

# Boosted Multiple Deformable Trees for Parsing Human Poses

Yang Wang and Greg Mori

School of Computing Science  
Simon Fraser University  
Burnaby, BC, Canada  
{[ywang12](mailto:ywang12@cs.sfu.ca),[mori](mailto:mori@cs.sfu.ca)}@cs.sfu.ca

**Abstract.** *Tree-structured models have been widely used for human pose estimation, in either 2D or 3D. While such models allow efficient learning and inference, they fail to capture additional dependencies between body parts, other than kinematic constraints. In this paper, we consider the use of multiple tree models, rather than a single tree model for human pose estimation. Our model can alleviate the limitations of a single tree-structured model by combining information provided across different tree models. The parameters of each individual tree model are trained via standard learning algorithms in a single tree-structured model. Different tree models are combined in a discriminative fashion by a boosting procedure. We present experimental results showing the improvement of our model over previous approaches on a very challenging dataset.*

## 1 Introduction

Estimating human body poses from still images is arguably one of the most difficult object recognition problems in computer vision. The difficulties of this problem are manifold – humans are articulated objects, and can bend and contort their bodies into a wide variety of poses; the parts which make up a human figure are varied in appearance (due to clothing), which makes them difficult to reliably detect; and parts often have small support in the image or are occluded. In order to reliably interpret still images of human figures, it is likely that multiple cues relating different parts of the figure will need to be exploited.

Many existing approaches to this problem model the human body as a combination of rigid parts, connected together in some fashion. The typical configuration constraints used are kinematic constraints between adjacent parts, such as torso-upper half-limb connection, or upper-lower half-limb connection (e.g. Fig. 1). This set of constraints has a distinct computational advantage – since the constraints form a tree-structured model, inferring the optimal pose of the person using this model is tractable.

However, this computational advantage comes at a cost. Simply put, the single tree model does not adequately model the full set of relationships between parts of the body. Relationships between parts not connected in the kinematic tree cannot be directly captured by this model.

In this paper, we develop a framework for modeling human figures as a collection of trees. We argue that this framework has the advantage of being able to locally capture constraints between the parts which constitute the model. With a collection of trees, a global set of constraints can be modeled. In our work, these constraints are spatial constraints, but this framework could be extended to other cues (e.g. color consistency, occlusion relationships). We demonstrate that the computational advantages of tree-structured models can be kept, and provide tractable algorithms for learning and inference in these multiple tree models.

The rest of this paper is organized as follows. Section 2 reviews previous work. Section 3 gives the details of our approach. Section 4 shows some experimental results. Section 5 concludes this paper and points to some future work.

## 2 Related Work

One of the earliest lines of research related to finding people from images is in the setting of detecting and tracking pedestrians. Starting with the work of Hogg [4], there have been a lot of work done in tracking with kinematic models in both 2D and 3D. Forsyth et al. [3] provide a survey of this work.

Some of these approaches are exemplar-based. For example, Toyama & Blake [26] use 2D exemplars for people tracking. Mori & Malik [13] and Sullivan & Carlsson [24] address the pose estimation problems as 2D template matching using pre-stored exemplars upon which joint locations have been marked. In order to deal with the complexity due to variations of pose and clothing, Shakhnarovich et al. [19] adopt a brute-force search, using a variant of locality sensitive hashing for speed. Exemplar-based models are effective when dealing with regular human poses. However, they cannot handle those poses that rarely occur. See Fig. 5 for some examples.

There are many approaches which explicitly model the human body as an assembly of parts. Ju et al. [7] introduce a “cardboard people” model, where body parts are represented by a set of connected planar patches. Felzenszwalb & Huttenlocher [2] develop a tree-structured model called pictorial structure (PS) and applied it to 2D human pose estimation. Lee & Cohen [10] present results on 3D pose estimation from a single image based on proposal maps, using skin and face detection as extra cues to guide the MCMC sampling of 3D models. Ramaman & Forsyth [16] describe a self-starting tracker that tracks people by building an appearance model from a stylized pose detected by a top-down PS method. Sudderth et al. [23] introduce a non-parametric belief propagation method with occlusion reasoning for hand tracking. Sigal & Black [20] use a similar idea for pose estimation. Ren et al. [18] use bottom-up detections of parallel lines as part hypotheses, and combine these hypotheses with various pairwise part constraints via an integer quadratic programming. Hua et al. [5] use bottom-up cues such as skin/face detection to guide a belief propagation inference algorithm. There is also some work on using segmentation as a pre-processing step [12, 14, 22].

