

# Attribute Recognition from Adaptive Parts

Luwei Yang<sup>1</sup>  
luweiy@sfu.ca

Ligeng Zhu<sup>2</sup>  
zhuligeng@zju.edu.cn

Yichen Wei<sup>3</sup>  
yichenw@microsoft.com

Shuang Liang<sup>4</sup>  
shuangliang@tongji.edu.cn

Ping Tan<sup>1</sup>  
pingtan@sfu.ca

<sup>1</sup> Simon Fraser University  
Vancouver, Canada

<sup>2</sup> Zhejiang University  
Hangzhou, China

<sup>3</sup> Microsoft Research Asia  
Beijing, China

<sup>4</sup> Tongji University  
Shanghai, China

---

## Abstract

Previous part-based attribute recognition approaches perform part detection and attribute recognition in separate steps. The parts are not optimized for attribute recognition and therefore could be sub-optimal. We present an end-to-end deep learning approach to overcome the limitation. It generates object parts from key points and perform attribute recognition accordingly, allowing adaptive spatial transform [10] of the parts. Both key point estimation and attribute recognition are learnt jointly in a multi-task setting. Extensive experiments on two datasets verify the efficacy of the proposed end-to-end approach.

## 1 Introduction

Object attribute recognition is of central importance for object retrieval [12]. It has been extensively studied in vision for face [2, 14], person [4, 8, 24], animals [15] or more general objects [12]. The definition of ‘attribute’ is loose. Some attributes are more abstract and can only be observed from a holistic view, such as ‘face is attractive’, or ‘animal eats fish’. Some are more concrete and well associated with object parts. For example, ‘face has beard’ can be observed around mouth, ‘person wears a shirt’ around torso, ‘bird has long tail’ around tail. We call such attributes *localized attributes*. Their recognition is critical for fine-grained object classification [7].

Attribute recognition in the first category is usually treated as classification of the whole object. This is undesirable for localized attributes. Because only local regions are useful, using the whole object incurs the risk of over-fitting and is less robust due to object pose variations. Most previous works address the problem with a two-step approach: object parts are firstly detected and then used for attribute recognition. Works in [4, 24] find dozens of human body poselets [3], concatenate the features on them, and then train attribute classifiers. Gkioxari *et al.* [8] trains three part detectors (head, torso, legs) using deep features, and combines the features on the parts similarly. Lin *et al.* [17] firstly detects two parts (head,

