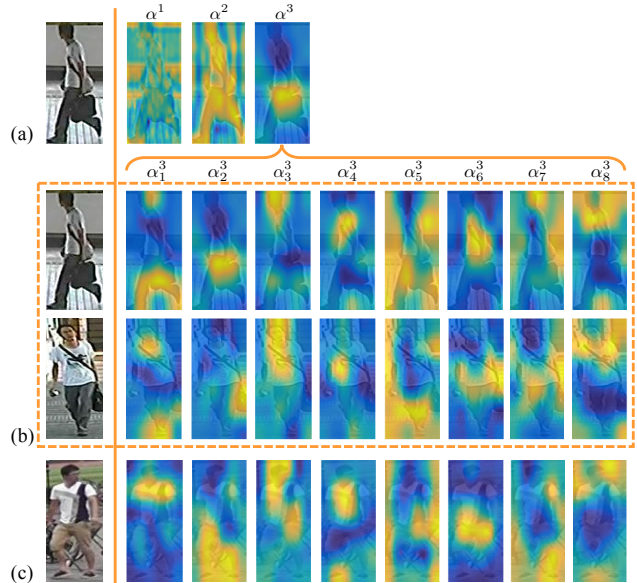


Figure 4. The diversity and semantic selectiveness of attention maps. (a) The attention maps generated from three adjacent blocks respond to visual patterns in different scales and levels, in which $\alpha^3$ owns the power to highlight semantic patterns in object level. (b-c) Different channels in attention maps capture different visual patterns related to body parts, salient objects and background.

utilizing pose and body part information in attribute recognition. Zhang *et al*. [37] proposed a pose aligned network to capture the pose-normalized appearance differences. Different from previous works, we propose an attention structure which can attend to important areas and align body parts without prior knowledge on body parts or poselets.

**Person re-identification** Feature extraction and metric learning [12, 17] are two main components for person re-identification. The success of deep learning in image classification inspired lots of studies on person ReID [5, 16, 30, 29, 26, 23, 25, 15, 31]. The filter pairing neural network (FPNN) proposed by Li *et al*. [16] jointly handles misalignment, transforms, occlusions and background clutters. Cheng *et al*. [5] presented a multi-channel parts-based CNN to learn body features from the input image. In this paper, we mainly target on feature extraction and cosine distance is directly adopted for metric learning. Moreover, attention masks are utilized in our pipeline to locate discriminative regions which can better describe each individual.

## 3. HydraPlus-Net Architecture

The design of the **HydraPlus network**[2] (HP-net) is motivated by the necessity to extract multi-scale features from multiple levels, so as not only to capture both global and local contents of the input image but also assemble its features with different levels of semantics. As shown in Fig. 2, the HP-net consists of two parts, one is the *Main Net* (M-net) that is a plain CNN architecture, the other is the *Attentive Feature Net* (AF-net) including multiple branches of multi-directional attention (MDA) modules applied to different semantic feature levels. The AF-net shares the same basic convolution architectures as the M-net except the added

---

[2]Hydra is a water monster with nine heads. In this work, the network consists of a 9-branch in AF-net (3 MDA modules and 3 attention sub-branches for each MDA), plus an M-net, so it is called HydraPlus Net.

MDA modules. Their outputs are concatenated and then fused by global average pooling (GAP) and fully connected (FC) layers. The final output can be projected as the attribute logits for attribute recognition or feature vectors for re-identification. In principle, any kind of CNN structure can be applied to construct the HP-net. But in our implementation, we design a new end-to-end model based on `inception_v2` architecture [10] because of its excellent performance in general image-related recognition tasks. As sketched in Fig. 2, each network of the proposed framework contains several low-level convolutional layers and is followed by three `inception` blocks. This model seems simple but is non trivial as it achieves all required abilities and brings them together to boost the recognition capability.

### 3.1. Attentive Feature Network

The Attentive Feature Network (AF-net) in Fig. 2 comprises three branches of sub-networks augmented by the multi-directional attention (MDA) modules, namely $\mathcal{F}(\alpha^i), i \in \Omega = \{1, 2, 3\}$, where $\alpha^i$ are the attention maps generated from the output features of the `inception` block $i$ marked by black solid lines, and are applied to the output of the $k^{\text{th}}$ block ($k \in \Omega = \{1, 2, 3\}$) in hot dash lines. For each MDA module, there is one link of attention generation and three links for attentive feature construction. Different MDA modules have their attention maps generated from different inception blocks and then been multiplied to feature maps of different levels to produce multi-level at-

tentive features. An example of a MDA module $\mathcal{F}(\alpha^2)$ is shown in Fig. 3. The main stream network of each AF-net branch is initialized exactly as the M-net, and thus the attention maps approximately distill similar features as what the M-net extracts.