Our work is closely related to some recent work on learning discriminative models for localization. Ramanan & Sminchisescu [17] use a variant of conditional random fields (CRF) [8] for training localization models for articulated objects, such as human figures, horses, etc. Ramanan [15] extends their work by iteratively building a region model based on color cues.

Our work is also related to boosting on structured outputs. Boosting was originally proposed for classification problems. Recently people have adopted it for various tasks where the outputs have certain structures (e.g., chains, trees, graphs). For example, Torralba et al. [25] use boosted random fields for object detection with contextual information. Truyen et al. [27] use a boosting algorithm on Markov Random Fields for multilevel activity recognition.

Another line of research related to our work is on various extensions of tree models in both the computer vision and the machine learning literature. Song et al. [21] detect corner features in video sequences and model them using a decomposable triangulated graph, where the graph structure is found by a greedy search. Ioffe & Forsyth [6] propose a sampling method based on body part candidates found by a rectangle detector. Meila & Jordan [11] propose “mixtures-of-trees” that combine multiple tree models. The parameters of such models are learned by an EM algorithm in either maximum likelihood or Bayesian framework. We would like to point out that although our work seems similar to “mixtures-of-trees”, there are some important differences. Instead of using the maximum likelihood criterion, our method optimizes a loss function that is directly tied to inference. And our model is learned by an efficient boosting procedure.

### 3 Our Approach

Our method is a combination of tree-structured deformable models for human pose estimation [15, 17] and boosting on MRFs [27]. The basic idea is to model a human figure as a weighted combination of several tree-structured deformable models. The parameters of each tree model, and the weights of different trees are learned from training data in a discriminative fashion using boosting. In this section, we first review deformable models for human pose estimation (Sect. 3.1), followed by the learning and inference algorithms in such models (Sect. 3.2). Then we introduce the boosted multiple trees (Sect. 3.3).

#### 3.1 Deformable Model

Consider a human body model with  $K$  parts, where each part is represented by an oriented rectangle with fixed size. We can construct an undirected graph  $G = (V, E)$  to represent the  $K$  parts. Each part is represented by a vertex  $v_i \in V$  in  $G$ , and there exists an undirected edge  $e_{ij} = (v_i, v_j) \in E$  between vertices  $v_i$  and  $v_j$  if  $v_i$  and  $v_j$  has a spatial dependency. Let  $l_i = (x_i, y_i, \theta_i)$  be a random variable encoding the image position and orientation of the  $i$ -th part. We denote

the configuration of the  $K$  part model as  $L = (l_1, l_2, \dots, l_K)$ . Given the model parameters  $\Theta$ , the conditional probability of  $L$  in an image  $I$  can be written as:

$$Pr(L|I, \Theta) \propto \exp \left( \sum_{(i,j) \in E} \psi(l_i - l_j) + \sum_{i=1}^K \phi(l_i) \right) \quad (1)$$

$\psi(l_i - l_j)$  corresponds to a spatial prior on the part geometry, and  $\phi(l_i)$  models the local image evidence at each part located at  $l_i$ . Most previous approaches use Gaussian shape priors  $\psi(l_i - l_j) \propto \mathcal{N}(l_i - l_j; \mu_{ij}, \Sigma_{ij})$  [2, 17]. However, since we are dealing with images with a wide range of poses and aspects, Gaussian shape priors seem too rigid. Instead we choose a spatial prior using discrete binning (Fig. 2) similar to the one used in Ramanan [15]:

$$\psi(l_i - l_j) = \alpha_i^T \text{bin}(l_i - l_j) \quad (2)$$

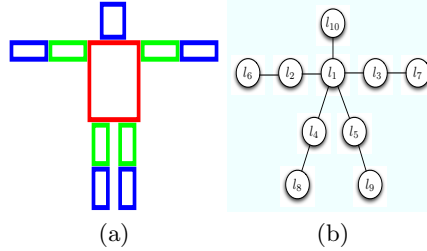
$\alpha_i$  is a parameter that favors certain relative spatial and angular bins for part  $i$  with respect to its parent  $j$ . This spatial prior captures more intricate distributions than a Gaussian prior.