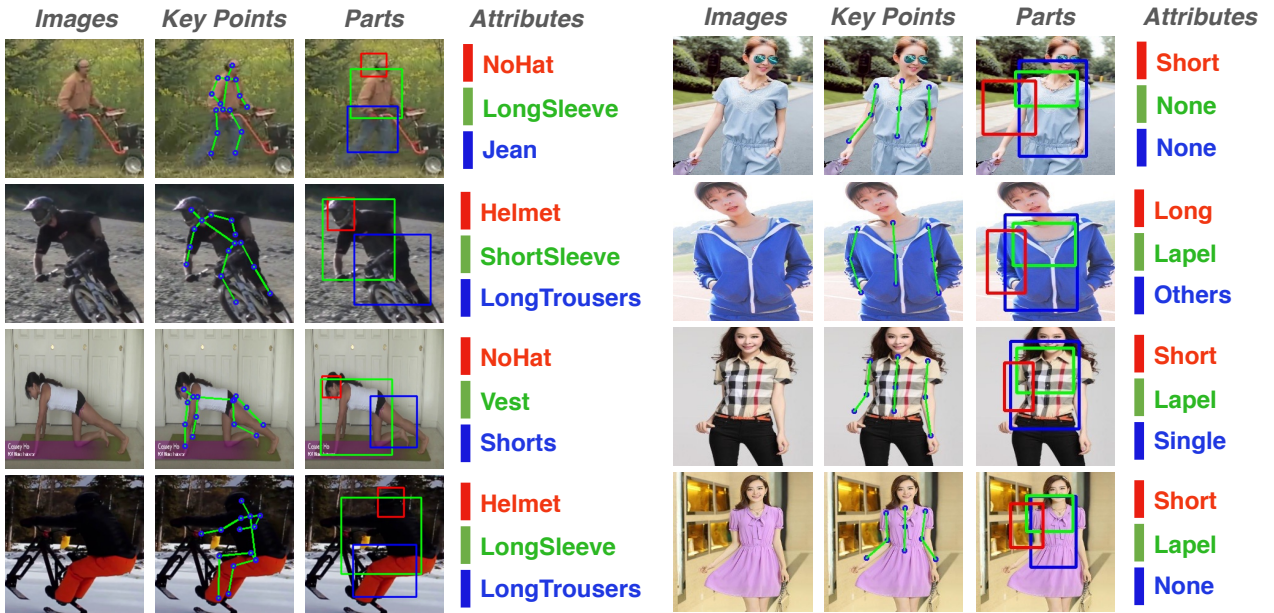


Figure 1: Overview of our approach. Given an image, we estimate the key points, generate object parts accordingly, and predict attributes of the parts. The learning is end-to-end in a single deep neural network. See experiment for details of the results.

body), rectifies the two parts to pose-aligned parts by comparing against a pre-defined part template database, and performs classification similarly.

These part-based approaches are superior than using the whole object. Yet, they still have flaws that are originated from the loose relation between the parts and the goal of attribute recognition. The part detectors are trained from bounding boxes that are either manually annotated [8, 17] or from heuristic clustering (poselet in [4, 24] and part template database in [17]). Such parts are not directly optimized for attribute classification and the training is not end-to-end. They could be sub-optimal for different types of attributes. For example, ‘has long hair’ attribute may require a large bounding box around head while ‘wear short sleeve’ attribute may only need a tight region covering the upper arms. It is challenging to obtain good parts in advance without knowing the attribute recognition goal.

In this work, we propose an end-to-end learning approach for localized attribute recognition, for the first time up to our knowledge. Instead of training part detector separately, we firstly estimate object key points as an auxiliary task. Because the definition of key point is clear, their annotation is less ambiguous than part bounding boxes. From the key points, the object parts are generated adaptively, with free parameters to adjust its spatial extent. This adaptive part generation is inspired by the recent spatial transformer network [10], which can learn image spatial transformation from the image classification goal. Instead of applying the transform to the whole image [10], we apply a spatial transform for each part and use the bilinear sampler [10] to warp the image features for subsequent attribute recognition. The whole network is learnt end-to-end in a multi-task setting, with attribute classification as the main task and key point prediction as the auxiliary one. Our pipeline is exemplified in Figure 1. The network framework is illustrated in Figure 2.

Very few public dataset has both complex object pose and rich attribute annotation. Therefore, we create two new datasets from the existing ones. The first is augmented from the currently largest human pose dataset MPII [1]. We labeled 11 clothing attributes on three body parts: head, torso and legs. It is larger and richer than the previous human attribute datasets used in [4, 24]. The second is refined from a recent garment database [6]. It con-

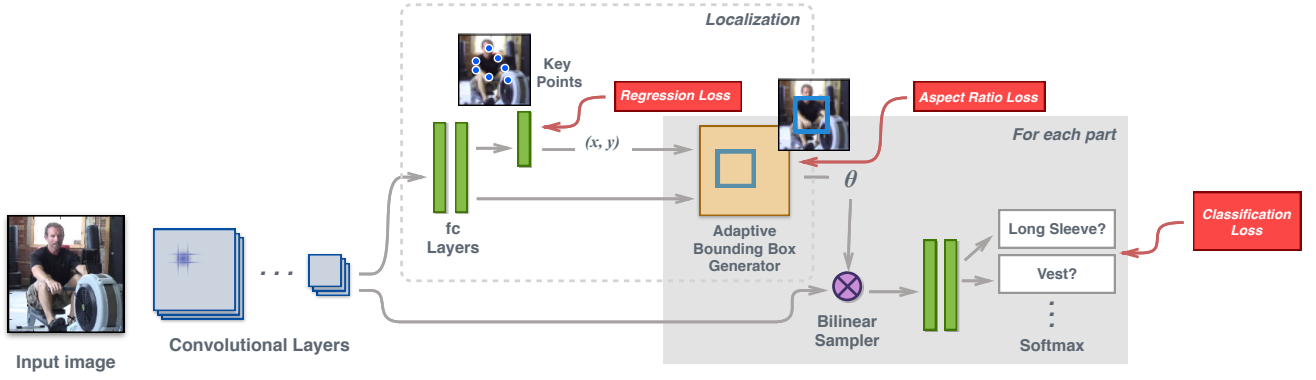


Figure 2: Overview of the network. It consists of initial convolutional feature extraction layers, a key point localization network, an adaptive bounding box generator for each part, and the final attribute classification network for each part. Besides the final classification loss, there is also intermediate key point regression loss and a regularization loss on bounding box aspect ratio. See text for details. We note that only one part is visualized here for clarify.

tains fine-grained and highly localized attributes on collar, sleeve and button types. *Both datasets will be released.* Extensive experiment comparison results verify the efficacy of our end-to-end learning approach.

