

Unified Structured Learning for Simultaneous Human Pose Estimation and Garment Attribute Classification

Jie Shen, Guangcan Liu, *Member, IEEE*, Jia Chen, Yuqiang Fang, Jianbin Xie, *Member, IEEE*, Yong Yu, and Shuicheng Yan, *Senior Member, IEEE*

Abstract—In this paper, we utilize structured learning to simultaneously address two intertwined problems: 1) human pose estimation (HPE) and 2) garment attribute classification (GAC), which are valuable for a variety of computer vision and multimedia applications. Unlike previous works that usually handle the two problems separately, our approach aims to produce an optimal joint estimation for both HPE and GAC via a unified inference procedure. To this end, we adopt a preprocessing step to detect potential human parts from each image (i.e., a set of candidates) that allows us to have a manageable input space. In this way, the simultaneous inference of HPE and GAC is converted to a structured learning problem, where the inputs are the collections of candidate ensembles, outputs are the joint labels of human parts and garment attributes, and joint feature representation involves various cues such as pose-specific features, garment-specific features, and cross-task features that encode correlations between human parts and garment attributes. Furthermore, we explore the strong edge evidence around the potential human parts so as to derive more powerful representations for oriented human parts. Such evidences can be seamlessly integrated into our structured learning model as a kind of energy function, and the learning process could be performed by standard structured support vector machines algorithm. However, the joint structure of the two problems is a cyclic graph, which hinders efficient inference. To resolve this issue, we compute instead approximate optima using an iterative

procedure, where in each iteration, the variables of one problem are fixed. In this way, satisfactory solutions can be efficiently computed by dynamic programming. Experimental results on two benchmark data sets show the state-of-the-art performance of our approach.

Index Terms—Human pose estimation, garment attribute classification, joint inference, structured learning.

I. INTRODUCTION

HUMAN-ORIENTED technologies play important roles in many computer vision and multimedia applications that require interactions between persons and electronic devices. The significance of human-oriented technologies naturally drives the research community to extensively investigate human-related topics, such as face recognition [1], human tracking [2], pose estimation [3], clothing technology [4], etc. In this work, we are interested in two of them: human pose estimation (HPE) and clothing technology (CT). Both problems have been studied extensively, and a review of previous work is presented in the following section.

A. Previous Work

1) *Human Pose Estimation*: The literature about HPE can trace back to 40 years ago. Fischler and Elschlager [6] proposed to represent the articulated human pose as a collection of rigid body parts. This classical model, called *pictorial structure* (PS) in [7], provides a straightforward representation for articulated objects and owns a tree structure that can facilitate efficient inference. Hence, PS is still adopted as a basic tool by recently established approaches, see [8]–[13]. These recent works mainly pursue better feature description, more efficient computation, more complex human structures and more effective contextual information.

Feature description is one of the key elements for various vision tasks and so HPE [14]. In [12], an iterative parsing paradigm was introduced to obtain an increasingly finer feature scheme for describing human parts. Other works, see [8], [10], employed shape-based feature descriptors such as Shape Context [15] and Histogram of Oriented Gradients (HOG) [16], which proved to be more effective than color-based features. In [9], Eichner and Ferrari considered the appearance consistence between connected/symmetric limbs and developed a better appearance model.

Manuscript received December 6, 2013; revised April 16, 2014 and June 20, 2014; accepted August 26, 2014. Date of publication September 15, 2014; date of current version September 30, 2014. The work of J. Xie was supported by the National Natural Science Foundation of China under Grant 61303188. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. David Frakes.

J. Shen is with the Department of Computer Science and Engineering, Shanghai Jiao Tong University, Shanghai 200240, China (e-mail: jieshen@apex.sjtu.edu.cn).

G. Liu is with the School of Information and Control Engineering, Nanjing University of Information Science and Technology, Nanjing 210044, China (e-mail: gcliu@nuist.edu.cn).

J. Chen and Y. Yu are with the Department of Computer Science and Engineering, Shanghai Jiao Tong University, Shanghai 200240, China (e-mail: chenjie@apex.sjtu.edu.cn; yyyu@apex.sjtu.edu.cn).

Y. Fang is with the College of Mechatronic Engineering and Automation, National University of Defense Technology, Changsha 410073, China (e-mail: fangyuqiang@nudt.edu.cn).

J. Xie is with the College of Electronic Science and Engineering, National University of Defense Technology, Changsha 410073, China (e-mail: jbxie@126.com).

S. Yan is with the Department of Electrical and Computer Engineering, National University of Singapore, Singapore 119077 (e-mail: eleyans@nus.edu.sg).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TIP.2014.2358082

The majority of early works on HPE emphasized on detection performance, i.e., a more effective inference schema, and the focus of later works was placed on the efficiency of inference. Typically, the search space of human pose is the main bottleneck for improving the efficiency, because one must search over the location and orientation space for each human part, as well as over the image pyramid. To reduce the search space of human poses, Ferrari et al. [11] utilized a generic upper-body detector and the GrabCut [17] algorithm, targeting at a reduced search space. In [8], [18], and [19], the tree structure is used to make the inference procedure more efficient, which can also well handle the spatial association between different human parts. One advantage of the tree structure is the ability to allow a fast computation via a dynamic programming [7]. Furthermore, deformable cost computation between connected human parts can be accelerated by performing a distance transform if the pair-wise cost suffices some specific conditions [7].

Although the tree-based PS models can draw a general representation for human body, it does not explore the implicit connections between rigid parts that do not share joints [20]. Therefore, graph-based structures are further proposed and explored by recent works. Such models are hard to be learnt exactly because of their high computational cost over the graph structure. Usually, Markov Random Field (MRF) is used to achieve an approximate inference. By taking a branch and bound step in [20], inference on a graph is nearly as efficient as the tree models. Recently, Yang and Ramanan [3] proposed a method that represents each human part by a mixture of templates. Unlike the previous limb models with articulated orientations, their templates are non-oriented and can well capture near-vertical and near-horizontal limbs. By tuning the part type and location, their model can handle the in-plane rotation and foreshortening.

Another drawback of the original PS model is that the contextual information is not explicitly considered. In the work of Sapp et al. [19], various visual cues (e.g., boundary and segmentation) were employed in their coarse-to-fine model. In [13], Rothrock et al. incorporated background context into the PS model. Their model encourages a high contrast of a part region from its surroundings. Experimental results in their paper demonstrated the effectiveness of such contextual information.

2) *Clothing Technology*: The clothing technologies, mainly including clothing segmentation [21]–[23], clothing recommendation [24] and garment attribute classification [4], play an important role in many multimedia systems such as clothing search engines [25], online shopping [26] and human recognition [27].

The work of Chen et al. [28] is one of the representative works in the related literature, which introduced an And-Or graph to generate a large set of composite clothing components for further recognition. In [26], Liu et al. proposed a cross-scenario clothing retrieval system which can search similar garments from the online shop by using a person's daily life photo as input. In their work, a well trained human part detector was employed for part alignment and an offline transfer learning scheme was introduced to handle the discrepancy

between images from daily life and online shops. Recently, they proposed a practical system called “magic closet” in [24] which can automatically recommend garments according to occasions. They utilized the latent SVM algorithm and modeled the clothing attributes as latent variables to provide mid-level features, rather than directly bridging the raw image features to the occasions.

Recent progress in clothing techniques has witnessed the benefit of HPE due to the strong relations between human parts and garment attributes. In other words, it is a popular way to perform clothing study based on the results of human part detection [4], [26]. Moreover, Yamaguchi et al. [23] and Chen et al. [4] used well trained human part detectors to produce a large number of garment types, aiming to achieve precise clothing recognition. Bourdev et al. [29] trained an SVM classifier over each human part with the purpose of indicating whether or not a human part has specific garment attributes.

3) *HPE and CT*: There are few works addressing the interrelations between HPE and CT. Yamaguchi et al. [23] tried to refine both HPE and clothing parsing by using a three-stage scheme: first, they obtained some initial results of HPE; second, they used those initial HPE results as the basis to gain more reliable clothing segmentation; finally, the produced segmentation results were used to further refine HPE. However, since the quality of clothing segmentation largely depends on the success of HPE and vice versa, such a separate modeling approach may fail to capture the correlations between the two tasks and cannot achieve significant improvements over the competing baseline, as can be seen from the reports in [23]. In the recent work of Ladicky et al. [30], they combined the part-based approach of pose estimation and pixel-based approach of image labeling into a principle way so as to inherit advantages of both. Inference for their model was performed with two steps: first, they iteratively added the next best pose candidate by computing an energy function of their model; second, they refined the final solution over the selected candidates of the first step.

B. Contributions of This Work

It is indeed natural to anticipate that HPE and CT are intertwined problems and can help each other. For example, in Fig. 2, depending on the garment type, one can erase a large number of incorrect HPE candidates. However, existing approaches that individually use pose information to refine CT or use clothing knowledge to help HPE cannot fully capture the advantages of modeling the correlations between the two tasks. In this paper, we therefore propose to integrate HPE and garment attribute classification (GAC) into a unified framework, with the purpose of making effective use of the possible correlations between human parts and garment attributes.