For the appearance model  $\phi(l_i)$ , we follow the one used in Ramanan [15].  $\phi(l_i)$  corresponds to the local image evidence for a part and is defined as:

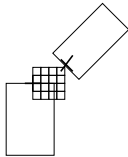
$$\phi(l_i) = \beta_i^T f_i(I(l_i)) \quad (3)$$

$f_i(I(l_i))$  is the part-specific feature vector extracted from the oriented image patch at location  $l_i$ . We use a binary vector of edges for all parts.  $\beta_i$  is a part-specific parameter that favors certain edge patterns for an oriented rectangle patch  $I(l_i)$  in image  $I$ , where  $l_i$  defines the location and orientation of the patch.

To facilitate tractable learning and inference,  $G$  is usually assumed to form a tree  $T = (V, E_T)$  [2, 15, 17]. In particular, most work uses the kinematic tree (see Fig. 1) as the underlying tree model.



**Fig. 1.** Representation of a human body. (a) human body represented as a 10-part model; (b) corresponding kinematic tree structured model.



**Fig. 2.** Discrete binning for spatial prior

### 3.2 Learning and Inference in a Single Tree Model

**Inference:** Given the model parameters  $\Theta = \{\alpha_i, \beta_i\}$ , parsing an image  $I$  of a human body involves computing the posterior distribution over part locations  $L$ , i.e.,  $P(L|I, \Theta)$ . Then the optimal part locations can be found by the *maximum a posterior* estimation  $L_{MAP} = \arg \max_L Pr(L|I, \Theta)$ . We use message-passing to carry out this computation (see [15, 17] for details).

We first pick a node (e.g., the torso) in the tree model as the root and make a directed graph from the tree structure. Then we pass messages “upstream” starting from leaf nodes to their parents. The message from part  $i$  to part  $j$  is:

$$m_i(l_j) \propto \sum_{l_i} \psi(l_i - l_j) a_i(l_i) \quad (4)$$

$$a_i(l_i) \propto \phi(l_i) \prod_{k \in kids_i} m_k(l_i) \quad (5)$$

$\phi(l_i)$  is obtained by convolving the edge image with the filter  $\beta_i$ .  $m_i(l_j)$  can be computed by convolving  $a_i(l_i)$  with a 3D spatial filter (with coefficient  $\alpha_i$ ) extending the bins from Fig. 2.  $a_i(l_i)$  is obtained by multiplying the response image  $\phi(l_i)$  together with messages from its child nodes  $m_k(l_i)$ . At the root,  $a_i$  is the true conditional marginal distribution  $Pr(l_i|I)$ . Then starting from the root, we pass messages “downstream” from part  $j$  to part  $i$  to compute the true conditional marginal of each node:

$$Pr(l_i|I) \propto a_i(l_i) \sum_{l_j} \psi(l_i - l_j) Pr(l_j|I) \quad (6)$$

It can be shown that in a tree structure, the inference is exact, and converges to the true conditional marginal distributions after this message-passing scheme. Similar to previous work [15, 17], we normalize each  $a_i$  to 1 for numerical stability, and keep track of the normalizing constants, which are needed for computing the partition function of the posterior  $Pr(L|I, \Theta)$ .

**Learning  $\Theta_{ML}$ :** If we are given a set of training images  $I^t$  where part locations  $L^t$  have been labeled, one way of learning the model parameters  $\Theta = \{\alpha_i, \beta_i\}$  is to maximize the joint likelihood of the labeled data:

$$\Theta_{ML} = \max_{\Theta} \prod_t Pr(I^t, L^t | \Theta) \quad (7)$$

$$= \max_{\Theta} \prod_t Pr(L^t | \Theta) \prod_t Pr(I^t | L^t, \Theta) \quad (8)$$

$\Theta_{ML}$  is also known as the maximum likelihood (ML) estimate of the model parameter  $\Theta$ .  $\Theta_{ML}$  can be found by independently fitting the ML estimate of each factor [17].

**Learning  $\Theta_{CL}$ :** It has been noticed [17] that the ML estimate is not directly tied to the inference, and a better criterion is to optimize the posterior distribution:

$$\Theta_{CL} = \max_{\Theta} \prod_t Pr(L^t | I^t, \Theta) \quad (9)$$

Finding  $\Theta_{CL}$  is equivalent to learning a Conditional Random Field (CRF) [8]. There are standard algorithms to learn  $\Theta_{CL}$  using gradient ascent methods.