It is well known that the attention maps learned from different blocks vary in scale and detailed structure. For example, the attention maps from higher blocks (*e.g.* $\alpha^3$) tend to be coarser but usually figure out the semantic regions like $\alpha^3$ highlights the handbag in Fig. 4(a). But those from lower blocks (*e.g.* $\alpha^1$) often respond to local feature patterns and can catch detailed local information like edges and textures, just as the examples visualized in Fig. 4(a). Therefore, if fusing the multi-level attentive features by MDA modules, we enable the output features to gather information across different levels of semantics, thus offering more selective representations. Moreover, the MDA module also differs from the traditional attention-based models [21, 34] that push the attention map back to the same block, and it extends this mechanism by applying the attention maps to *adjacent* blocks, as shown in lines with varying hot colors in Fig. 3. Applying one single attention map to multiple blocks naturally let the fused features encode multi-level information within the same spatial distribution which is illustrated in Section 4.2.

More specifically, for a given `inception` block $i$, its output feature map is denoted as $\mathbf{F}^i \in \mathbb{R}^{C \times H \times W}$ with the width $W$, height $H$ and $C$ channels. The attention map $\alpha^i$ is generated from $\mathbf{F}^i$ by a $1 \times 1$ `conv` layer with `BN` and `ReLU` activation function afterwards, noted as

$$\alpha^i = g_{\text{att}}(\mathbf{F}^i; \boldsymbol{\theta}^i_{\text{att}}) \in \mathbb{R}^{L \times H \times W}, \qquad (1)$$

where $L$ means the channels of the attention map. In this paper, we fix $L = 8$ for both tasks. And the attentive feature map to the `inception` block $k$ is an element-wise multiplication

$$\tilde{\mathbf{F}}^{i,k}_l = \alpha^i_l \circ \mathbf{F}^k, \ \ l \in \{1, \ldots, L\}. \qquad (2)$$

Each attentive feature map $\tilde{\mathbf{F}}^{i,k}_l$ is then passed through the following blocks thereafter, and at the end of MDA module we concatenate the $L$ attentive feature maps as the final feature representation. We visualized the detailed structure of an MDA module $\mathcal{F}(\alpha^2)$ in Fig. 3. $\alpha^2$ is generated from the `inception` block 2 and then applied to feature maps indexed by $k \in \Omega = \{1, 2, 3\}$, as shown in Fig. 3(b). Note that we prefer the `ReLU` activation function rather than the `sigmoid` function to constrain the attention maps so that the attentive regions receive more weights, and the contrast of the attention map is enlarged. More examples and analyses are shown in Sec. 4 to illustrate the MDA's effectiveness.

### 3.2. HP-Net Stage-wise Training

We train the HP-net in a stage-wise fashion. Initially, a plain M-net is trained to learn the fundamental pedes-
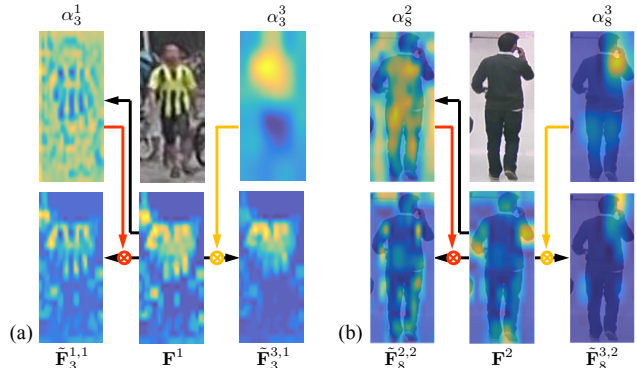


Figure 5. Examples of multi-directional attentive features. (a) The identification of low-level attributes like "upper-clothing pattern" requires the low-level attention connections, for example, applying $\alpha^1_3$ to extract $\tilde{\mathbf{F}}^{1,1}_3$ indicating textures onto the T-shirt. (b) But the semantic or object-level attributes like "phone" require high-level attention connections such as applying $\alpha^3_8$ to extract $\tilde{\mathbf{F}}^{3,2}_8$ for the detection of the phone near the ear.

trian features. Then the M-net is duplicated three times to construct the AF-net with adjacent MDA modules, each of which following the framework shown in Fig. 3. Since each MDA module consists of three branches where the attention map masks adjacent `inception` blocks, thus in each branch we only fine-tune the blocks after the attention-operated block. After separately fine-tuning three MDA modules in AF-net, we fix both the M-net and AF-net and train the remaining `GAP` and `FC` layers. The output layer to minimize losses defined by different tasks, in which the cross-entropy loss $\mathcal{L}_{\text{att}}$ is applied for pedestrian attribute recognition, and softmax loss for person re-identification.