## 2 End-to-End Learning with Adaptive Parts

Our approach is demonstrated on human attribute recognition. Our network architecture is illustrated in Figure 2. It consists of a convolutional network for feature extraction, a localization network for key point estimation, an adaptive bounding box generator for each part, and part based feature sampler and attribute classifier. They are elaborated as follows.

**Convolutional Feature Exaction** Given an input image, a convolutional neural network is used to extract feature maps. The features are shared for all subsequent tasks for computational efficiency. We use AlexNet [13] (first 5 convolutional layers) and VGG-16 [21] (first 13 convolutional layers) in our experiments. The input image size is  $227 \times 227$  for AlexNet and  $224 \times 224$  for VGG-16.

**Key Point Estimation** As shown in Figure 2 and 3, there are three fully-connected layers after the convolutional features, with output dimensions 2048, 2048 and  $2N$ , respectively. Here  $N$  is the number of key points. After each fc layer, ReLU activation function and drop out layer (ratio 0.5) is used. We use L2 distance loss for key point estimation,  $\sum_i^N \|\hat{p}_i - p_i\|_2^2$ , where  $\hat{p}_i, p_i \in \mathbb{R}^2$  are the normalized ground truth and estimation for key point  $i$ .

**Adaptive Part Generation** Some attributes are clearly associated with certain object parts. We encode such prior knowledge by specifying a subset of key points  $\mathcal{P}_t$  for each part  $t$ . For example, in Figure 2, for part torso we have  $\mathcal{P}_{torso} = \{shoulders, elbows, wrists, hips\}$ . The initial part bounding box  $b_t = [w_t, h_t, x_t, y_t]$  is defined as an enlarged bounding box of key

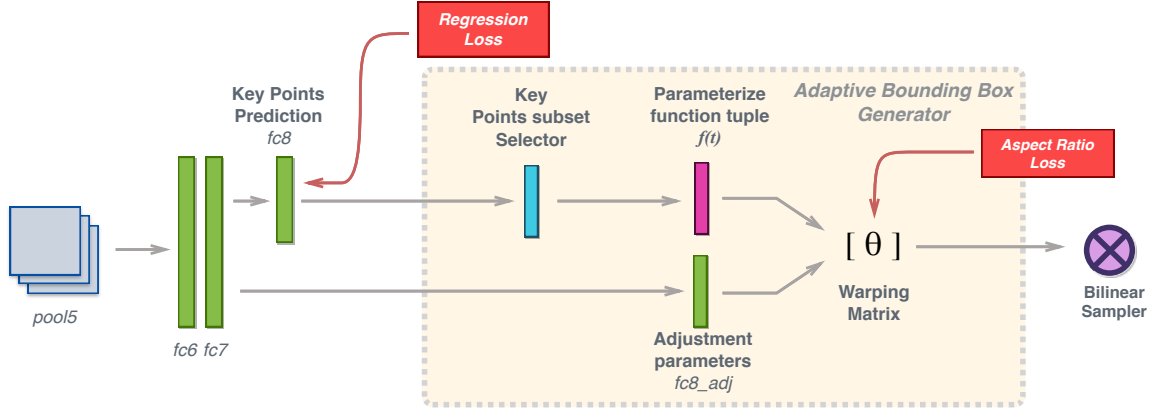


Figure 3: The details of Bounding Box Generator in Figure 2. It has two streams. The first one takes a subset of key point prediction as input and outputs an initial part bounding box. The second one predicts the bounding box adjustment parameters, using the previous fc7 layer’s features. The output of Bounding Box Generator is a 2x3 warping transformation matrix  $\theta$ . The key point regression loss is for learning the key points. The aspect ratio loss is to regularize the part bounding box estimation. See text for details.

points in  $\mathcal{P}_t$ ,

$$w_t = s \cdot (\max_x(\mathcal{P}_t) - \min_x(\mathcal{P}_t)), \quad h_t = s \cdot (\max_y(\mathcal{P}_t) - \min_y(\mathcal{P}_t)), \quad (1)$$

$$x_t = \frac{1}{2} (\min_x(\mathcal{P}_t) + \max_x(\mathcal{P}_t) - w_t), \quad y_t = \frac{1}{2} (\min_y(\mathcal{P}_t) + \max_y(\mathcal{P}_t) - h_t). \quad (2)$$

Here,  $\max_x(\mathcal{P}_t)$  is the maximum  $x$  coordinate in key points of  $\mathcal{P}_t$ . Other notations are similar.  $w/h$  is the initial box’s width/height.  $x/y$  is initial box’s upper left corner.  $s$  is a constant scalar larger than 1. It is set to 1.5 in our experiment.