### 3.3 Boosted Multiple Trees

There is a trade-off between representational power and computational complexity amongst different forms of spatial priors. A complete graph captures all the possible spatial dependencies between all the parts, but the learning and inference of such models are intractable. On the other hand, tree-structured models are appealing due to their tractability. However, previous work [9, 21] has shown that tree models fail to capture some additional dependencies between body parts.

In order to alleviate the limitation of tree models, various classes of graph structures that allow tractable learning and inference have been studied, e.g., mixture of trees [11], triangulated graph [21],  $k$ -fan [1], common-factor model [9]. In this section, we present our algorithm on boosting multiple trees for human pose estimation. Our algorithm is based on AdaBoost.MRF proposed in Truyen et al. [27] with some modifications. The basic idea of this method is to combine multiple tree-structured models. Since each component of the combined model is still a tree, learning and inference will be tractable. At the same time, since we are using several trees, we can capture additional spatial dependencies that are missing from a single tree model. Although our model is similar to “mixtures of trees” at a first glance, there are some importance differences. “Mixtures of trees” is trained by the EM algorithm to maximize the likelihood of the training data, while our model is trained by boosting to minimize a loss function directly tied to inference.

Given an image  $I$ , the problem of pose estimation is to find the best part labeling  $L^*$  that maximize some function  $F(L, I)$ , i.e.  $L^* = \arg \max_L F(L, I)$ .  $F(L, I)$  is known as the “strong learner” in the boosting literature. Given a set of training examples  $(I^i, L^i), i = 1, 2, \dots, N$ .  $F(L, I)$  is found by minimizing the following loss function:

$$L_O = \sum_i \sum_L \exp (F(I^i, L) - F(I^i, L^i)) \quad (10)$$

We assume  $F(L, I)$  is a linear combination of a set of so-called “weak learners”, i.e.,  $F(I, L) = \sum_t \alpha_t f_t(L, I)$ . The  $t$ -th weak learner  $f_t(L, I)$  and its corresponding weight  $\alpha_t$  are found by minimizing the loss function defined in Equation 10, i.e.  $(f_t, \alpha_t) = \arg \max_{f, \alpha} L_O$ .

Since we are interested in finding the distribution  $p(L|I)$ , we can choose the weak learner as  $f(L, I) = \log p(L|I)$ . To achieve computational tractability, we assume each weak learner is defined on a tree model.

If we can successfully learn a set of tree-based weak learners  $f_t(L, I)$  and their weights  $\alpha_t$ , the combination of these weak learners captures more spatial dependencies than a single tree model. At the same, the inference in this model is still tractable, since each component is a tree.

Optimizing  $L_O$  is difficult, Truyen et al. [27] suggest optimizing the following alternative loss function:

$$L_H = \sum_i \exp(-F(L^i, I^i)) \quad (11)$$

It can be shown that  $L_H$  is an upper bound of the original loss function  $L_O$ , provided that we can make sure  $\sum_j \alpha_j = 1$ . In Truyen et al. [27], the requirement  $\sum_j \alpha_j = 1$  is met by scaling down each previous weak learner’s weight by a factor of  $1 - \alpha_t$  as  $\alpha'_j \leftarrow \alpha_j(1 - \alpha_t)$ , for  $j = 1, 2, \dots, t-1$ , so that  $\sum_{j=1}^{t-1} \alpha'_j + \alpha_t = \sum_{j=1}^{t-1} \alpha_j(1 - \alpha_t) + \alpha_t = 1$ , since  $\sum_{j=1}^{t-1} \alpha_j = 1$ .

In practice, we find this trick sometimes has the undesirable effect of scaling down previous weak learners to have zero weights. So we use another method by scaling down each weak learner’s weight up to  $t$  by a factor of  $1/(1 + \alpha_t)$ , i.e.,  $\alpha'_j \leftarrow \frac{\alpha_j}{1 + \alpha_t}$  for  $j = 1, 2, \dots, t$ . It can be easily shown that we still have  $\sum_{j=1}^t \alpha'_j = \sum_{j=1}^{t-1} \frac{\alpha_j}{1 + \alpha_t} + \frac{\alpha_t}{1 + \alpha_t} = 1$ , since  $\sum_{j=1}^{t-1} \alpha_j = 1$ .