## 4. Ablation Study On Attentive Deep Features

The advantages of HP-net are its capability of learning both multi-level attentions and multi-scale attentive features for a comprehensive feature representation of a pedestrian image. To better understand these advantages, we analyze the effectiveness of each component in the network with qualitative visualization and quantitative comparisons.

### 4.1. Multi-level Attention Maps

**The level of attention maps.** The compared exemplars of attention maps from three layers (*i.e.* the outputs of the `inception` blocks $i \in \Omega = \{1, 2, 3\}$) are shown in Fig. 4(a). We observe that the attention map from earlier layer $i = 1$ prefers grasping low-level patterns like edges or textures, while those from higher layers $i = 2$ or $3$ are more likely to capture semantic visual patterns corresponding to a specific object (*e.g.* handbag) or human identity.

**The quantity of attention maps.** Most previous studies [34, 21] merely demonstrated the effectiveness of an attention-based model with a limited number of channels

Figure 6. Results of discarding partial attention modules or connections compared with that of the complete network fed with all MDA modules on VIPeR dataset. The $3 \times 3$ boxes in (a) indicates the indices of different attention maps and their mask directions. The hollow white in each box means the corresponding attentions or directional links have been cut down. Bars are plot by the Top-1 accuracy. (b) and (c) present the qualitative results by the complete network compared with two kinds of partial networks in (a). For a query image shown in the middle, Top-5 results are shown aside with the correct marked by green and the false alarm are red. Best viewed in color.

(*i.e.*, $L = 1$ or 2). In this study, we explore the potential performance of an attention model with increasing channels in both diversity and consistency.

*1) Attention Diversity.* Fig. 4(b) shows two images of one single pedestrian captured by two cameras, alongside with $L = 8$ attention channels of $\alpha^3$ are presented. From the raw image, it is hard to distinguish these images due to the large intra-class variations from cluttered background, varying illumination, viewpoint changes and *etc*. Nevertheless, benefited from the discriminative localization ability of multiple attention channels from one level, the entire features can be captured separately with respect to different attentive areas. Compared to a single attention channel, the diversity of multiple attention channels enriches the feature representations and improves the chance to accurately analyze both the attributes and identity of one pedestrian.

*2) Attention Consistency.* We also observe that one attention map generated upon different input samples might be similarly distributed in spatial domain since they highlight the same semantic parts of a pedestrian. Notwithstanding different pedestrians, shown in Fig. 4(b-c), their attention channels $\alpha_3^3$ capture the head-shoulder regions and the channels $\alpha_5^3$ infer the background area. Since the consistent attention maps are usually linked to salient objects, the selectiveness of these attention maps is thus essential on identifying the pedestrian.

### 4.2. Multi-Directional Attentive Features

Apart from the benefits of the multi-level attention maps, the effectiveness of the proposed method also lies on the novel transition scheme. For instance, the pedestrian in Fig. 5(b) holds a phone near the right ear that cannot be directly captured neither by the feature map $\mathbf{F}^2$ in a lower layer $i = 2$, nor by the naïve attentive feature maps $\tilde{\mathbf{F}}_8^{2,2}$. Surprisingly, with the help of a higher level attention map $\alpha_8^3$, the attentive feature map $\tilde{\mathbf{F}}_8^{2,3}$ can precisely attend the

region around the phone. On the other hand, the high-level attention map $\alpha_3^3$ might not be able to capture lower-level visual patterns related to attributes like "upper-clothing pattern". For example, the attention map $\alpha_3^3$ shown in Fig. 5(a) does not point out the local patterns onto the T-shirt, while on the contrary, the low-level attention map $\alpha_3^1$ filters out $\tilde{\mathbf{F}}_3^{1,1}$ that typically reflects these texture patterns.

### 4.3. Component Analysis

We also demonstrate the cases when dropping partial attention modules or connections in comparison with the complete AF-net. As an example, the person ReID on VIPeR dataset [8] with six typical configurations are compared in Fig. 6(a). The orange bar shown in its bottom indicates the performance with the complete AF-net, while the yellow one is the M-net which is considered as the baseline model without the attention modules. The rest four bars are configured as:

*(1) Blue: naïve attention modules per branch.* In each branch of AF-net, a naïve attention module is applied to extract the attentive features $\tilde{\mathbf{F}}^{i,i}, i \in \Omega = \{1, 2, 3\}$.

*(2) Cyan: discarding the middle-level attention maps and attentive features.* We discard both the attention maps and attentive features of the block $2^{\text{nd}}$, *i.e.* prune the modules that produce $\tilde{\mathbf{F}}^{2,k}$ and $\tilde{\mathbf{F}}^{i,2}, \forall i, k \in \{1, 2, 3\}$.