Also, we aim to provide an effective way to jointly model multiple visual cues, including the features specific for human parts (pose-specific features) and for garment attributes (garment-specific features), as well as the cross-task features that encode the correlations between human parts and

garment attributes. For a more informative description for oriented human part, we explore the “strong edge” evidence as an energy function so as to incorporate the contextual information around a human part. This motivation is based on the following observation: Since our representation for a human part is an oriented bounding box, it generally holds that there exist parallel edges sharing a similar orientation with an underlying correct part candidate, as illustrated in Fig. 5.

To this end, we use the HPE algorithm presented in [3] to obtain from each image a set of bounding boxes (called “candidates”) that have potentials to be correct human parts, resulting in a basic representation for each image – one image is represented by one set of candidate ensembles. In this way, the joint inference for HPE and GAC is converted to a *structured learning* [31] problem, where the input is the image represented by a collection of candidate ensembles, the output is the joint labels of human parts and garment attributes, and the joint feature representation involves the aforementioned multiple visual cues.

Given a set of annotated training images, the prediction function of structured learning is learnt by using the structured Support Vector Machines (SVM) algorithm. Inference for a new test samples can be performed efficiently by iteratively computing a local optimum on a tree with the dynamic programming algorithm. Experimental results on two benchmark datasets show the state-of-the-art performance of our approach.

One may want to take HPE and GAC into the multi-task learning (MTL) framework [32]. We remark here that our problem cannot be solved via existing MTL methods since models of MTL always assume that the underlying tasks share the same feature space. In our case, however, this assumption is not valid as we address two different tasks: human pose estimation (a detection task) and garment attribute classification (a recognition task). Each of the two tasks has its own feature space that can not be shared with the other, i.e. pose-specific features and garment specific features (see Section II-B). Also, methods from domain adaption (DA) [33] cannot be applied as DA algorithms deal with the variations in some combinations of factors, including scene, object location and pose, view angle, resolution etc. Obviously, our problem is not under the setting of DA algorithms.

The rest of the paper is organized as follows. Section II elaborates on our approach for combined HPE and GAC, including feature design, parameter learning and inference algorithm. Section III presents the experiments and results. Section IV concludes this paper and discusses our future work.

II. JOINT INFERENCE OF HPE AND GAC

As Fig. 1 shows, our approach contains three major procedures, including a preprocessing step that detects candidates from each image, an engineering step that forms a joint feature representation from various visual cues, and an inference step that uses structured SVM learned from a set of training images. In the following, we shall detail each step one by one.

A. Preprocessing and Problem Formulation

We do not build our approach by directly utilizing the images represented by raw pixels, and instead, we use the

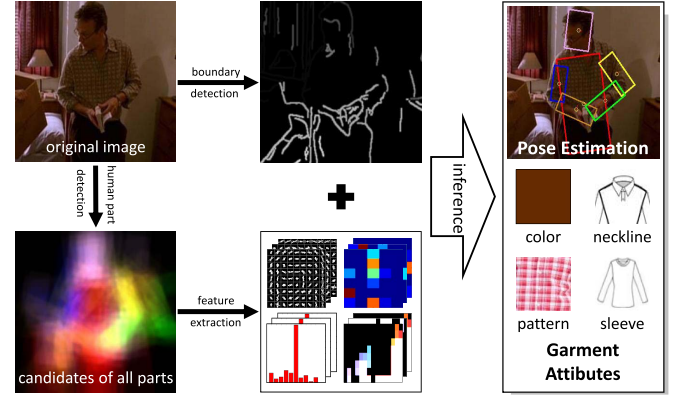


Fig. 1. Overview of the proposed approach. For a given image, first we detect candidates that have potentials to be valid human parts using the HPE algorithm proposed in [3]. This preprocessing step simplifies the representation of the image and converts the joint inference of HPE and GAC into a structured learning problem. Second, we design the joint feature representation for structured learning using various visual cues. Meanwhile, potential strong edges in the image are detected by utilizing a well-established algorithm presented in [5]. Finally, we use the prediction function learned from structured SVM to produce the joint labels of human parts and garment attributes.

existing HPE method [3] to produce some initial results as input to our approach. More precisely, for each human part i (e.g., head), we perform a non-maximum suppression on the output of [3] and take the top K_i (each $K_i = 40$ in our work) candidates (denoted by \mathbf{b}_i) from each image, where each candidate is a bounding box (x, y, θ, s) , with (x, y) , θ and s denoting the center coordinates, the angle and the size of the bounding box, respectively.¹ This step allows us to obtain a manageably sized state space and simplifies the representation of a given instance. Suppose there are m human parts in total ($m = 6$ in this work), and then each image is represented by m candidate ensembles, each of which contains K_i candidates respectively. Thus, the input space (i.e., sample space) \mathcal{X} of our approach is defined as

$$\mathcal{X} = \{\mathbf{x} \mid \mathbf{x} = (\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_m)\}, \quad (1)$$

where \mathbf{x} refers to an image and \mathbf{b}_i denotes the candidate ensemble for the i -th human part (there are K_i candidates in \mathbf{b}_i). Furthermore, we introduce the following notation:

$$\mathcal{P} = \{\mathbf{p} \mid \mathbf{p} = (p_1, p_2, \dots, p_m), \forall i, 1 \leq p_i \leq K_i\}, \quad (2)$$

where p_i is a positive integer that indicates the index of the candidate for the i -th human part. In this way, the task of HPE is formulated as the problem of learning a prediction function from \mathcal{X} to \mathcal{P} .

The goal of GAC in our work is to determine the garment attributes possessed by each image. We consider five types of attributes, including “Collar”, “Color”, “Neckline”, “Pattern” and “Sleeve” that are most relevant with the upper body limbs. Each attribute has multiple styles, e.g., short sleeve and long sleeve for the “Sleeve” attribute. We use T_k to denote the number of attribute values for the k -th attribute. The attribute values we consider in this paper are given in

¹The original output of [3] is a set of non-oriented bounding boxes. We transform them to the oriented ones using the online code that [3] provides.

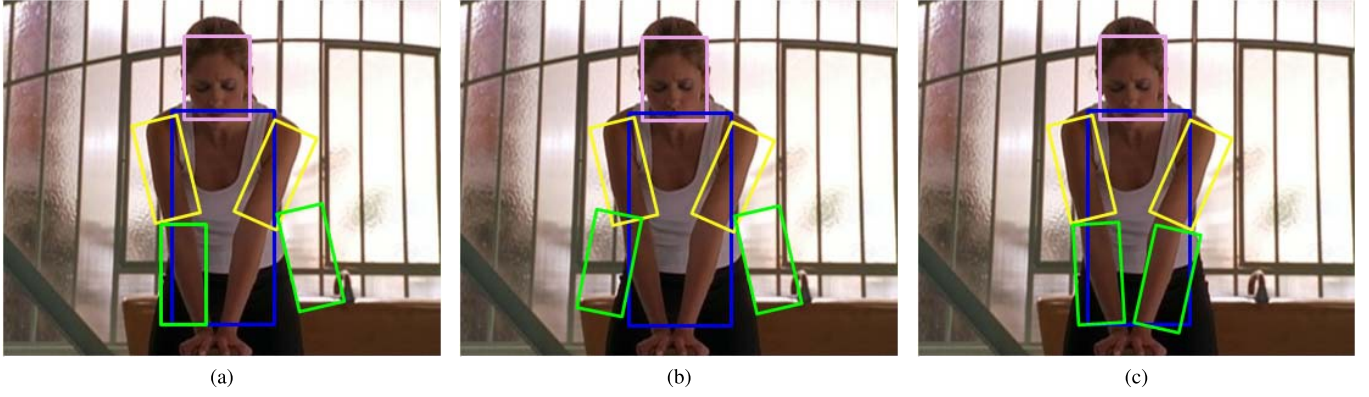


Fig. 2. Examples to illustrate that GAC can help HPE. (a) Result with incorrect left-lower arm. (b) Result with both incorrect lower arms. (c) Result with all correct arms. For (a), one can immediately assert that the HPE result can hardly be correct because the appearances of right-lower arm and left-lower arm differ greatly. Such prior knowledge about limb appearance was considered in [9]. However, it can't distinguish which human estimation is right from (b) and (c) in this way, where the lower arms' appearances differ slightly. Given that the garment attribute is known, e.g., sleeve type is sleeveless, one can easily exclude (b) because the lower arms in (b) have few skin colors.














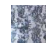
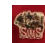






Upper Body Attributes					
Collar					
	none	has			
Color					
	Red	Yellow	Green	Blue	Purple
					
	Brown	Black	White	Gray	Multi-color
Neckline					
	V-shape	Round	Pointed	Strapless	
Pattern					
	Plain	Floral	Graphics	Plaid	Horizontal
					
	Vertical				
Sleeve					
	Sleeveless	Short	Long		

Fig. 3. Garment attributes definition. We list all the garment attributes and types. In this work, we focus on the upper body dressing style. Part of the icons in this figure are quoted from [26].

Fig. 3. For the ease of presentation, we introduce the following notation:

$$\mathcal{C} = \{\mathbf{c} \mid \mathbf{c} = (c_1, c_2, \dots, c_n), \forall k, 1 \leq c_k \leq T_k\}, \quad (3)$$

where n is the number of garment attributes ($n = 5$ in this work), and c_k is the label for the k -th attribute (e.g., $c_5 = 1$ means short sleeve, and $c_5 = 2$ means long sleeve). In this way, similar with HPE, the task of GAC can also be formulated as a problem of learning a prediction function from \mathcal{X} to \mathcal{C} .