The initial bounding box is then adaptively adjusted by free parameters  $\Delta = [\Delta_w, \Delta_h, \Delta_x, \Delta_y]$ . The final bounding box is defined as  $[w_t(1 + \Delta_w), h_t(1 + \Delta_h), x_t + \Delta_x, y_t + \Delta_y]$ . To learn the adjustment parameters, we add one more fully connected layer ( $fc8\_adj$ ) to the previous layer (fc7), with 4 output values. This is depicted in Figure 3.

The final bounding box could have too distorted dimensions due to the free adjustment parameters. Inspired by [20], in order to alleviate this issue, we introduce a bounding box aspect ratio loss as

$$L_r^t = \begin{cases} \frac{1}{2} \{ [\alpha [h_t(1 + \Delta_h)]^2 - w_t(1 + \Delta_w)]^2 \}_+ & \text{if } h_t(1 + \Delta_h) > w_t(1 + \Delta_w) \\ \frac{1}{2} \{ [\alpha [w_t(1 + \Delta_w)]^2 - h_t(1 + \Delta_h)]^2 \}_+ & \text{if } w_t(1 + \Delta_w) > h_t(1 + \Delta_h) \end{cases}, \quad (3)$$

where  $\alpha$  is a ratio threshold (set to 0.6 in experiment). The loss is 0 when the value in bracket  $\{ \}_+$  is less than 0.

**Bilinear Feature Sampling and Attribute Recognition** For each part bounding box, the convolutional feature maps are warped accordingly. We use the Bilinear Sampler in [10] that warps a local region via bilinear interpolation. It allows the gradient flow into the localization network, serving as a bridge to link the localization and attribute classification networks.

The warping employs a  $2 \times 3$  affine transformation, parameterized as

$$\theta_t = \begin{bmatrix} w_t(1 + \Delta_w) & 0 & x_t + \Delta_x \\ 0 & h_t(1 + \Delta_h) & y_t + \Delta_y \end{bmatrix}. \quad (4)$$

The transformation  $\theta_t$  warps the coordinates  $(x', y')$  in the target feature map  $U$  back to the coordinates  $(x, y)$  in the source feature map  $V$ . Specifically, we have

$$V_{(x,y)} = \sum_{m=1}^W \sum_{n=1}^H U_{(m,n)} \max(0, 1 - |x' - m|) \max(0, 1 - |y' - n|), \quad (5)$$

where  $(x, y) \in \mathbb{R}^2$  and  $(m, n) \in \mathbb{R}^2$  are coordinates in  $V$  and  $U$ ,  $W/H$  are feature map dimensions. The  $(x', y') \in \mathbb{R}^2$  in  $U$  are warped coordinates that satisfy  $[x' \ y']^T = \theta_t [x \ y \ 1]^T$ . For more details, we refer readers to [10].

After bilinear feature sampling, there are two fully connected layers with output dimensions 512 and 256, respectively. After each, the ReLU activation and dropout regularization (ratio 0.5) are used. The softmax classification loss is appended at last.

All equations above are differentiable. The learning is end-to-end using stochastic gradient descent. For Eq. (4), we compute the gradient  $\frac{\partial \theta_t}{\partial b_t}$  and  $\frac{\partial \theta_t}{\partial \Delta}$ , respectively. For  $b_t$ , we record the max and min index of the key points during feed forward, and pass the error back to corresponding channels during backward propagation, similarly as max pooling.

### 3 Experiments

Our approach is implemented in Caffe [11]. We use SGD for network training. Mini-batch size is 128 for AlexNet [13] and 16 for VGG-16 [21]. The network is initialized with pre-trained models on ImageNet. The initial base learning rate is 0.0008. We train 60K iterations (107 epochs), decrease the learning rate by 0.1, and train another 60K iterations. Note that for the bounding box generator part we use a learning rate that is  $\frac{1}{10}$  of the base learning rate, similarly as in [10]. The weight of losses is 1.0 for key point regression, 0.3 for attribute classification of each part, and 0.1 for bounding box aspect ratio of each part. For the training data augmentation, we randomly translate, scale and mirror the images for 6 times.