*(3) Purple: pruning one branch.* It discards the first MDA module $\mathcal{F}(\alpha^1)$.

*(4) Light purple: pruning two branches.* The first two MDA modules $\mathcal{F}(\alpha^1)$ and $\mathcal{F}(\alpha^2)$ are discarded.

The results clearly prove that either cutting down the number of MDA modules or connections within this module will pull down the performance, and it is reasonable that these attention components complement each other to generate the comprehensive feature representation and thus gain a higher accuracy. Two examples with Top-5 identification results shown in Fig. 6(b-c) further demonstrate the
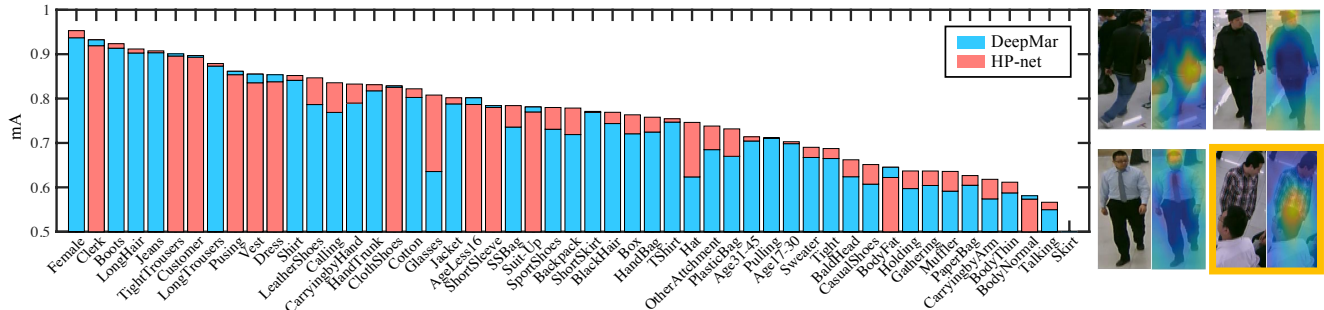
Figure 7. Mean accuracy scores for all attributes of RAP dataset by HP-net and DeepMar marked with red and blue bars respectively. The bars are sorted according to the larger mAs between two methods. The HP-net outperforms DeepMar especially on "glasses" and "hat" which have the exemplar sample listed aside. The sample in orange provides a failure case of predicting the attribute "talking".

| | PETA | RAP | **PA-100K** |
|---|---|---|---|
| # scene | - | 26 | **598** |
| # sample | 19,000 | 41,585 | **100,000** |
| # attribute | 61 (+4) | 69 (+3) | **26** |
| # tracklet | - | - | **18,206** |
| resolution | from $17 \times 39$ to $169 \times 365$ | from $36 \times 92$ to $344 \times 554$ | **from** $50 \times 100$ **to** $758 \times 454$ |

Table 1. Comparison of the proposed PA-100K dataset with existing datasets. The attribute number listed in the parentheses indicates the multi-class attributes while the one outside means the number of binary attributes.

effectiveness and indispensability of each component of the entire AF-net. The complete network is superior to both the multi-level naïve attention modules (Fig. 6(b)) and the single MDA module (Fig. 6(c)).

## 5. Pedestrian Attribute Recognition

We evaluate our HP-net on two public datasets comparing the state-of-the-art methods. In addition, we further propose a new large-scale pedestrian attribute dataset PA-100K with larger scene diversities and amount of samples.

### 5.1. PA-100K Dataset

Most of existing public pedestrian attribute datasets [7, 14] only contain a limited number of scenes (at most 26) with no more than $50,000$ annotated pedestrians. To further evaluate the generality of the proposed method, we construct a new large-scale pedestrian attribute (PA) dataset named as PA-100K with $100,000$ pedestrian images from 598 scenes, and therefore offer a superiorly comprehensive dataset for pedestrian attribute recognition. To our best knowledge, it is to-date the largest dataset for pedestrian attribute recognition. We compare our PA-100K dataset with the other two publicly available datasets in Table 1.

The samples of one person in **PETA dataset** [7] are only annotated once by randomly picking one exemplar image, and therefore share the same annotated attributes even though some of them might not be visible and some other

attributes are ignored. Another limitation is that the random partition of the training, validation and test sets are conducted in the whole dataset with no consideration of the person's identity across images, which leads to unfair image assignment of one person in different sets. In **RAP dataset** [14], the high-quality indoor images with controlled lighting conditions contain much lower variances than those under unconstrained real scenarios. Moreover, some attributes are even highly imbalanced.