Hence, the task of performing combined HPE and GAC can be formulated as follows:

$$f: \mathcal{X} \rightarrow \mathcal{Y}, \quad (4)$$

where \mathcal{Y} is the joint output space defined by

$$\mathcal{Y} = \{\mathbf{y} \mid \mathbf{y} = (\mathbf{p}, \mathbf{c}), \mathbf{p} \in \mathcal{P}, \mathbf{c} \in \mathcal{C}\}. \quad (5)$$

Regarding the prediction function f , we first presume that there is a compatibility function S that measures the fitness between an input-output pair (\mathbf{x}, \mathbf{y}) :

$$S(\mathbf{x}, \mathbf{y}; \mathbf{w}) = \mathbf{w} \cdot J(\mathbf{x}, \mathbf{y}) + \alpha Q(\mathbf{x}, \mathbf{p}), \quad (6)$$

where $\mathbf{w} \cdot J(\mathbf{x}, \mathbf{y})$ denotes the inner product of two vectors, $J(\mathbf{x}, \mathbf{y})$ is the joint feature representation (which should be designed carefully), \mathbf{w} is an unknown weight vector (which should be learned from training samples), $Q(\mathbf{x}, \mathbf{p})$ is the energy function indicating the response of a strong edge around the potential human parts and α is a parameter (which can be hand-tuned by cross-validation).

In this way, the mapping function f in Eq. (4) can be written as:

$$f(\mathbf{x}; \mathbf{w}) = \arg \max_{\mathbf{y} \in \mathcal{Y}} S(\mathbf{x}, \mathbf{y}; \mathbf{w}). \quad (7)$$

B. Joint Feature Representation

The joint feature representation $J(\mathbf{x}, \mathbf{y})$ is an important component of the prediction function. In our approach, $J(\mathbf{x}, \mathbf{y})$ consists of three types of features, including the *pose-specific* features denoted by $J_p(\mathbf{x}, \mathbf{p})$, the *garment-specific* features denoted by $J_c(\mathbf{c})$, and the *cross-task* features denoted by $J_{pc}(\mathbf{x}, \mathbf{y})$; that is,

$$\mathbf{w} \cdot J(\mathbf{x}, \mathbf{y}) = \mathbf{w}_p \cdot J_p(\mathbf{x}, \mathbf{p}) + \mathbf{w}_c \cdot J_c(\mathbf{c}) + \mathbf{w}_{pc} \cdot J_{pc}(\mathbf{x}, \mathbf{y}). \quad (8)$$

In the following, we shall present our techniques used to design each type of feature.

1) *Pose-Specific Features*: Given an image represented as m candidate ensembles (each candidate is a bounding box), we extract some features specifically useful for HPE, called pose-specific features, as follows:

$$\begin{aligned} \mathbf{w}_p \cdot J_p(\mathbf{x}, \mathbf{p}) = & \sum_{i=1}^m \mathbf{w}_p^{u,i} \cdot \phi_p(\mathbf{x}, p_i) \\ & + \sum_{(i,j) \in E_p^d} \mathbf{w}_p^{d,ij} \cdot \psi_p^d(\mathbf{x}, p_i, p_j) \\ & + \sum_{(i,j) \in E_p^c} \mathbf{w}_p^{c,ij} \cdot \psi_p^c(\mathbf{x}, p_i, p_j), \end{aligned} \quad (9)$$

where $\phi_p(\mathbf{x}, p_i)$ denotes the unary feature for the i -th human part, $\psi_p^d(\mathbf{x}, p_i, p_j)$ models the pairwise relations between spatially connected human parts, E_p^d is the collection of all

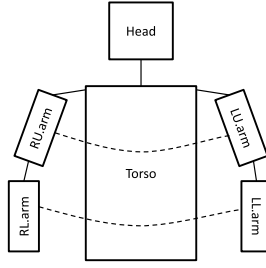


Fig. 4. Part-part relations. The solid lines are used to indicate spatially connected parts, e.g., torso and right upper arm (RU.arm). We mainly consider the geometry constraint for this case, e.g., relative position and relative rotation. The dashed lines mean the relations between symmetric parts, e.g., right upper arm (RU.arm) and left upper arm (LU.arm). For these parts, we model the appearance constraint, i.e., appearance similarity in color and texture descriptors.

pairs of connected parts, $\psi_p^c(\mathbf{x}, p_i, p_j)$ contains the appearance consistency message between symmetric parts and E_p^c is the set of all symmetric parts (see Fig. 4 for details). The three terms are called *unary score*, *deformation score* and *consistency score* respectively.

In our approach, the unary feature $\phi_p(\mathbf{x}, p_i)$ is chosen as the HOG descriptor [16] which has been proved quite effective for object detection [3], [10].

The design of $\psi_p^d(\mathbf{x}, p_i, p_j)$ concerns some basic geometry constraints between connected parts, including relative position, rotation and distance of part candidate p_i with respect to p_j . More concretely, we divide the image space into 3 by 3 regions, with p_j at the central region. Then we use a 9 dimensional one-zero vector as the relative position feature, where there is only one element with value “1” that indicates the region where p_i is located. To describe the relative rotation, we divide the range of angles $[0, 360]$ into 20 bins and use a 20 dimensional one-zero vector as the feature. Our relative distance feature is the euclidean distance between the center of p_i and p_j .

For the symmetric parts, we assume some consistency constraints between them which hold for most cases; that is, they should share similar appearance. In this work, we compute the divergence of the color histogram in RGB and LAB space and take it as the feature descriptors.

In Fig. 4, we mark all the parts that are related to each other.

2) *Garment-Specific Features*: There are some features only specific for garment, i.e., garment-specific features. In this work, we consider the co-occurrences between different garment attribute values:

$$\mathbf{w}_c \cdot J_c(\mathbf{c}) = \sum_{k,l} \mathbf{w}_c^{kl} \cdot \psi_c(c_k, c_l), \quad (10)$$

where $\psi_c(c_k, c_l)$ is a binary vector that indicates whether or not c_k and c_l co-occur in an image. For example, the texture type “drawing” (usually belongs to T-shirt style) often co-occurs with the collar type “round”.

3) *Cross-Task Features*: The cross-task features encode the correlations between human parts and garment attributes. In our approach, we model the part-garment relations manually specified as in Table I. For a given attribute k , we denote the

TABLE I
THE INTERDEPENDENCY BETWEEN HUMAN PARTS AND GARMENT ATTRIBUTES. FIRST COLUMN: A GARMENT ATTRIBUTE; SECOND COLUMN: THE HUMAN PARTS ASSOCIATING WITH THE GARMENT ATTRIBUTE SHOWN IN THE FIRST COLUMN; THIRD COLUMN: THE CORRESPONDING FEATURE DESCRIPTORS USED TO DESCRIBE THE PARTS (OR ATTRIBUTE)

Attribute	Human Parts	Low-level Features
Collar	Torso+Head	HOG
Color	Torso	Color Histogram
Neckline	Torso+Head	HOG
Pattern	Torso	LBP [34]
Sleeve	All arms	Color Histogram

human part(s) associated with it as $\hat{\mathbf{p}}(k)$ and the corresponding configuration(s) as $\hat{\mathbf{p}}_k$. Then the cross-task features are formulated as:

$$\mathbf{w}_{pc} \cdot J_{pc}(\mathbf{x}, \mathbf{y}) = \sum_{k=1}^n \mathbf{w}_{pc}^k \cdot \Psi_{pc}^k(\mathbf{x}, \hat{\mathbf{p}}_k, c_k), \quad (11)$$

where $\Psi_{pc}^k(\mathbf{x}, \hat{\mathbf{p}}_k, c_k)$ denotes the features extracted from \mathbf{x} under the constraints of part configuration $\hat{\mathbf{p}}_k$ and attribute label c_k . Note that here we write the cross-task score as a summary by the attribute order. Since the dependency between part and attribute is cyclic, one can also write it by the human part order.

To describe the design of $\Psi_{pc}^k(\mathbf{x}, \hat{\mathbf{p}}_k, c_k)$, we first convert the garment attribute label c_k to a T_k dimension vector $\mathbf{I}(c_k)$, with only one dimension assigned with value one and others with zeros. From Table I, the low-level feature descriptors of the k -th garment attribute depend on two aspects: the corresponding human part(s) and the feature type (denoted by F_k and specified in Table I). We use $F_k(\hat{\mathbf{p}}_k)$ to denote features of the k -th garment attribute under the part candidate(s) $\hat{\mathbf{p}}_k$. Then our cross-task feature $\Psi_{pc}^k(\mathbf{x}, \hat{\mathbf{p}}_k, c_k)$ is represented as follows:

$$\Psi_{pc}^k(\mathbf{x}, \hat{\mathbf{p}}_k, c_k) = F_k(\hat{\mathbf{p}}_k) \otimes \mathbf{I}(c_k), \quad (12)$$

where the “ \otimes ” operator is the outer product of two vectors. In fact, we map the resulting matrix to a vector by the row order. Note that in Table I, a garment attribute depends exclusively on the some of the limbs, not all ones. This technique that feature descriptors draw from both the labels of human parts and garment attributes, provides us a simple way to capture the correlations between HPE and GAC and makes it a unified approach towards the two intertwined problems.