**Datasets** There are few dataset with both complex object pose and rich attribute annotation. We augment the currently largest human pose dataset MPII [1] by labeling 11 clothing attributes. The augmented dataset has about 28K human instances, on only training images with pose ground truth annotation. The 11 attributes are multi-classes grouped for three parts: head, torso and legs. They are  $\{NoHat, HasHat, Helmet\}$  for head,  $\{LongSleeve, ShortSleeve, Vest, Naked\}$  for torso, and  $\{LongTrousers, Jean, Dress, Shorts\}$  for legs. This dataset is larger and richer than previous human attribute datasets [4, 24]. Dataset in [4] has about 8K instances and 9 binary attributes (6 about clothing). Dataset in [24] has about 25K instances and 8 binary attributes (3 about clothing).

We also used a recent Garment dataset [6]. Its images are from online shopping site and has fine-grained clothing attributes. The labels, however, are quite unbalanced on lower body since many persons are only upper body. Therefore, we select a subset of 4K images of upper body humans and 12 attributes in 3 multi-classes groups for collar, sleeve and button types. The attributes are  $\{Stand, Lapel, None, Others\}$  for collar,  $\{Long, Medium, Short,$

Attributes	AlexNet (8 Layers)					VGG-16 (16 Layers)				
	Full	Stn	Separa.	Ours	Oracle	Full	Stn	Separa.	Ours	Oracle
Helmet	68.30	<b>68.31</b>	51.23	67.69	76.99	81.76	80.58	57.60	<b>83.53</b>	84.04
HasHat	57.53	61.91	53.57	<b>64.57</b>	80.01	79.16	78.18	57.51	<b>81.21</b>	83.75
NoHat	92.06	93.21	89.11	<b>93.80</b>	96.86	96.39	<b>96.78</b>	88.30	96.45	97.15
Accuracy	76.27	77.40	70.43	<b>78.00</b>	84.02	84.25	83.96	74.00	<b>86.02</b>	86.16
LongSleeve	76.31	78.88	76.53	<b>81.64</b>	84.88	83.62	84.22	83.51	<b>87.89</b>	88.52
Vest	72.05	71.48	<b>74.44</b>	72.32	78.72	80.38	80.64	80.86	<b>81.57</b>	83.49
ShortSleeve	80.57	80.26	82.40	<b>84.34</b>	89.46	88.99	88.67	88.43	<b>91.35</b>	92.76
Naked	46.44	47.15	40.38	<b>55.42</b>	45.86	49.73	54.99	47.43	<b>61.18</b>	49.41
Accuracy	68.69	68.92	68.38	<b>73.20</b>	74.69	73.60	75.68	74.51	<b>79.68</b>	78.94
Jean	59.19	62.08	63.88	<b>65.38</b>	68.02	67.75	69.18	67.31	<b>69.58</b>	73.55
Dress	18.77	19.78	9.94	<b>34.33</b>	23.97	26.30	34.81	20.98	<b>38.45</b>	37.81
Shorts	88.44	88.29	86.52	<b>89.44</b>	90.50	91.46	91.21	89.42	<b>93.42</b>	92.97
LongTrousers	85.72	87.18	85.13	<b>87.72</b>	88.01	89.04	89.51	86.96	<b>90.83</b>	90.90
Accuracy	77.03	76.57	76.06	<b>77.74</b>	79.99	79.00	80.82	78.71	<b>82.22</b>	82.25

Table 1: Average Precision (%) of all attributes on our augmented MPII dataset. The three groups are for head, torso, and legs from top to bottom. For each group, the multi-classification accuracy is shown in the last ‘Accuracy’ row. The best results in each row are bold. Note that we exclude the ‘Oracle’ column from consideration.

*Others*} for sleeve, and {*Single, DoubleButton, None, Others*} for button. Such attributes are highly localized. It is hard to manually annotate parts for their recognition.

**Baselines** To validate our approach, we compare with various baselines. For fairness, all methods use the same input images, the same initial networks. The learning parameters are separately tuned for each method.

The first baseline is called **Full**. It does not use part information and directly learns attribute from the whole input image. The network architecture is a subset of Figure 2, by removing the localization sub-network, the adaptive bounding box generator and the bilinear sampler. It is straightforward and serves as the *lower bound* of all methods.