In practice, the algorithm could be very slow, since learning CRF parameters requires gradient ascent on a high dimensional space. To speed up the learning process, we employ several simple tricks. Firstly, we learn  $\Theta_{CL} = \{\alpha_i, \beta_i\}$  using the kinematic tree structure, and fix the appearance parameters  $\{\beta_i\}$  during the boosting process. The rational behind this is that multiple tree structures should only affect the spatial prior, not the appearance model. Secondly, during each boosting iteration, we learn  $\Theta_{ML}$  instead of  $\Theta_{CL}$ . Thirdly, instead of selecting the best tree structure in each iteration, we simply sequentially select a tree from a set of pre-specified tree structures. We also allow re-selecting a tree.

## 4 Experiments

We test our algorithm on the people dataset used in previous work [15, 17]. This dataset contains 305 images of people in various interesting poses. First 100 images are used for training, and the remaining 205 images for testing. We manually select three tree structures shown in Fig. 4, although it will be an interesting future work on how to automatically learn the tree structure at each iteration in an efficient way. The results are obtained by running 15 boosting

---

**Input:**  $i = 1, 2, \dots, D$  data pairs, graphs  $\{G_i = (V_i, E_i)\}$   
**Output:** set of trees with learned parameters and weights  
 Select a set of spanning trees  $\{\tau\}$   
 Choose the number of boosting iterations  $T$   
 Initialize  $\{w_{i,0} = \frac{1}{D}\}$ , and  $\alpha_1 = 1$   
**for** each boosting round  $t = 1, 2, \dots, T$   
   Select a spanning tree  $\tau_t$   
   /\* Add a weak learner \*/  
    $\Theta_t = \arg \max_{\Theta} \sum_i w_{i,t-1} \log Pr_{\tau_t}(L_i, I_i | \Theta)$   
    $f_t = \log Pr_{\tau_t}(L | I, \Theta_t)$   
   **if**  $t > 1$  **then**  
     select the step size  $0 < \alpha_t < 1$  using line searches  
   **end if**  
   /\* Update the strong learner \*/  
    $F_t = \frac{1}{1+\alpha_t} F_{t-1} + \frac{\alpha_t}{1+\alpha_t} f_t$   
   /\* Scale down the previous learners' weights \*/  
    $\alpha_j \leftarrow \frac{\alpha_j}{1+\alpha_t}$ , for  $j = 1, 2, \dots, t$   
   /\* Re-weight training data \*/  
    $w_{i,t} \propto w_{i,t-1} \exp(-\alpha_t f_{i,t})$   
**end for**  
 Output  $\{\tau_t\}, \{\Theta_t\}$  and  $\{\alpha_t\}$ ,  $t = 1, 2, \dots, T$

---

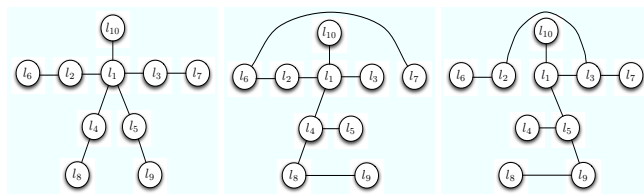
**Fig. 3.** Algorithm of boosted multiple trees

iterations. We visualize the posterior distribution  $Pr(L|I)$  on a 2D image using the same technique in Ramanan [15], where the torso is represented as red, upper-limbs as green, and lower-limbs and the head as blue. Some of the parsing results are shown in Fig. 5. We can see that our parsing results are much clearer than the one using the kinematic tree. In many images, the body parts are almost clearly visible from our parsing results. In the parsing results of using the kinematic tree, there are many white pixels, indicating high uncertainty about body parts at those locations. But with multiple trees, most of the white pixels are cleaned up. We can imagine if we sample the part candidates  $l_i$  according to  $Pr(l_i|I; \Theta)$  and use them as the inputs to other pose estimation algorithms (e.g., Ren et al. [18]), the samples generated from our parsing results are more likely to be the true part locations.

## 5 Conclusion and Future Work

We have presented a framework for modeling human figures as a collection of tree-structured models. This framework has the computational advantages of previous tree-structured models used for human pose estimation. At the same





**Fig. 4.** Three tree structures used for boosting

time, it models a richer set of spatial constraints between body parts. We demonstrate our results on a challenging dataset with substantial pose variations.