C. Learning With Structured SVM

We perform our joint estimation for HPE and GAC using the prediction function f in Eq. (7). The weight vector \mathbf{w} is a critical component of the prediction function. Given N training samples $\{(\mathbf{x}_r, \mathbf{y}_r)\}_{r=1}^N$, we compute \mathbf{w} by solving the

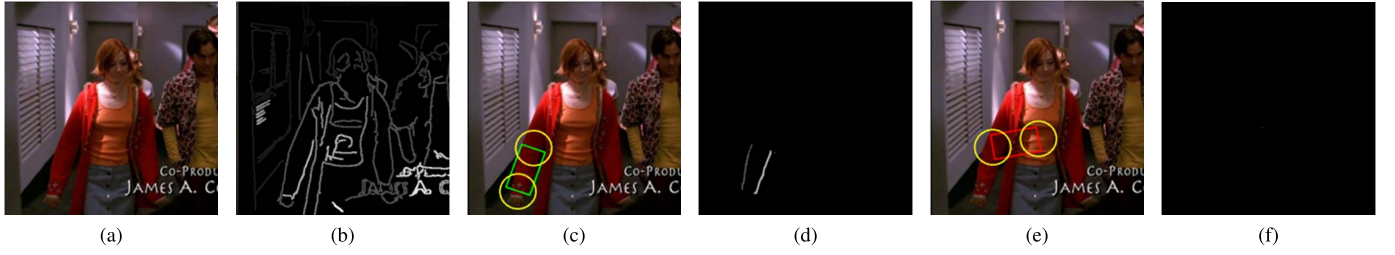


Fig. 5. Illustration for Strong Edge Evidence. (a) original image; (b) boundary detection results; (c) a correct candidate for the right-lower arm (green rectangle); (d) strong edge evidence for the correct candidate; (e) an incorrect prediction for the same arm (red rectangle); (f) strong edge evidence for the incorrect candidate. Given a part candidate (denoted by an oriented rectangle), we extract the strong edges that connect the regions of the joints of the part (denoted by two yellow rounds).

following structured SVM problem:

$$\begin{aligned} & \underset{\mathbf{w}, \xi}{\text{minimize}} && \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{r=1}^N \xi_r, \\ & \text{subject to} && \forall r \in \text{pos}, \mathbf{w} \cdot J(\mathbf{x}_r, \mathbf{y}_r) \geq 1 - \xi_r, \\ & && \forall r \in \text{neg}, \mathbf{w} \cdot J(\mathbf{x}_r, \mathbf{y}_r) \leq -1 + \xi_r. \end{aligned} \quad (13)$$

where C is the parameter that controls the trade-off between margin and accuracy, and $\xi_r \geq 0$ is a slack variable. In our experiments, we set C with a fixed value 0.01 that allows a soft margin.

D. Strong Edge Evidence

Now we give details on the design of our energy function $Q(\mathbf{x}, \mathbf{p})$ in Eq. (6). First, we utilize a boundary detector [5] to detect all potential edges in an image. Then for a candidate human part, we try to find the long edges that connect the two joint regions as the strong edge evidence (see Fig. 5). Our energy function mainly considers three factors: the consistency of orientation between the part p_i and the strong edges \mathbf{se}_i (denoted by $Q^o(p_i, \mathbf{se}_i)$), the distance of the part away from the strong edges (denoted by $Q^d(p_i, \mathbf{se}_i)$) and the strength of the strong edges themselves. That is,

$$Q(\mathbf{x}, \mathbf{p}) = \sum_{i=1}^m Q^o(p_i, \mathbf{se}_i) + \beta \sum_{i=1}^m Q^d(p_i, \mathbf{se}_i), \quad (14)$$

where β is a parameter which can be tuned by cross validation. Given a part candidate $p_i = (x, y, \theta, s)$, the first term in Eq. (14) is computed as:

$$Q^o(p_i, \mathbf{se}_i) = \frac{1}{Z} \sum_{e \in \mathbf{se}_i} \cos(\theta - \theta_e) \cdot \text{strge}_e, \quad (15)$$

where e is an image pixel on the strong edges \mathbf{se}_i , Z is the number of pixels on \mathbf{se}_i , θ_e is the orientation of the strong edge at pixel e and strge_e is the edge strength (which is produced by the algorithm [5]). The measurement of the distance from the given part to the strong edges is represented as follows:

$$Q^d(p_i, \mathbf{se}_i) = \frac{1}{Z} \sum_{e \in \mathbf{se}_i} \min\{\text{dt}(e, l), \text{dt}(e, r)\} \cdot \text{strge}_e, \quad (16)$$

where l and r are the two parallel edges of the part bounding box whose angles are θ , and $\text{dt}(e, l)$ is the closest distance of pixel e from edge l which can be efficiently computed by

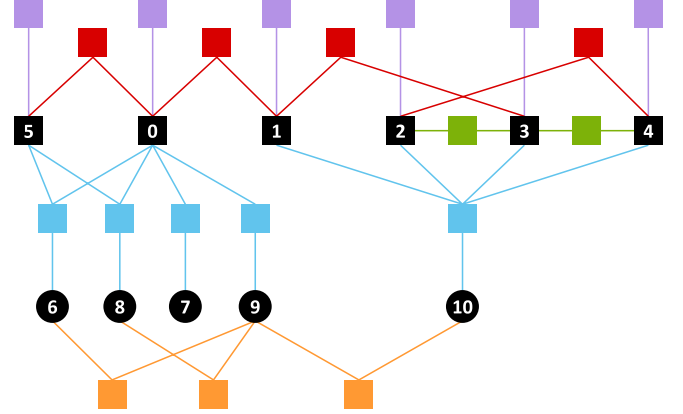


Fig. 6. The factor graph representation for our problem. We denote our variables with black nodes, those of which with number 0–5 represent the human parts: torso, RU/LU/RL/LL arm and head, while those with number 6–10 denote the garment attributes: collar, color, neckline, pattern and sleeve. We denote our potentials with colored nodes, with purple ones denoting the unary potential, red denoting the deformation potential, green denoting the consistency potential, orange denoting the attribute co-occurrence potential and cyan denoting the cross potential.

a distance transform algorithm [35]. Note that if there is no strong edge associated with the underlying part, we force the energy to be zero.

E. Inference

Now that we have clarified how to design the joint feature, the strong edge energy function, as well as the learning algorithm for weight vector \mathbf{w} in Eq. (6), we propose our inference algorithm which is quite efficient (for each input sample, our algorithm only needs 2 seconds for the joint estimation) and effective.

In Fig. 6, we represent our problem as a factor graph \mathcal{G} , where the black-rectangle node denotes a human part, the black-circle node denotes a garment attribute and the colored node denotes a potential. As our original problem is a cyclic graph, it cannot be optimized exactly and efficiently. Therefore, in Algorithm 1, we propose an iterative algorithm to search for an approximate solution. Our algorithm receives a sample $@$ (defined in Eq. (1)), the SVM weight \mathbf{w} , parameter α and β as inputs and outputs the optima for the joint problem. In each iteration, by fixing one type of the variable (either human part or garment attribute, see step 5 and step 6),

Algorithm 1 Approximate Inference for Joint Estimation

Input: An input sample \mathbf{x} , the weight vector \mathbf{w} , parameter α and β .

Output: Optimal joint estimation \mathbf{y}^* .

- 1: Set the optimal joint estimation $\mathbf{y}^* = \emptyset$.
- 2: Set the optimal score $S^* = -\infty$.
- 3: Initialize the parts estimation:
 $\mathbf{p}_0 = \arg \max_{\mathbf{p} \in \mathcal{P}} \mathbf{w}_p \cdot J_p(\mathbf{x}, \mathbf{p}) + \alpha Q(\mathbf{x}, \mathbf{p})$.
- 4: **repeat**
- 5: Compute the garment attributes:
 $\mathbf{c}_t = \arg \max_{\mathbf{c} \in \mathcal{C}} \mathbf{w}_c \cdot J_c(\mathbf{c}) + \mathbf{w}_{pc} \cdot J_{pc}(\mathbf{x}, \mathbf{p}_{t-1}, \mathbf{c})$.
- 6: Compute the parts estimation:
 $\mathbf{p}_t = \arg \max_{\mathbf{p} \in \mathcal{P}} \mathbf{w}_p \cdot J_p(\mathbf{x}, \mathbf{p}) + \alpha Q(\mathbf{x}, \mathbf{p}) + \mathbf{w}_{pc} \cdot J_{pc}(\mathbf{x}, \mathbf{p}, \mathbf{c}_t)$.
- 7: Compute the local score: $S = S(\mathbf{x}, \mathbf{y}_t; \mathbf{w})$.
- 8: **if** $S > S^*$ **then**
- 9: $S^* = S$, $\mathbf{y}^* = \mathbf{y}_t$.
- 10: **end if**
- 11: **until** S^* not change

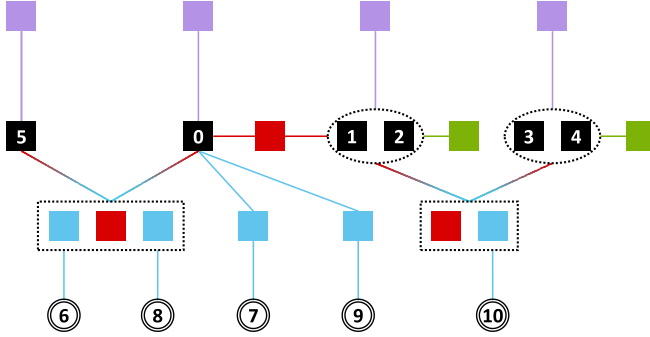


Fig. 7. The factor graph representation for inferring human pose. Circle nodes with double boundaries are assigned with fixed values. Symmetric parts are grouped into a super-node, denoted by a dashed oval. For some part nodes, their deformation and cross potential can also be grouped as the associated attribute nodes are now fixed. Note that we don't draw the garment-specific potentials as they don't contribute for searching a best pose.

our inference procedure can be performed on a tree structure which yields an efficient computation by dynamic programming [7]. This procedure is also illustrated in Fig. 7 and Fig. 8.