**PA-100K dataset** surpasses the the previous datasets both in quantity and diversity, as shown in Table 1. We define 26 commonly used attributes including global attributes like gender, age, and object level attributes like handbag, phone, upper-clothing and *etc*. The PA-100K dataset was constructed by images captured from real outdoor surveillance cameras which is more challenging. Different from the existing datasets, the images were collected by sampling the frames from the surveillance videos, which makes some future applications available, such as video-based attribute recognition and frame-level pedestrian quality estimation. We annotated all pedestrians in each image and abandoned pedestrians with blurred motion or extreme low resolution (lower than $50 \times 100$). The whole dataset is randomly split into training, validation and test sets with a ratio of $8 : 1 : 1$. The samples of one person was extracted along its tracklets in a surveillance video, and they are randomly assigned to one of these sets, in which case PA-100K dataset ensures the attributes are learned independent of the person's identity. All these sets are guaranteed to have positives and negatives of the 26 attributes. Note that this partition based on tracklets is fairer than the partition that randomly shuffles the images in PETA dataset.

In the following experiments, we employ five evaluation criteria[3] including a label-based metric mean accuracy (mA), and four instance-based metrics, *i.e.* accuracy, precision, recall and F1-score. To address the issue of imbalanced classes, we adopt a weighted cross-entropy loss func-

---
[3]Criterion definitions are the same as those in [14].

| Dataset | RAP | | | | | | | PETA | | | | | | | PA-100K | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Method | ELF-mm | FC7-mm | FC6-mm | ACN | Deep-Mar | M-net | HP-net | ELF-mm | FC7-mm | FC6-mm | ACN | Deep-Mar | M-net | HP-net | Deep-Mar | M-net | HP-net |
| mA | 69.94 | 72.28 | 73.32 | 69.66 | 73.79 | 74.44 | **76.12** | 75.21 | 76.65 | 77.96 | 81.15 | **82.6** | 80.58 | 81.77 | 72.7 | 72.3 | **74.21** |
| Accu | 29.29 | 31.72 | 33.37 | 62.61 | 62.02 | 64.99 | **65.39** | 43.68 | 45.41 | 48.13 | 73.66 | 75.07 | 75.68 | **76.13** | 70.39 | 70.44 | **72.19** |
| Prec | 32.84 | 35.75 | 37.57 | **80.12** | 74.92 | 77.83 | 77.33 | 49.45 | 51.33 | 54.06 | 84.06 | 83.68 | 84.81 | **84.92** | 82.24 | 81.7 | **82.97** |
| Recall | 71.18 | 71.78 | 73.23 | 72.26 | 76.21 | 77.89 | **78.79** | 74.24 | 75.14 | 76.49 | 81.26 | 83.14 | 82.9 | **83.24** | 80.42 | 81.05 | **82.09** |
| F1 | 44.95 | 47.73 | 49.66 | 75.98 | 75.56 | 77.86 | **78.05** | 59.36 | 61 | 63.35 | 82.64 | 83.41 | 83.85 | **84.07** | 81.32 | 81.38 | **82.53** |

Table 2. Quantitative results(%) on three datasets for pedestrian attribute recognition, compared with previous benchmark methods.
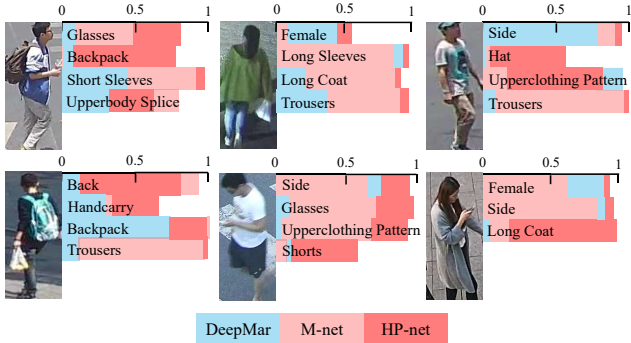


Figure 8. Comparison results between DeepMar, M-net and HP-net on partial ground truth attributes annotated for the given examples. Different colors represent different methods. Bars are plot by the prediction probabilities.

tion introduced by [13].

## 5.2. Comparison with the Prior Arts

We quantitatively and qualitatively compare the performance of the proposed method with the previous state-of-the-art methods on the previously mentioned three datasets. The following comparisons keep the same settings as the prior arts on different datasets respectively.

**Quantitative Evaluation.** We list the results of each method on RAP, PETA and PA-100K datasets in Table 2. Six reference methods are selected to be compared with the proposed model. The first three models are based on SVM classifier with hand-crafted features (ELF-mm [9, 22]) and deep-learned features (FC7-mm and FC6-mm) respectively. ACN [24] and DeepMar [13] are CNN models that achieved good performances by joint training the multiple attributes.