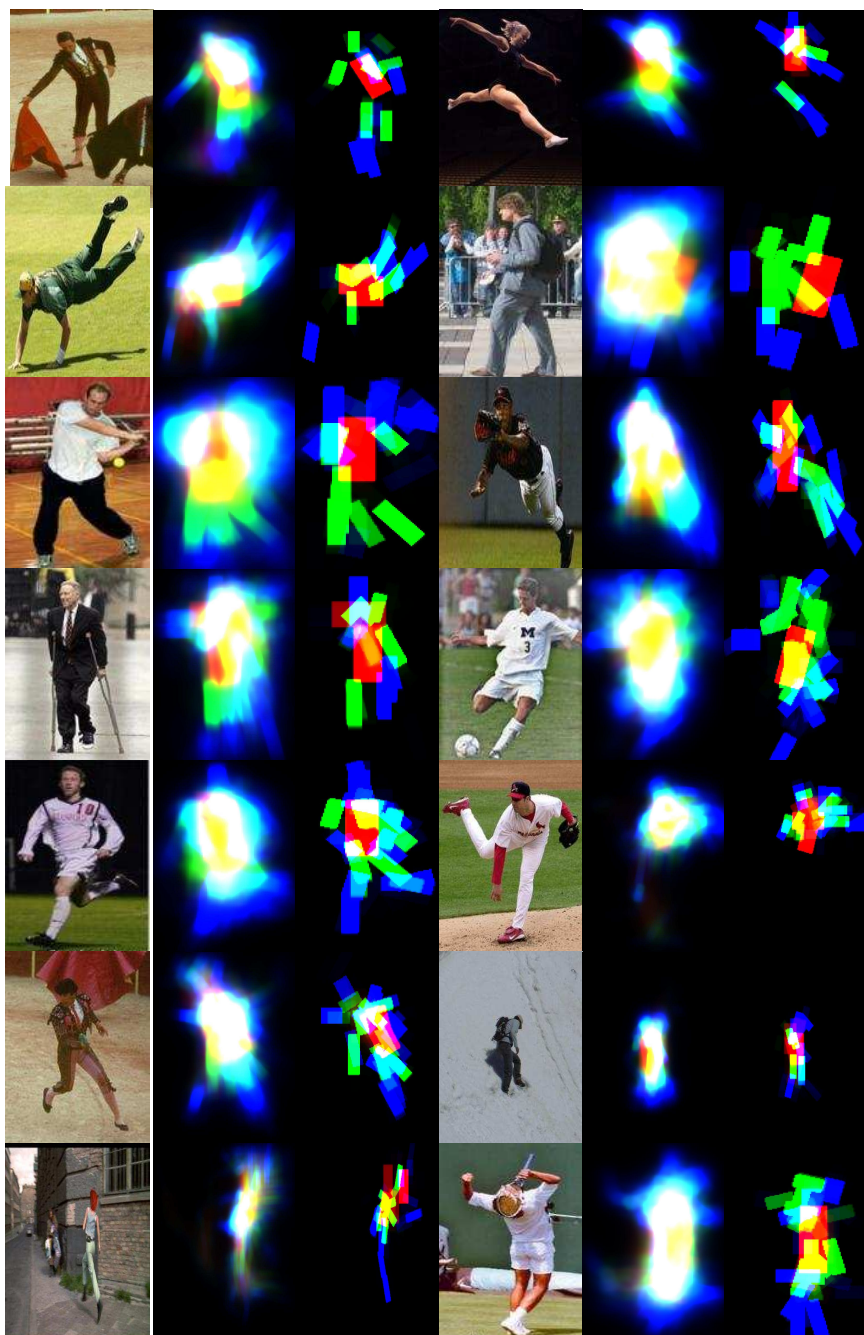
Human pose estimation is an extremely difficult computer vision problem. The solution of this problem probably requires the symbiosis of various kinds of visual cues. This paper represents our first step in that direction. Our framework nicely solves the problem of modeling spatial dependencies between non-connected body parts. In the future, we would like to extend our framework to other cues (e.g., color consistency, occlusion relationships). We would also like to combine our framework with the iterative color parsing [15].