1) *Inference for Pose:* In the work of [7], the PS model is restricted as a tree: each node has a unary term that describes how suitable a configuration is assigned to this part, and each edge encodes the deformation cost for a pair of connected parts. In Fig. 6, we demonstrate our extension for the traditional PS model:

- appearance consistency between symmetric parts (green nodes)
- joint compatibility across the human part(s) and the garment attribute(s) (blue nodes)

Adding the edges connecting symmetric parts will destroy the tree structure. In Fig. 7, however, we propose a trick to group the symmetric parts as a *super-node* so that the global structure remains to be a “tree”. On the other hand, an edge across human part and garment attribute is used to

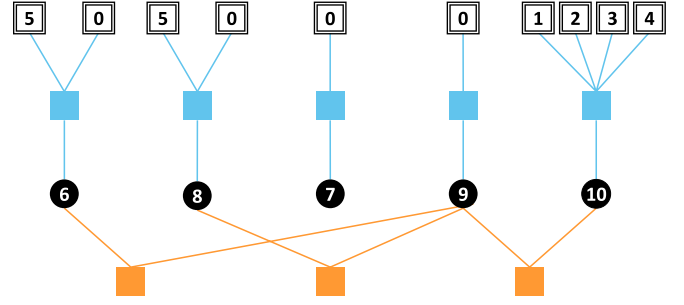


Fig. 8. The factor graph representation for inferring attributes. Part nodes are with fixed values and stretching them will not affect the optimal solution. Unary, deformation and consistency potentials of part nodes are not drawn here, as they don't contribute to searching for a best attribute solution.

measure how compatible a human part configuration is with a given attribute. We call this kind of cost as *cross score*. When the attribute variables are fixed, we can remove the garment-specific potentials as they do not contribute to searching for the best pose. In addition, we can group some deformation and cross potentials for a more concise representation, e.g. as what we do for node 1 and 2 in Fig. 7.

In Algorithm 2, we describe our computation procedure. For a super-node \mathbf{i} ,² we denote its children nodes as C_i . In the line 3–13, we first compute the scores with respect to a single node \mathbf{i} . This step involves calculation of unary score, strong edge score, consistency score and cross score. Note that we have grouped the symmetric parts as one node. In this way, the consistency score is a self description towards the node \mathbf{i} . In line 15, we compute the deformation score of node \mathbf{i} and \mathbf{j} . For example, if the super-node $\mathbf{i} = \{1, 2\}$ whereas the super-node $\mathbf{j} = \{0\}$, the deformation score between \mathbf{i} and \mathbf{j} is the sum of deformation score of $(1, 0)$ and $(2, 0)$. In line 16, we compute the cross score for all attributes whose associated human parts are exactly $\mathbf{i} \cup \mathbf{j}$. For example, the attribute node 6 is associated with part node 5 and 0, so we will compute the cross score with respect to nodes $\{5, 0, 6\}$. Line 18–27 is a conventional message passing procedure that can be computed efficiently by dynamic programming [7].

2) *Inference for Attributes:* Referring to Fig. 8, we stretch the part variables and remove redundant edges associated with the stretched variables from the original graph as they contribute nothing to this inference step. Note that for the attribute-attribute pairs (i.e. the garment-specific feature), we manually model them as a tree structure. In this way, we can still perform an efficient computation like Algorithm 2.

3) *Implementation and Computation Complexity:* We write $K = \max_{1 \leq i \leq m} K_i$ and $T = \max_{1 \leq k \leq n} T_k$. Now we propose a computation analysis for our Algorithm 2 and give some optimization tricks. In line 3–13, one has to loop over all possible configuration \mathbf{p}_i for the super-node \mathbf{i} , and there are at most 2 nodes in a super-node (see Fig. 7); this makes the computation $O(K^2)$. In fact, note that the computation of unary and strong edge score can be *decomposed* into a summation of each node (line 4 and 5), which indicates that we can separately

²For simplicity, here we call all variable nodes as super-nodes.

Algorithm 2 Exact Inference for Human Pose Estimation (Extended Pictorial Structure Inference)

Input: An input sample \mathbf{x} , the weight vector \mathbf{w} , parameter α , β and garment attributes \mathbf{c} .

Output: Optimal pose estimation \mathbf{p}^* .

```

1: Set the optimal joint estimation  $\mathbf{p}^* = \emptyset$ .
2: Set the node 0 as the root node.
3: for each configuration  $\mathbf{p}_i$  of super-node  $\mathbf{i}$  do
4:    $m_1 = \sum_{i \in \mathbf{i}} \mathbf{w}_p^{u,i} \cdot \phi_p(\mathbf{x}, p_i)$ .
5:    $m_2 = \alpha \sum_{i \in \mathbf{i}} Q(\mathbf{x}, p_i)$ .
6:   if  $\mathbf{i} \in E_p^c$  then
7:      $m_3 = \mathbf{w}_p^{s,i} \cdot \psi_p^s(\mathbf{x}, \mathbf{p}_i)$ ,
8:   else
9:      $m_3 = 0$ .
10:  end if
11:   $m_4 = \sum_{k, \hat{\mathbf{p}}(k)=\mathbf{i}} \mathbf{w}_{pc}^k \cdot \Psi_{pc}^k(\mathbf{x}, \hat{\mathbf{p}}_k, c_k)$ .
12:  set  $m(\mathbf{p}_i) = m_1 + m_2 + m_3 + m_4$ .
13: end for
14: for each configuration of parent-child pair  $\mathbf{p}_j$  and  $\mathbf{p}_i$  do
15:    $l_1 = \sum_{i \in \mathbf{i}, j \in \mathbf{j}, (i,j) \in E_p^d} \mathbf{w}_p^{d,ij} \cdot \psi_p^d(\mathbf{x}, p_i, p_j)$ .
16:    $l_2 = \sum_{k, \hat{\mathbf{p}}(k)=\mathbf{i} \cup \mathbf{j}} \mathbf{w}_{pc}^k \cdot \Psi_{pc}^k(\mathbf{x}, \hat{\mathbf{p}}_k, c_k)$ .
17:   set  $l(\mathbf{p}_i, \mathbf{p}_j) = l_1 + l_2$ .
18:   if  $\mathbf{i}$  is a leaf node then
19:      $B_i(\mathbf{p}_j) = \max_{\mathbf{p}_i} (m(\mathbf{p}_i) + l(\mathbf{p}_i, \mathbf{p}_j))$ ,
20:   else
21:      $B_i(\mathbf{p}_j) = \max_{\mathbf{p}_i} (m(\mathbf{p}_i) + l(\mathbf{p}_i, \mathbf{p}_j) + \sum_{\mathbf{v} \in C_i} B_v(\mathbf{p}_i))$ .
22:   end if
23: end for
24: Compute the best configuration for the root node:
    $\mathbf{p}_0^* = \arg \max_{\mathbf{p}_0} (m(\mathbf{p}_0) + \sum_{\mathbf{v} \in C_0} B_v(\mathbf{p}_0))$ .
25: for each parent-child pair  $(\mathbf{p}_j^*, \mathbf{p}_i)$  do
26:    $\mathbf{p}_i^* = \arg \max_{\mathbf{p}_i} B_i(\mathbf{p}_j^*)$ .
27: end for

```

compute these scores for each node configuration, yielding a computation $O(K)$. Also note that actually we only compute cross score for node 0 in line 11 (see Fig. 6). Based on this observation, computation on $m(\mathbf{p}_i)$ is reduced from $O(K^2) + O(K^2) + O(K^2) + O(K^2)$ to $O(K) + O(K) + O(K^2) + O(K)$. In line 14–16, we compute the deformation and cross score for each pair $(\mathbf{p}_i, \mathbf{p}_j)$, which yields a computation complexity $O(K^4)$ if without any optimization. For the deformation score, as the decomposition property still holds, the computation is $O(K^2)$. For the cross score in line 16, when $k \in \{6, 8\}$, the computation is $O(K^2)$ since we have to loop over all the configurations for nodes 0 and 5. When $k = 10$, however, it is not necessary to enumerate the $O(K^4)$ combinations of node 1, 2, 3 and 4 if we design a suitable cross-task feature. In our case, $F_{10}(\hat{\mathbf{p}}_{10})$ is the concatenation of the color histogram of each $p_i \in \hat{\mathbf{p}}_{10}$, which implicitly owns the decomposition property. Thus, the computation can be reduced to $O(K)$, if we omit the summation operation of these separate scores.