The second baseline resembles the spatial transformer network [10]. It is called **Stn**. It extends the *full* baseline by introducing a spatial transform for each part’s attribute classification. The network architecture is similarly to Figure 2, without the localization sub-network and the adaptive bounding box generator. Instead, it uses a fixed initial bounding box for the bilinear sampler for each part. The initial bounding box is located at the input image center. Its width/height is set to 60% of the image width/height. Such parameters are determined via cross-validation. Other values such as using the whole image as in [10] are found inferior.

The third baseline is similar to the previous approaches [4, 8, 24]. We call it **Separate**. The network is similar as Figure 2. The learning is performed in two separate steps. Firstly, the convolutional layers and pose localization network is trained to predict key points. From the fixed key points, the adaptive bounding box generator, the subsequent bilinear sampler, attribute classification layers as well as the initial convolutional layers are trained.

The last baseline is called **Oracle**. It is similar to the *separate* baseline. The difference is that it directly uses the ground truth key points instead of predicted key points, so we do not train the localization network. It performs attribute recognition with ‘perfect’ part localization, and serves as an *upper bound* of all methods.

	Adaptive× RatioLoss×	Adaptive√ RatioLoss×	Adaptive√ RatioLoss√
Head	84.73	85.30	<b>86.02</b>
Torso	78.45	77.68	<b>79.68</b>
Legs	81.34	81.65	<b>82.22</b>

Table 2: Attribute multi-classification accuracy (%) of three parts on MPII subset, using VGG-16. The symbol  $\sqrt$  indicates the module is enabled, otherwise  $\times$ . Note that we do not have Adaptive $\times$ -RatioLoss $\sqrt$  combination.

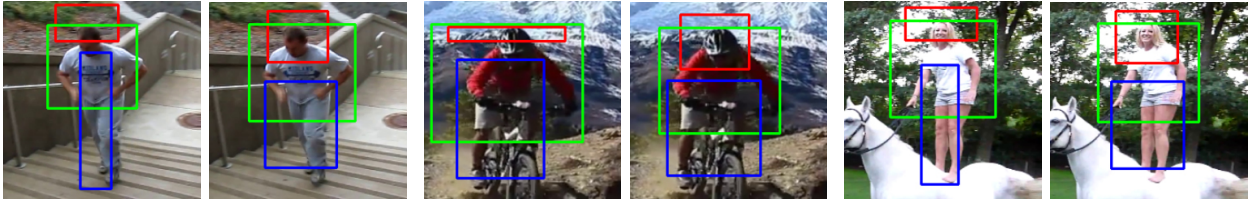


Figure 4: Side-by-side comparison of part bounding boxes without using (left) and using (right) aspect ratio loss. (Red: Head, Blue: Leg, Green: Torso)

**Results on MPII Dataset** Some persons in MPII dataset have invisible body parts. We exclude such images for simplicity. The remaining ones are divided into 12K for training and 3.5K for test. There are in total 14 key points. The three subsets of key points are  $\{head, neck, shoulders\}$  for head,  $\{shoulders, elbows, wrists, hips\}$  for torso, and  $\{hips, knees, ankles\}$  for legs. Note that both *shoulders* and *hips* are used for two parts.

Example part localization and attribute recognition results are shown in Figure 6. Table 1 reports the Average Precision of all attributes and the multi-classification accuracy of all parts. We make a few conclusions. Firstly, *Full* and *Stn* consistently perform worse than our approach. This shows that explicitly exploit pose information is beneficial for attribute recognition. Secondly, when pose is used, *Oracle* is the best and *Separate* is the worst, even worse than *Full* and *Stn* in most cases. This indicates that pose estimation could be hard to learn, and inaccurate part hurts attribute recognition. Lastly, when both tasks are jointly trained in our approach, the performance becomes much better and is almost on par with *Oracle* for most attributes, especially for VGG-16 network. This verifies the effectiveness of end-to-end multi task learning. Interestingly, our approach outperforms *Oracle* occasionally, for example, on torso-Naked attribute. This may be an evidence of the power of end-to-end learning over human knowledge (ground truth pose annotation and key point - attribute association).

To demonstrate the effectiveness of adaptive part learning (adjustment parameters in Figure 3) and the bounding box aspect ratio loss, we trained two more models by removing the corresponding modules respectively. As reported in Table 2, using both can clearly boost the performance. Note that only using adaptive part learning is not always better. For example, it decreases the accuracy on Torso. This is probably due to too much degrees of freedom of the part. Applying aspect ratio loss alleviates this problem and always helps. As illustrated in Figure 4, the detected part boxes are better when aspect ratio loss is used.