The baseline M-net and the proposed final model significantly outperform the state-of-the-art methods. We are also interested in the performance of each attribute. The bar in Fig. 7 shows the overlapped histograms of the mean accuracy (mA) for all attributes by DeepMar and HP-net. The bars are sorted in descending order according to the larger mA between these methods at one attribute. We find that the envelope superimposing the histogram is always supported by the HP-net with prominent performance gain against DeepMar, and is extremely superior on attributes which require fine-grained localization, like glasses and handbags.

**Qualitative Evaluation.** Besides the quantitative results in

| Dataset | Market-1501 [38] | CUHK03 [16] | VIPeR [8] |
|---|---|---|---|
| # identities | 1501 | 1360 | 632 |
| # images | 32643 | 13164 | 1264 |
| # cameras | 6 | 2 | 2 |
| # training IDs | 750 | 1160 | 316 |
| # test IDs | 751 | 100 | 316 |
| # probe images | 3368 | 100 | 316 |
| # gallery images | 19732 | 100 | 316 |

Table 3. The specifications of three evaluated ReID datasets.

Table 2, we also conduct qualitative evaluations for exemplar pedestrian images. As shown in the examples in Fig. 7, sample images from RAP dataset and their attention maps demonstrate the localizability of the learned attention maps. Especially in the first image, the attention map highlights two bags simultaneously. We also notice a failure case on the attribute "talking" which is irrelevant to a certain region but requires a global understanding of the whole image.

For the PA-100K dataset, we show attribute recognition results for several exemplar pedestrian images in Fig. 8. The bars indicate the prediction probabilities. Although the probabilities of one attribute do not directly imply its actual recognition confidences, they uncover the discriminative power of different methods as the lower probability corresponds to ambiguity or difficulty in correctly predicting one attribute. The proposed HP-net reliably predicts these attributes with region-based saliency, like "glasses", "back-pack", "hat", "shorts" and "handcarry".

## 6. Person Re-identification

Referring to the person re-identification, we also evaluate the HP-net with several reference methods on three publicly available datasets, quantitatively and qualitatively.

## 6.1. Datasets and Setups

The proposed approach is evaluated on three publicly standard datasets, including CUHK03 [16], VIPeR [8], and Market-1501 [38]. A summary about the statistical information of the three datasets are listed in Table 3. For the Market-1501 dataset, the same data separation strategy is used as [38]. For the other datasets, the training, validation and testing images are sampled based on the strategy introduced in [30]. The training and validation identities are guaranteed to have no overlaps with the testing ones for all

| CUHK03 | Top-1 | Top-5 | Top-10 | Top-20 |
|---|---|---|---|---|
| PersonNet [29] | 64.8 | 89.4 | 94.9 | 98.2 |
| JSTL [30] | 75.3 | - | - | - |
| Joint ReID [1] | 54.7 | - | - | - |
| LOMO-XQDA [17] | 52.2 | - | - | - |
| M-net | **88.2** | 98.2 | 99.1 | 99.5 |
| HP-net | **91.8** | 98.4 | 99.1 | 99.6 |
| **VIPeR** | **Top-1** | **Top-5** | **Top-10** | **Top-20** |
| NFST [36] | 51.2 | 82.1 | 90.5 | 96.0 |
| SCSP [3] | **53.5** | 82.6 | 91.5 | 96.7 |
| GOG+XQDA [19] | 49.7 | 79.7 | 88.7 | 94.5 |
| TCP [5] | 47.8 | 74.7 | 84.8 | 91.1 |
| M-net | 51.6 | 73.1 | 81.6 | 88.3 |
| HP-net | **56.6** | 78.8 | 87.0 | 92.4 |
| **Market-1501** | **Top-1** | **Top-5** | **Top-10** | **Top-20** |
| WARCA-L [11] | 45.2 | 68.1 | 76.0 | 84.0 |
| LOMO+CN [27] | 61.6 | - | - | - |
| S-CNN [26] | **65.9** | - | - | - |
| BoW-best [38] | 44.4 | 63.9 | 72.2 | 79.0 |
| M-net | 73.1 | 89.5 | 93.4 | 96.0 |
| HP-net | **76.9** | 91.3 | 94.5 | 96.7 |

Table 4. Experimental results(%) of the proposed HP-net and other comparisons on three datasets. The CMC Top-1-5-10-20 accuracies are reported. The Top-1 accuracies of two best performing approaches are marked in bold.

evaluated datasets. Following the pipeline of JSTL [30], all the training samples are combined together to train a single ReID model from scratch, which can be directly evaluated on all the testing datasets.

The widely applied cumulative match curve (CMC) metric is adopted for quantitative evaluation. While in the matching process, the cosine distance is computed between each query image and all the gallery images, and the ranked gallery list is returned. All the experiments are conducted under the setting of single query and the testing procedure is repeated 100 times to get an average result.