F. Parameter Sharing

Our work is distinct from other works which address pose estimation and garment attribute in two aspects. First, our

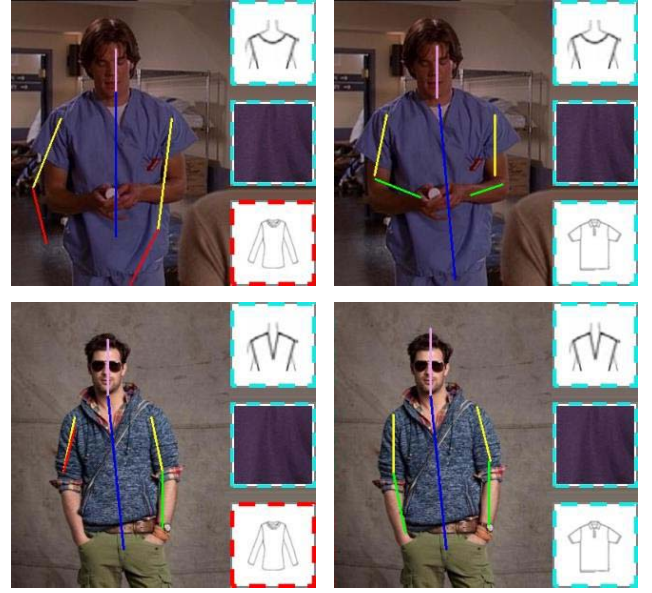


Fig. 9. Our joint approach v.s. YDR [3]. Left column: results from [3]. Right column: results from our joint approach. [3] estimates incorrect lower arm(s), which subsequently results in an incorrect prediction for the sleeve attribute. Our joint approach captures the co-relations between the arms and the sleeve attribute and makes a correct estimation for *both*.

carefully designed structured learning model integrates the pose feature and garment attributes into a principle fashion, facilitating a global optimal model. Second, although we derive an iterative inference algorithm to approximate the optima, we allow the *parameter sharing* between the two steps, i.e. the cross-task features are shared and contribute to both pose estimation and attribute classification (line 5 and 6 in Algorithm 1). Therefore, our approach is a paradigm of learning globally and inferring locally, achieving both effectiveness and efficiency (see Section III for experimental justification).

III. EXPERIMENTS

A. Experimental Settings

In this section, we introduce our experimental settings, including the used datasets, the baselines, the evaluation metrics and the scheme for training structured SVM and inferring for a testing image.

1) *Datasets*: We conduct experiments on two datasets. The first one is the widely used Buffy dataset [11] consisting of 748 annotated video frames from Buffy TV show. This dataset is proposed as a standard one for HPE task but not originally for GAC task. We manually annotate the garment attributes for the Buffy dataset. The second dataset, called “DL”, contains 1000 daily life photos we collect from websites. Compared with Buffy, the DL dataset possesses more various garment attribute values. In order to obtain quantitative evaluation results, we also manually annotate the human parts and garment attributes for the images in the DL dataset.

Some garment attributes cannot be labeled for the images in which the person does not wear any garment or some attributes’ visual cues cannot be described. In Fig. 11, we illustrate some of such samples and list

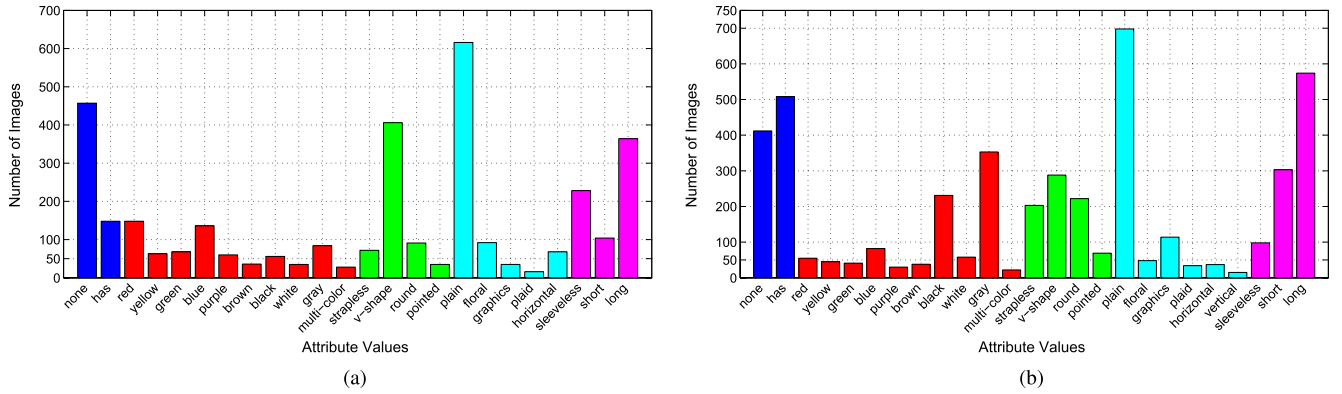


Fig. 10. Statistical information for the clothing attribute annotation. Left: Buffy dataset. Right: DL dataset. We use bars with different colors to denote different attributes. Orders from left to right defer to the orders in Fig. 3. (a) Buffy. (b) DL.



Fig. 11. Examples lacking some garment attributes. Characters in these images all lack visual cues for specific garment attributes. Persons of the first row wear no garment, thus cannot be labeled with any garment attribute while those of the second row cannot be labeled with part of attributes, e.g., neckline types.

the statistical information for the attributes we annotate in Fig. 10.

2) *Baselines*: We select four state-of-the-art HPE methods as our baselines: Andriluka et al. (ARS) [8], Sapp et al. (STT) [19], Yang and Ramanan (YDR) [3] and Ladicky et al. (LTZ) [30]. As the code of LTZ is not publicly available, we only evaluate our DL dataset by the first three methods.

Although all these algorithms are primarily designed for HPE, they can actually produce results for GAC: as discussed in [4] and [26], the results of HPE can be used for part alignment, which enables the extraction of attribute features. Then for each attribute, we individually train an SVM multi-class model [36]. We use the features described in Table I to train each SVM model. We also compare our GAC results with CGG [4],³ which designed a specific pipeline to recognize semantic clothing attributes.

³This code is not publicly available. We thank the authors Chen et al. [4] for providing us the source code for performance comparison.

3) *Evaluation Metrics*: We evaluate the HPE results with the standard metric of Probability of Correct Part (PCP) [9]. The GAC results are evaluated by the Garment Attribute Precision (GAP) criterion, i.e., the classification accuracy for each garment attribute (there are 5 attributes in this work).

4) *Training/Testing*: For the Buffy dataset, like [3], [11] and [37], we select the images from Episode 3, 4 for training, and Episode 2, 5 and 6 for testing. For our DL dataset, we select randomly 300 images for training and use the remaining 700 images for testing.

As we have discussed in Section III-A1, some images cannot be annotated with some garment attributes. For an image without the garment attribute c_j , we set all the features related to c_j to a zero vector in Eq. (8) when training our structured SVM model and skip evaluation on such attributes. For a person wearing two (or more) garments, we label each garment's attribute values. Thus the image has several groups of labels in terms of these garments. When training the model, all the groups of labels are used to construct the constraints. When testing a new instance, any attribute value which the algorithm produces is acceptable if the value belongs to any of the groups.

B. Results

Fig. 12 shows some exemplar results produced by our approach. In the following, we shall analyze our approach and compare it with the baselines.

1) *Examining the Advantages of Joint Learning*: To show the advantages of combining HPE and GAC together, we compare our joint learning approach with its two separated versions: one is an HPE algorithm created by removing the garment-related features in Eq. (8); the other is a GAC algorithm created by ignoring all part-related features. Fig. 13 shows the comparison results, which demonstrate the significant advantages of the joint learning over the separated schemes.

Note that the basic appearance constraints of human parts (as [9] considered) have been modeled in the pose-specific features (see Section II-B). Thus, the improvement on HPE is attributed to the integration of garment attribute evidence modeled in Eq. (11), i.e., joint inference. The improvement on GAC can be examined in the same way.



Fig. 12. Examples of our results obtained on the Buffy and DL datasets. We demonstrate some good results from Buffy and DL in the first and second panels respectively. Some failure cases are shown in the bottom panel. We use the oriented line to denote the pose estimation. If an HPE result is incorrect, the line is red. We visualize three attributes (neckline, pattern and sleeve) of our GAC results by some icons (see Fig. 3 for the icon definition). If a GAC result is incorrect, we use a dashed red rectangle to mark it. Examining the failure cases, we find our algorithm is confused when some human parts are occluded or the human pose is largely variational. Attributes are misclassified when the corresponding parts are mis-detected, or occluded by some objects.