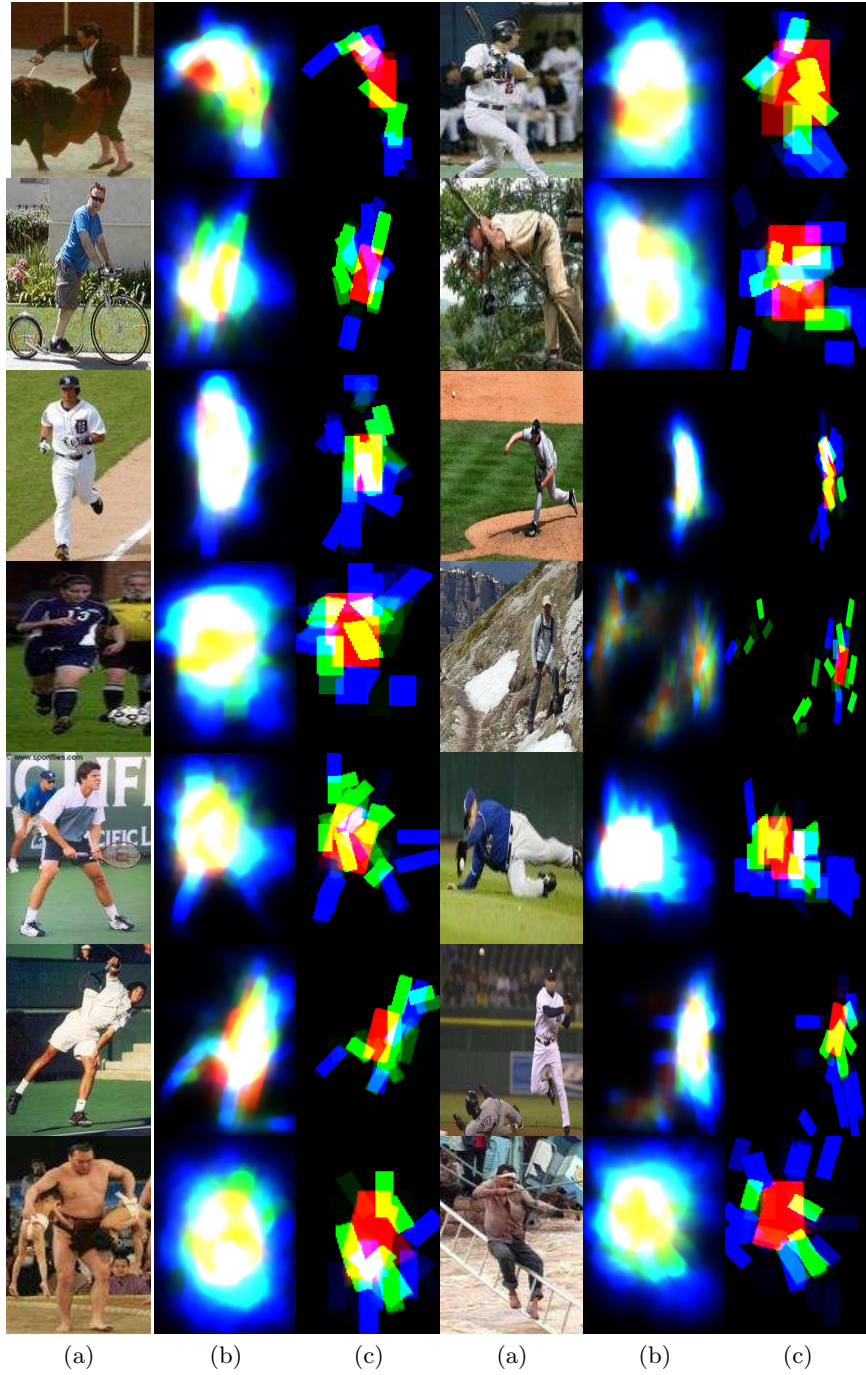
## Acknowledgement

We would like to thank Deva Ramanan for providing source code, the dataset, and many helpful discussions.

## References

1. Crandell, D., Felzenszwalb, P.F., Huttenlocher, D.P.: Spatial priors for part-based recognition using statistical models. In: IEEE Conference on Computer Vision and Pattern Recognition. Volume 1. (2005) 10–17
2. Felzenszwalb, P.F., Huttenlocher, D.P.: Pictorial structures for object recognition. *International Journal of Computer Vision* **61**(1) (January 2003) 55–79
3. Forsyth, D.A., Arikian, O., Ikemoto, L., O’Brien, J., Ramanan, D.: Computational studies of human motion: Part 1, tracking and motion synthesis. *Foundations and Trends in Computer Graphics and Vision* **1**(2/3) (July 2006) 77–254
4. Hogg, D.: Model-based vision: a program to see a walking person. *Image and Vision Computing* **1**(1) (1983) 5–20
5. Hua, G., Yang, M.H., Wu, Y.: Learning to estimate human pose with data driven belief propagation. In: IEEE Conference on Computer Vision and Pattern Recognition. Volume 2. (2005) 747–754
6. Ioffe, S., Forsyth, D.: Finding people by sampling. In: IEEE International Conference on Computer Vision. Volume 2. (1999) 1092–1097
7. Ju, S.X., Black, M.J., Yacobi, Y.: Cardboard people: A parameterized model of articulated image motion. In: International Conference on Automatic Face and Gesture Recognition. (1996) 38–44





**Fig. 5.** Some results of our algorithm: (a) original images; (b) results of using one kinematic tree; (c) results of using multiple trees.