## 6.2. Performance Comparisons

**Quantitative Evaluation.** As shown in Table 4, the proposed approach is compared with a series of the deep neural networks like PersonNet [29], the multi-domain CNN JSTL [30], the Joint ReID method [1], and the horizontal occurrence model LOMO-XQDA [17] on CUHK03 [16]. As for the VIPeR [8] dataset, the null space semi-supervised learning method NFST [36], the similarity learning method SCSP [3], the hierarchical Gaussian model GOG+XQDA [19], and the triplet loss model TCP [5] are selected for comparison. The Market-1501 [38] dataset is also evaluated with the metric learning WARCA-L [11], a novel Siamese LSTM architecture LOMO+CN [27], the Siamese CNN with learnable gate S-CNN [26], and the bag of words model BoW-best [38].

Besides the results of the proposed approach with complete HP-net, the results of the M-net are also listed as the baseline for the three datasets. From Table 4, we can ob-
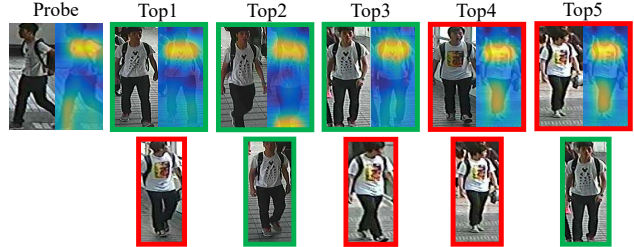


Figure 9. Comparison results between HP-net and M-net. For the probe images, the Top-5 retrieval results of HP-net together with attention maps are shown in the first row, and the results of M-net are shown in the second row.

serve that the proposed approach achieves the Top-1 accuracies of 91.8%, 56.6% and 76.9% on the CUHK03, ViPeR and Market-1501 datasets, respectively, and it achieves the state-of-the-art performance on all the three datasets. Moreover, even though the M-net can achieve quite satisfactory results on all datasets, the proposed pipeline can further improve the Top-1 accuracies by 3.6%, 5.0%, and 3.8% for each dataset, respectively.

**Qualitative Evaluation.** To highlight the performance of the proposed method on extracting localized semantic features, one query image together with its Top-5 gallery results by the proposed method and the M-net are visualized in Fig. 9. We observe that the proposed approach improves the rankings of the M-net and gets the correct results. By visualizing the attention maps from HP-Net of the query images and the Top-5 gallery images of both methods, we observe that the proposed attention modules can successfully locate the T-shirt patterns, in which the fine-grained features are extracted and discriminatingly identify the query person against the other identities with similar dressing.

## 7. Conclusion

In this paper, we present a new deep architecture called HydraPlus network with a novel multi-directional attention mechanism. Extensive ablation studies and experimental evaluations have manifested the effectiveness of the HP-net to learn multi-level and multi-scale attentive feature representations for fine-grained tasks in pedestrian analysis, like pedestrian attribute recognition and person re-identification. In the end, a new large-scale attribute dataset PA-100K is introduced to facilitate various pedestrian analysis tasks.

# References

[1] E. Ahmed, M. Jones, and T. K. Marks. An improved deep learning architecture for person re-identification. In *CVPR*, 2015. 8

[2] J. Ba, V. Mnih, and K. Kavukcuoglu. Multiple object recognition with visual attention. *arXiv preprint arXiv:1412.7755*, 2014. 2

[3] D. Chen, Z. Yuan, B. Chen, and N. Zheng. Similarity learning with spatial constraints for person re-identification. In *CVPR*, 2016. 8

[4] L.-C. Chen, Y. Yang, J. Wang, W. Xu, and A. L. Yuille. Attention to scale: Scale-aware semantic image segmentation. In *CVPR*, 2016. 2

[5] D. Cheng, Y. Gong, S. Zhou, J. Wang, and N. Zheng. Person re-identification by multi-channel parts-based cnn with improved triplet loss function. In *CVPR*, 2016. 3, 8

[6] M. Cornia, L. Baraldi, G. Serra, and R. Cucchiara. A deep multi-level network for saliency prediction. *arXiv preprint arXiv:1609.01064*, 2016. 2

[7] Y. Deng, P. Luo, C. C. Loy, and X. Tang. Pedestrian attribute recognition at far distance. In *Proceedings of the 22nd ACM international conference on Multimedia*, 2014. 6

[8] D. Gray, S. Brennan, and H. Tao. Evaluating appearance models for recognition, reacquisition, and tracking. In *Proc. IEEE International Workshop on Performance Evaluation for Tracking and Surveillance (PETS)*, 2007. 5, 7, 8

[9] D. Gray and H. Tao. Viewpoint invariant pedestrian recognition with an ensemble of localized features. In *ECCV*, 2008. 7