2) *Examining the Effectiveness of Strong Edge*: We demonstrate the effectiveness of the strong edge evidence in this section. By setting the α in Eq. (6) with value zero, our

inference algorithm 2 produces the results without strong edge evidence. Then we use 3-fold cross validation to tune the parameters α and β in Eq. (14). The error reduction rate

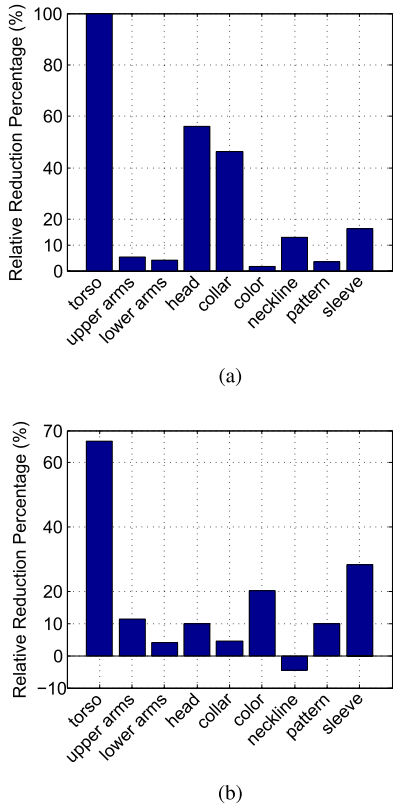


Fig. 13. Demonstrating the advantages of combining HPE and GAC together. X-axis: human part and garment attribute. Y-axis: error reduction rate. We compare our joint learning approach with its separated versions on the Buffy and DL dataset, which deal with HPE and GAC individually. The joint approach improves all of the human parts and a majority of garment attributes. This is because there exists a strong correlation between human parts and garment attributes, and our algorithm captures their inter-dependency that improves both simultaneously. (a) Buffy. (b) DL.

for employing the strong edge evidence on two datasets is demonstrated in Fig. 14. For the Buffy dataset, by using the strong edge evidence, the lower arm accuracy is refined and for the DL dataset, *all* of the human parts are predicted more precisely.

3) *Comparisons With State-of-the-Art Algorithms*: In this section, we compare our joint approach (with strong edge evidence) with the state-of-the-art algorithm. Note that the results of GAC produced by HPE algorithms have been explained in Section III-A2. Fig. 9 gives the exemplar comparison of HPE and GAC results from YDR [3] and our joint approach. On the Buffy dataset, Table II shows that our approach consistently outperforms YDR [3] which is a recently established algorithm and provides the candidates for our approach. We also compare our approach with LTZ [30] which combines pose estimation and segmentation for computation. We improve the lower arms performance and achieve the highest overall accuracy. Table III shows the comparison results on the DL dataset. It can be seen that our approach outperforms all the competing baselines on the task of HPE.

To examine the effectiveness of our approach on GAC, we also compare it with CGG [4], which is a real GAC method. On the Buffy dataset, there is a significant

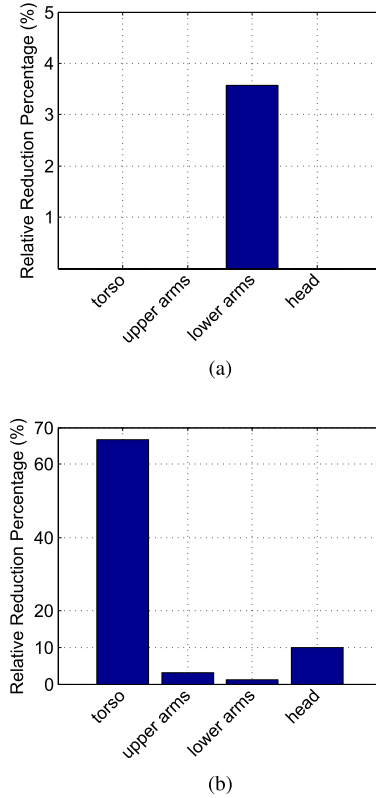


Fig. 14. Demonstrating the effectiveness of strong edge evidence. X-axis: human part. Y-axis: error reduction rate. We demonstrate the error reduction rate for employing the strong edge evidence on both Buffy and DL dataset. The strong edge evidence improves the detection rate because it captures the contextual information for a human part, which is complementary to other features. (a) Buffy. (b) DL.

improvement of our approach on the attributes of “color” and “sleeve”, and on the DL dataset, we also reach a competitive performance. The reason of the surprising gain on the Buffy dataset is that human pose of Buffy is more variational than DL. And our model can capture the inter-dependency between human part(s) and garment attributes. However, the pipeline of CGG is step by step, like the work of [23]. Thus it is depressed if the human pose is largely unconstrained.

One may notice that in Table II and Table III, compared with the baselines except CGG [4], the GAC results of our approach are significantly improved when we gain less improvement for HPE on average (see YDR [3] for example). The reasons are twofold. First, the training procedure of our approach is different from that of theirs. Our model is trained in a unified manner, which allows us to integrate more useful features. That is, during the training procedure, our model captures the dependency between human parts and garment attributes (i.e. cross-task features), as well as that between different garment attributes (i.e. garment-specific features). As a result, we finally have a global optimal model for HPE and GAC. However, the competing baselines are trained in a separate manner. That is, given the HPE, each garment attribute is trained individually. Therefore, neither the inter-dependency between human parts and garment attributes, nor that between different garment attributes can be utilized, resulting in a marginal model for multiclass SVM. Second, we search for

TABLE II
COMPARISON WITH STATE-OF-THE-ART ALGORITHMS ON THE BUFFY DATASET

Method	Torso	U. arms	L. arms	Head	Total
ARS [8]	90.7	79.3	41.2	95.5	73.5
STT [19]	100	95.3	63.0	96.2	85.5
YDR [3]	100	96.6	70.9	99.6	89.1
LTZ [30]	100	97.5	75.4	100	90.9
CGG [4]	—	—	—	—	—
Our Approach	100	96.4	78.4	98.9	91.4

Collar	Color	Neckline	Pattern	Sleeve	Total
70.3	71.7	68.7	80.9	46.6	67.7
77.8	71.2	73.4	80.1	49.2	70.3
82.8	70.8	68.3	80.9	51.5	70.9
—	—	—	—	—	—
89.1	58.4	69.4	80.9	46.2	68.8
88.3	73.1	76.1	81.6	61.7	76.2

TABLE III
COMPARISON WITH STATE-OF-THE-ART ALGORITHMS ON THE DL DATASET

Method	Torso	U. arms	L. arms	Head	Total
ARS [8]	89.4	80.3	60.6	85.0	76.0
STT [19]	99.9	91.1	69.2	97.0	86.2
YDR [3]	99.9	96.0	82.2	99.0	92.5
CGG [4]	—	—	—	—	—
Our Approach	99.9	96.9	83.4	99.1	93.3

Collar	Color	Neckline	Pattern	Sleeve	Total
70.0	55.8	50.9	77.9	60.7	63.1
55.5	58.1	35.9	77.7	61.9	57.8
75.0	58.6	60.0	77.7	64.1	67.1
78.1	69.5	59.2	78.7	68.4	70.8
78.5	67.1	60.0	79.9	68.2	70.7

the optimal prediction for GAC by iteratively updating HPE and GAC, reaching a (local) optimal state of HPE and GAC.⁴ However, the baselines can only make a prediction for GAC by the given HPE result.

IV. CONCLUSIONS AND FUTURE WORK

Based on the observation that there exist correlations between human parts and garment attributes, we propose to integrate HPE and GAC into a unified procedure and handle both tasks simultaneously. We show that such integration can be seamlessly achieved by using the framework of structured SVM. First, due to the joint feature representation, it is convenient to involve various visual cues such as pose-specific features, garment-specific features and cross-task features. Second, the structured nature of the output space of structured SVM provides us a straightforward way to jointly infer the solutions for several problems (e.g., HPE and GAC). Benefiting from these superiorities, our approach achieves state-of-the-art performance in both HPE and GAC problems, as demonstrated in the experiments. Obviously, the boosted performance can benefit quite many multimedia applications, e.g., online clothing retrieval, clothing recommendation, and we are planning to extend our proposed approach for these applications in our future work.

REFERENCES

- [1] X. Xiong and F. De la Torre, "Supervised descent method and its applications to face alignment," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2013, pp. 532–539.
- [2] M. Andriluka, S. Roth, and B. Schiele, "People-tracking-by-detection and people-detection-by-tracking," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2008, pp. 1–8.
- [3] Y. Yang and D. Ramanan, "Articulated pose estimation with flexible mixtures-of-parts," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2011, pp. 1385–1392.
- [4] H. Chen, A. Gallagher, and B. Girod, "Describing clothing by semantic attributes," in *Proc. 12th Eur. Conf. Comput. Vis. (ECCV)*, 2012, pp. 609–623.
- [5] P. Arbelaez, M. Maire, C. Fowlkes, and J. Malik, "Contour detection and hierarchical image segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 5, pp. 898–916, May 2011.
- [6] M. A. Fischler and R. A. Elschlager, "The representation and matching of pictorial structures," *IEEE Trans. Comput.*, vol. C-22, no. 1, pp. 67–92, Jan. 1973.
- [7] P. F. Felzenszwalb and D. P. Huttenlocher, "Pictorial structures for object recognition," *Int. J. Comput. Vis.*, vol. 61, no. 1, pp. 55–79, 2005.
- [8] M. Andriluka, S. Roth, and B. Schiele, "Pictorial structures revisited: People detection and articulated pose estimation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2009, pp. 1014–1021.
- [9] M. Eichner and V. Ferrari, "Better appearance models for pictorial structures," in *Proc. Brit. Mach. Vis. Conf. (BMVC)*, 2009, pp. 1–11.
- [10] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, "Object detection with discriminatively trained part-based models," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 9, pp. 1627–1645, Sep. 2010.
- [11] V. Ferrari, M. Marin-Jimenez, and A. Zisserman, "Progressive search space reduction for human pose estimation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2008, pp. 1–8.
- [12] D. Ramanan, "Learning to parse images of articulated bodies," in *Advances in Neural Information Processing Systems 19*, vol. 1. Cambridge, MA, USA: MIT Press, 2006, no. 6, p. 7.
- [13] B. Rothrock, S. Park, and S.-C. Zhu, "Integrating grammar and segmentation for human pose estimation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2013, pp. 3214–3221.
- [14] J. Wang, Z.-Q. Zhao, X. Hu, Y.-M. Cheung, M. Wang, and X. Wu, "Online group feature selection," in *Proc. 23rd Int. Joint Conf. Artif. Intell.*, 2013, pp. 1757–1763.
- [15] S. Belongie, G. Mori, and J. Malik, "Matching with shape contexts," in *Statistics and Analysis of Shapes*. Boston, MA, USA: Birkhäuser, 2006, pp. 81–105.
- [16] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, vol. 1. Jun. 2005, pp. 886–893.
- [17] C. Rother, V. Kolmogorov, and A. Blake, "'GrabCut': Interactive foreground extraction using iterated graph cuts," *ACM Trans. Graph.*, vol. 23, no. 3, pp. 309–314, 2004.
- [18] B. Sapp, C. Jordan, and B. Taskar, "Adaptive pose priors for pictorial structures," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2010, pp. 422–429.
- [19] B. Sapp, A. Toshev, and B. Taskar, "Cascaded models for articulated pose estimation," in *Proc. 11th Eur. Conf. Comput. Vis. (ECCV)*, 2010, pp. 406–420.
- [20] M. Sun, M. Telaprolu, H. Lee, and S. Savarese, "Efficient and exact MAP-MRF inference using branch and bound," in *Proc. Int. Conf. Artif. Intell. Statist.*, 2012, pp. 1134–1142.
- [21] B. Hasan and D. Hogg, "Segmentation using deformable spatial priors with application to clothing," in *Proc. BMVC*, 2010, pp. 1–11.