To evaluate key point prediction, we use the commonly adopted accuracy metric called Percentage of Detected Joints (PDJ) [19]. A key point prediction is considered correct if its Euclidean distance to ground-truth is smaller than a percentage of the ground truth torso length. Figure 5 shows the PDJ values over different thresholds for four joints of our approach and *Separate* baseline. It shows that learning part attribute also improves the key

Attributes	Full	Stn	Separate	Ours	Oracle
Single	55.79	55.26	54.93	<b>59.51</b>	59.75
None	90.07	91.1	89.75	<b>92.7</b>	90.62
DoubleButton	<b>13.63</b>	3.41	13.07	3.87	28.68
Others	44.1	<b>58.35</b>	47.55	46.23	51.8
Accuracy	75.12	75.99	74.75	<b>76.24</b>	76.61
Lapel	55.46	58.68	<b>64.93</b>	64.08	75.66
None	78.7	84.98	82.61	<b>89.85</b>	90.45
Stand	56.95	68.4	66.78	<b>73.16</b>	73.91
Others	23.93	21.95	26.01	<b>29.15</b>	46.65
Accuracy	61.63	65.35	65.72	<b>67.70</b>	73.51
Short	90.72	90.87	92.85	<b>92.9</b>	94.33
Medium	60.35	67.86	<b>69.02</b>	66.76	73.98
Long	65.22	72.16	70.23	<b>75.03</b>	73.6
Others	79.02	<b>80.69</b>	75.12	78.61	80.8
Accuracy	74.01	75.99	75.5	<b>77.48</b>	78.47

Table 3: Average Precision (%) of attribute groups in Garment dataset using AlexNet. The three groups are for collar, button, and sleeve from top to bottom. For each group, the multi-classification accuracy is shown in the last 'Accuracy' row. The best results in each row are bold. Note that we exclude the 'Oracle' column from consideration.

point localization, as the overfitting of key points regression may be reduced by the additional attribute classification task. The conclusion remains similar for other key points.

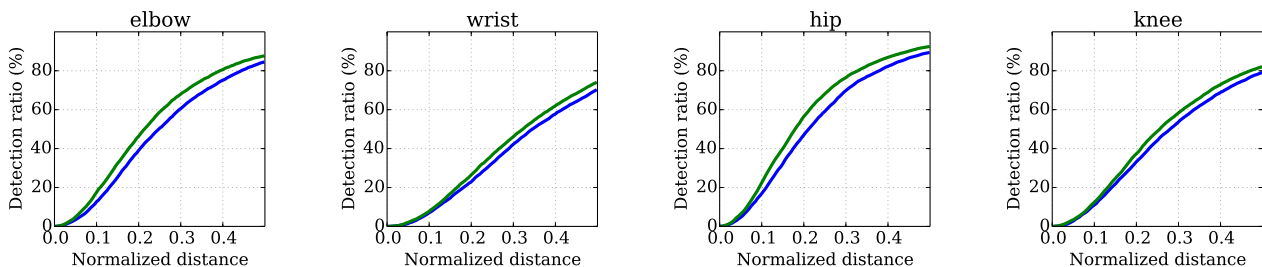


Figure 5: The Percentage of Detected Joints (PDJ) curves for four key points. The x-axis is the threshold normalized by torso length. The y-axis is the ratio of correct key joints. Results of our approach are in green. Results of *Separate* baseline are in blue.

We note that key point estimation is not the main focus of this work. We use a simple localization network as in Figure 2. There exists more sophisticated and better performing models for human pose estimation [5, 9, 16, 18, 22, 23]. It is yet unclear whether combining such models can improve attribute recognition, since using ground truth pose (*Oracle* baseline) is only marginally better than our approach, as shown in Table 1 (VGG-16). Such investigation is left as future work.

**Results on Garment Dataset** This dataset has in total 8 key points. The subsets for three parts are {*throat, shoulders, chest*} for Collar, {*right shoulders, right elbows, right wrists*} for Sleeve, and {*shoulders, elbows, hips*} for Button. The dataset is divided into training and testing sets with 3.2k and 0.9K images, respectively. As the number of images is small, we only use AlexNet. Example results are shown in Figure 6. Table 3 reports comparison results of our method and baselines. The conclusions are consistent with those of MPII dataset. Our method outperforms other baselines on most attributes. We note that there are





Figure 6: Examples results on MPII dataset (left) and Garment dataset (right). Last row shows some failure cases in red.

more fluctuations in the results. This is partially due to the small size of the dataset and imbalanced samples for different attributes.