[10] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015. 3

[11] C. Jose and F. Fleuret. Scalable metric learning via weighted approximate rank component analysis. *arXiv preprint arXiv:1603.00370*, 2016. 8

[12] M. Koestinger, M. Hirzer, P. Wohlhart, P. M. Roth, and H. Bischof. Large scale metric learning from equivalence constraints. In *CVPR*, 2012. 3

[13] D. Li, X. Chen, and K. Huang. Multi-attribute learning for pedestrian attribute recognition in surveillance scenarios. In *Pattern Recognition (ACPR), 2015 3rd IAPR Asian Conference on*, 2015. 1, 2, 7

[14] D. Li, Z. Zhang, X. Chen, H. Ling, and K. Huang. A richly annotated dataset for pedestrian attribute recognition. *arXiv preprint arXiv:1603.07054*, 2016. 6, 7

[15] S. Li, T. Xiao, H. Li, B. Zhou, D. Yue, and X. Wang. Person search with natural language description. In *CVPR*, 2017. 3

[16] W. Li, R. Zhao, T. Xiao, and X. Wang. Deepreid: Deep filter pairing neural network for person re-identification. In *CVPR*, 2014. 3, 7, 8

[17] S. Liao, Y. Hu, X. Zhu, and S. Z. Li. Person re-identification by local maximal occurrence representation and metric learning. In *CVPR*, 2015. 3, 8

[18] J. Lu, J. Yang, D. Batra, and D. Parikh. Hierarchical question-image co-attention for visual question answering. In *NIPS*, 2016. 2

[19] T. Matsukawa, T. Okabe, E. Suzuki, and Y. Sato. Hierarchical gaussian descriptor for person re-identification. In *CVPR*, 2016. 8

[20] V. Mnih, N. Heess, A. Graves, et al. Recurrent models of visual attention. In *NIPS*, 2014. 2

[21] A. Newell, K. Yang, and J. Deng. Stacked hourglass networks for human pose estimation. In *ECCV*, 2016. 2, 4, 5

[22] B. Prosser, W.-S. Zheng, S. Gong, T. Xiang, and Q. Mary. Person re-identification by support vector ranking. In *BMVC*, 2010. 7

[23] C. Su, S. Zhang, J. Xing, W. Gao, and Q. Tian. Deep attributes driven multi-camera person re-identification. In *ECCV*, 2016. 3

[24] P. Sudowe, H. Spitzer, and B. Leibe. Person attribute recognition with a jointly-trained holistic cnn model. In *CVPR*, 2015. 1, 2, 7

[25] E. Ustinova, Y. Ganin, and V. Lempitsky. Multiregion bilinear convolutional neural networks for person re-identification. *arXiv preprint arXiv:1512.05300*, 2015. 3

[26] R. R. Varior, M. Haloi, and G. Wang. Gated siamese convolutional neural network architecture for human re-identification. In *ECCV*, 2016. 3, 8

[27] R. R. Varior, B. Shuai, J. Lu, D. Xu, and G. Wang. A siamese long short-term memory architecture for human re-identification. In *ECCV*, 2016. 8

[28] R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992. 2

[29] L. Wu, C. Shen, and A. v. d. Hengel. Personnet: person re-identification with deep convolutional neural networks. *arXiv preprint arXiv:1601.07255*, 2016. 3, 8

[30] T. Xiao, H. Li, W. Ouyang, and X. Wang. Learning deep feature representations with domain guided dropout for person re-identification. In *CVPR*, 2016. 1, 3, 8

[31] T. Xiao, S. Li, B. Wang, L. Lin, and X. Wang. Joint detection and identification feature learning for person search. In *CVPR*, 2017. 3

[32] T. Xiao, Y. Xu, K. Yang, J. Zhang, Y. Peng, and Z. Zhang. The application of two-level attention models in deep convolutional neural network for fine-grained image classification. In *CVPR*, 2015. 2

[33] H. Xu and K. Saenko. Ask, attend and answer: Exploring question-guided spatial attention for visual question answering. In *ECCV*, 2016. 2

[34] K. Xu, J. Ba, R. Kiros, K. Cho, A. C. Courville, R. Salakhutdinov, R. S. Zemel, and Y. Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*, 2015. 2, 4, 5

[35] S. Zagoruyko and N. Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. *arXiv preprint arXiv:1612.03928*, 2016. 2

[36] L. Zhang, T. Xiang, and S. Gong. Learning a discriminative null space for person re-identification. *arXiv preprint arXiv:1603.02139*, 2016. 8

[37] N. Zhang, M. Paluri, M. Ranzato, T. Darrell, and L. Bourdev. Panda: Pose aligned networks for deep attribute modeling. In *CVPR*, 2014. 3

[38] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian. Scalable person re-identification: A benchmark. In *ICCV*, 2015. 7, 8