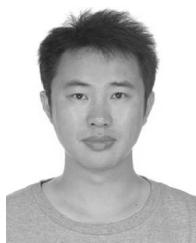
⁴Empirically, the local optima are good enough as we have demonstrated in our experiments.

- [22] N. Wang and H. Ai, "Who blocks who: Simultaneous clothing segmentation for grouping images," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Nov. 2011, pp. 1535–1542.
- [23] K. Yamaguchi, M. H. Kiapour, L. E. Ortiz, and T. L. Berg, "Parsing clothing in fashion photographs," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2012, pp. 3570–3577.
- [24] S. Liu et al., "Hi, magic closet, tell me what to wear!" in *Proc. 20th ACM Int. Conf. Multimedia*, 2012, pp. 619–628.
- [25] X. Wang and T. Zhang, "Clothes search in consumer photos via color matching and attribute learning," in *Proc. 19th ACM Int. Conf. Multimedia*, 2011, pp. 1353–1356.
- [26] S. Liu, Z. Song, G. Liu, C. Xu, H. Lu, and S. Yan, "Street-to-shop: Cross-scenario clothing retrieval via parts alignment and auxiliary set," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2012, pp. 3330–3337.
- [27] A. C. Gallagher and T. Chen, "Clothing cosegmentation for recognizing people," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2008, pp. 1–8.
- [28] H. Chen, Z. J. Xu, Z. Q. Liu, and S. C. Zhu, "Composite templates for cloth modeling and sketching," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, vol. 1, Jun. 2006, pp. 943–950.
- [29] L. Bourdev, S. Maji, and J. Malik, "Describing people: A poselet-based approach to attribute classification," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Nov. 2011, pp. 1543–1550.
- [30] L. Ladicky, P. H. S. Torr, and A. Zisserman, "Human pose estimation using a joint pixel-wise and part-wise formulation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2013, pp. 3578–3585.
- [31] I. Tsochantridis, T. Joachims, T. Hofmann, Y. Altun, and Y. Singer, "Large margin methods for structured and interdependent output variables," *J. Mach. Learn. Res.*, vol. 6, no. 9, pp. 1453–1484, 2005.
- [32] A. Argyriou, T. Evgeniou, and M. Pontil, "Multi-task feature learning," in *Advances in Neural Information Processing Systems 19*, vol. 19, Cambridge, MA, USA: MIT Press, 2007, p. 41.
- [33] J. Donahue, J. Hoffman, E. Rodner, K. Saenko, and T. Darrell, "Semi-supervised domain adaptation with instance constraints," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2013, pp. 668–675.
- [34] T. Ojala, M. Pietikainen, and D. Harwood, "Performance evaluation of texture measures with classification based on Kullback discrimination of distributions," in *Proc. 12th IAPR Int. Conf. Pattern Recognit., Conf. A, Comput. Vis. Image Process.*, vol. 1, Oct. 1994, pp. 582–585.
- [35] P. Felzenszwalb and D. Huttenlocher, "Distance transforms of sampled functions," *Theory Comput.*, vol. 8, no. 1, pp. 415–428, 2012.
- [36] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin, "LIBLINEAR: A library for large linear classification," *J. Mach. Learn. Res.*, vol. 9, pp. 1871–1874, Jun. 2008.
- [37] V. Ferrari, M. Marin-Jimenez, and A. Zisserman, "Pose search: Retrieving people using their pose," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2009, pp. 1–8.



USA.

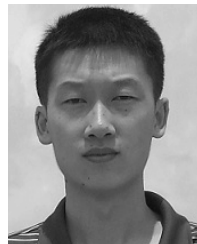
Jie Shen received the bachelor's degree in mathematics and the M.S. degree in computer science and engineering from Shanghai Jiao Tong University (SJTU), Shanghai, China, in 2011 and 2014, respectively. From 2013 to 2014, he was a Visiting Scholar with the Learning and Vision Research Group, National University of Singapore, Singapore. His research interests mainly include machine learning and computer vision. He is currently pursuing the Ph.D. degree with the Department of Computer Science, Rutgers University, New Brunswick, NJ, USA.



Guangcan Liu (M'13) is currently a Professor with the School of Information and Control Engineering, Nanjing University of Information Science and Technology, Nanjing, China. He received the bachelor's degree in mathematics and the Ph.D. degree in computer science and engineering from Shanghai Jiao Tong University, Shanghai, China, in 2004 and 2010, respectively. From 2011 to 2014, he was a Post-Doctoral Researcher in several universities, including the National University of Singapore, Singapore (2011–2012), the University of Illinois at Urbana-Champaign, Champaign, IL, USA (2012–2013), and Cornell University, Ithaca, NY, USA (2013–2014). His research interests mainly include machine learning, computer vision, and image processing.



Jia Chen received the dual bachelor's degree in mathematics and computer science from Shanghai Jiao Tong University, Shanghai, China, in 2008, where he is currently pursuing the Ph.D. degree with the Department of Computer Science and Engineering. His research interests include image annotation, content-based image retrieval, and machine learning.



Yuqiang Fang received the master's degree in control science and engineering from the National University of Defense Technology, Changsha, China, in 2010, where he is currently pursuing the Ph.D. degree in pattern recognition and intelligent systems. He was with the Learning and Vision Research Group, National University of Singapore, Singapore, as a Research Assistant in 2013. His research interests lie in machine learning and computer vision, and their applications in autonomous vehicle.

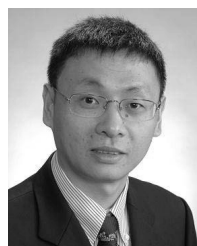


Jianbin Xie is currently a Professor with the National University of Defense Technology (NUDT), Changsha, China. He received the B.S., M.S., and Ph.D. degrees from NUDT in 1993, 1996, and 1999, respectively, where he is currently a Research Fellow. His main research interests include biometric identification, pattern recognition, and machine learning.



Web, Web mining, and information retrieval.

Yong Yu received the master's degree from the Department of Computer Science, East China Normal University, Shanghai, China. He was with Shanghai Jiao Tong University (SJTU), Shanghai, in 1986, where he is currently a Professor and Ph.D. Candidate Tutor with the Department of Computer Science and Engineering, and has taught the Data Structure course. As the Head Coach of the ACM-ICPC Team at SJTU, he and his team received the 2002, 2005, and 2010 ACM ICPC Championships. His research interests include semantic



Shuicheng Yan is currently an Associate Professor with the Department of Electrical and Computer Engineering, National University of Singapore (NUS), Singapore, where he is also the Founding Lead of the Learning and Vision Research Group. His research areas include machine learning, computer vision, and multimedia, and has authored or co-authored hundreds of technical papers over a wide range of research topics, with the Google Scholar citation of over 13 000 times and an H-index of 50. He has served as an Associate Editor of the

IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS FOR VIDEO TECHNOLOGY, and the *ACM Transactions on Intelligent Systems and Technology*. He was a recipient of the Best Paper Awards from ACM MM'13 (Best Paper and Best Student Paper), ACM MM12 (Best Demo), PCM'11, ACM MM10, ICME10, and ICIMCS'09, the Runner-Up Prize of ILSVRC'13, the Winner Prize of ILSVRC14 detection task, the Winner Prizes of the classification task in PASCAL VOC 2010–2012, the Winner Prize of the segmentation task in PASCAL VOC 2012, the Honorable Mention Prize of the detection task in PASCAL VOC'10, the 2010 TCSVT Best Associate Editor Award, the 2010 Young Faculty Research Award, the 2011 Singapore Young Scientist Award, and the 2012 NUS Young Researcher Award.