## 4 Conclusion

We propose an end-to-end deep learning approach to jointly learn key point estimation and object attribute recognition. An adaptive part generation serves as an intermediate representation to connect the two tasks. Through joint learning, we overcome the limitation in previous two-step approaches and explicitly optimize the part location for attribute recognition. Our approach is validated on human attribute recognition on two datasets, via extensive experiment comparison.

## 5 Acknowledgement

This project is supported by the Canada NSERC Discovery project 611664 and Discovery Acceleration Supplement 611633. S. Liang is supported by The National Science Foundation of China (No.61305091), and The Fundamental Research Funds for the Central Universities (No.2100219054).

## References

- [1] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2D human pose estimation: New benchmark and state of the art analysis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- [2] Thomas Berg and Peter N Belhumeur. Poof: Part-based one-vs.-one features for fine-grained categorization, face verification, and attribute estimation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013.
- [3] Lubomir Bourdev and Jitendra Malik. Poselets: Body part detectors trained using 3d human pose annotations. In *IEEE International Conference on Computer Vision (ICCV)*, 2009.
- [4] Lubomir Bourdev, Subhransu Maji, and Jitendra Malik. Describing people: A poselet-based approach to attribute classification. In *IEEE International Conference on Computer Vision (ICCV)*, 2011.
- [5] Joao Carreira, Pulkit Agrawal, Katerina Fragkiadaki, and Jitendra Malik. Human pose estimation with iterative error feedback. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [6] Xiaowu Chen, Bin Zhou, Feixiang Lu, Lin Wang, Lang Bi, and Ping Tan. Garment modeling with a depth camera. *ACM Transactions on Graphics (TOG)*, 34(6):203, 2015.
- [7] Kun Duan, Devi Parikh, David Crandall, and Kristen Grauman. Discovering localized attributes for fine-grained recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.
- [8] Georgia Gkioxari, Ross Girshick, and Jitendra Malik. Actions and attributes from wholes and parts. In *IEEE International Conference on Computer Vision (ICCV)*, 2015.
- [9] Peiyun Hu and Deva Ramanan. Bottom-up and top-down reasoning with hierarchical rectified gaussians. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [10] Max Jaderberg, Karen Simonyan, and Andrew Zisserman. Spatial transformer networks. In *Advances in Neural Information Processing Systems (NIPS)*, 2015.
- [11] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014.

- [12] Adriana Kovashka, Devi Parikh, and Kristen Grauman. Whittlesearch: Interactive image search with relative attribute feedback. *International Journal on Computer Vision (IJCV)*, 115:185–210, 2015.
- [13] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems (NIPS)*, 2012.
- [14] Neeraj Kumar, Alexander C Berg, Peter N Belhumeur, and Shree K Nayar. Attribute and simile classifiers for face verification. In *IEEE International Conference on Computer Vision (ICCV)*, 2009.
- [15] Christoph H. Lampert, Hannes Nickisch, and Stefan Harmeling. Learning to detect unseen object classes by between-class attribute transfer. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- [16] Ita Lifshitz, Ethan Fetaya, and Shimon Ullman. Human pose estimation using deep consensus voting. *arXiv preprint arXiv:1603.08212*, 2016.
- [17] Di Lin, Xiaoyong Shen, Cewu Lu, and Jiaya Jia. Deep lac: Deep localization, alignment and classification for fine-grained recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [18] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. *arXiv preprint arXiv:1603.06937*, 2016.
- [19] Ben Sapp and Ben Taskar. Modec: Multimodal decomposable models for human pose estimation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013.
- [20] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [21] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *International Conference on Learning Representations (ICLR)*, 2015.
- [22] Jonathan Tompson, Arjun Jain, Yann LeCun, and Christoph Bregler. Joint training of a convolutional network and a graphical model for human pose estimation. In *Advances in Neural Information Processing Systems (NIPS)*, 2014.
- [23] Shih-En Wei, Varun Ramakrishna, Takeo Kanade, and Yaser Sheikh. Convolutional pose machines. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [24] Ning Zhang, Manohar Paluri, Marc’Aurelio Ranzato, Trevor Darrell, and Lubomir Bourdev. Panda: Pose aligned networks for deep attribute modeling. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.