8. Lafferty, J., McCallum, A., Pereira, F.: Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In: International Conference on Machine Learning (ICML). (2001) 282–289
9. Lan, X., Huttenlocher, D.P.: Beyond trees: Common-factor models for 2d human pose recovery. In: IEEE International Conference on Computer Vision. Volume 1. (2005) 470–477
10. Lee, M.W., Cohen, I.: Proposal maps driven mcmc for estimating human body pose in static images. In: IEEE Conference on Computer Vision and Pattern Recognition. Volume 2. (2004) 334–341
11. Meila, M., Jordan, M.I.: Learning with mixtures of trees. *Journal of Machine Learning Research* **1** (2000) 1–48
12. Mori, G.: Guiding model search using segmentation. In: IEEE International Conference on Computer Vision. Volume 2. (2005) 1417–1423
13. Mori, G., Malik, J.: Estimating human body configurations using shape context matching. In: European Conference on Computer Vision. Volume 3. (2002) 666–680
14. Mori, G., Ren, X., Efros, A., Malik, J.: Recovering human body configuration: Combining segmentation and recognition. In: IEEE Conference on Computer Vision and Pattern Recognition. Volume 2. (2004) 326–333
15. Ramanan, D.: Learning to parse images of articulated bodies. In: Advances in Neural Information Processing Systems. Volume 19. (2007) 1129–1136
16. Ramanan, D., Forsyth, D.A., Zisserman, A.: Strike a pose: Tracking people by finding stylized poses. In: IEEE Conference on Computer Vision and Pattern Recognition. Volume 1. (2005) 271–278
17. Ramanan, D., Sminchisescu, C.: Training deformable models for localization. In: IEEE Conference on Computer Vision and Pattern Recognition. Volume 1. (2006) 206–213
18. Ren, X., Berg, A., Malik, J.: Recovering human body configurations using pairwise constraints between parts. In: IEEE International Conference on Computer Vision. Volume 1. (2005) 824–831
19. Shakhnarovich, G., Viola, P., Darrell, T.: Fast pose estimation with parameter sensitive hashing. In: IEEE International Conference on Computer Vision. Volume 2. (2003) 750–757
20. Sigal, L., Black, M.J.: Measure locally, reason globally: Occlusion-sensitive articulated pose estimation. In: IEEE Conference on Computer Vision and Pattern Recognition. Volume 2. (2006) 2041–2048
21. Song, Y., Goncalves, L., Perona, P.: Unsupervised learning of human motion. *IEEE Transaction on Pattern Analysis and Machine Intelligence* **25**(7) (July 2003) 814–827
22. Srinivasan, P., Shi, J.: Bottom-up recognition and parsing of the human body. In: IEEE Conference on Computer Vision and Pattern Recognition. (2007)
23. Sudderth, E.B., Mandel, M.I., Freeman, W.T., Willsky, A.S.: Distributed occlusion reasoning for tracking with nonparametric belief propagation. In: Advances in Neural Information Processing Systems. MIT Press (2004) 1369–1376
24. Sullivan, J., Carlsson, S.: Recognizing and tracking human action. In: European Conference on Computer Vision LNCS 2352. Volume 1. (2002) 629–644
25. Torralba, A., Murphy, K.P., Freeman, W.T.: Contextual models for object detection using boosted random fields. In: Advances in Neural Information Processing Systems 17. MIT Press (2005) 1401–1408
26. Toyama, K., Blake, A.: Probabilistic exemplar-based tracking in a metric space. In: IEEE International Conference on Computer Vision. Volume 2. (2001) 50–57

27. Truyen, T.T., Phung, D.Q., Bui, H.H., Venkatesh, S.: AdaBoost.MRF: Boosted markov random forests and application to multilevel activity recognition. In: IEEE Conference on Computer Vision and Pattern Recognition. Volume 2. (2006) 1686–1693