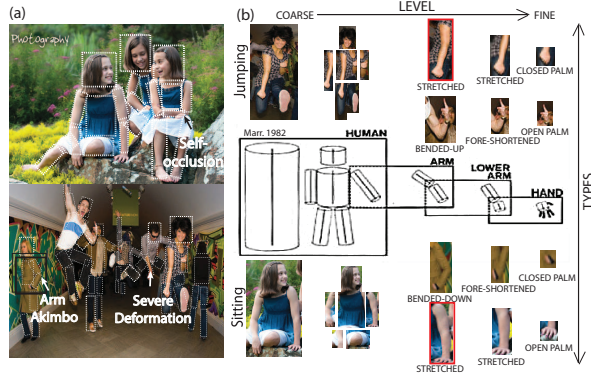


1. Overview

-Goal: Joint Object Detection and Pose Estimation

-Articulated Part-based Model:

- recursive coarse-to-fine representation
- multiple part-types
- parents-child relationship



3. Learning

-Linear Weights

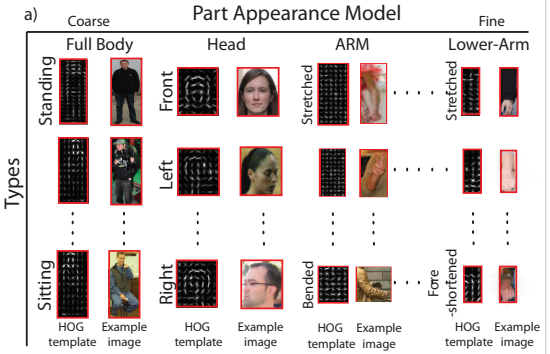
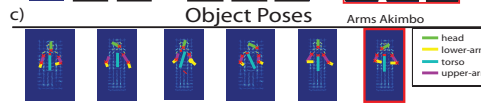
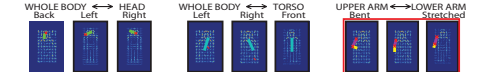
$$w^T \Psi(H; I) = \sum_i A_{(i_s, i_e)}^T \psi_a(h_i; I) + \sum_{ij} (b_{ij}^{(s, i, j)} - d_{(i_s, i_e)}^T \psi_d(h_j; T(h_i, t_{ij}^{(s, i, j)})))$$

-Struct SVM

$$\min_{w, \xi^n > 0} w^T w + C \sum_n \xi^n(H)$$

$$s. t. \xi^n(H) = \max_H (D(H; H^n) + w^T \Psi(H; I^n) - w^T \Psi(H^n; I^n))$$

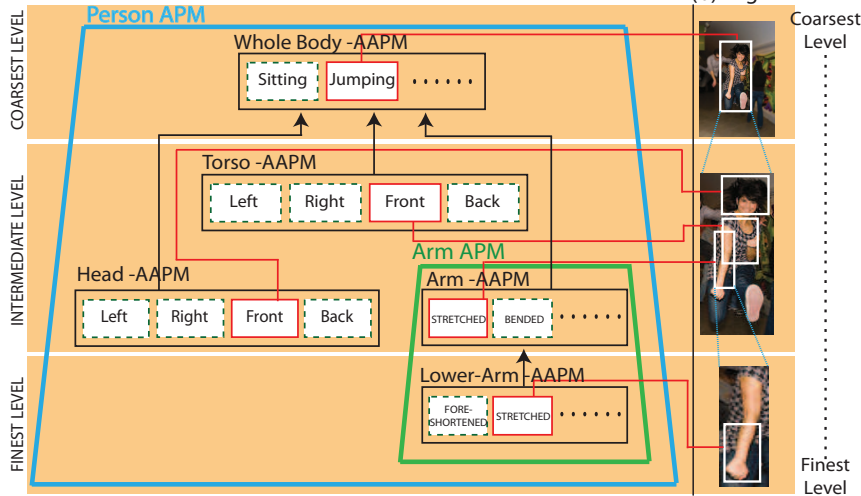
(b) Parent-Child Relationship



2. Model

(a) Model Structure

(b) Img Evidence



-Matching Score

- Appearance Score: $f^A(h; I) = A^T \psi_a(h, I)$
- Deformation Score: $f^D(h, \hat{h}) = -d^T \psi_d(h, \hat{h})$

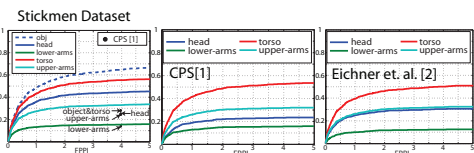
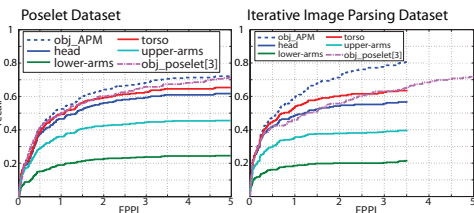
-Model Properties

- Sublinearity
- Efficient Exact Inference

-Score Aggregation

- Child location selection: $f_{c_s c}(\hat{h}_c, I) = \max_{h_c} f_{c_s c}(h_c, I) + f^D(h_c, \hat{h}_c)$
- Child alignment: $f_{c_s c}(T(h_i, t_{i_c}^{s, i, c}), I); T(h, t) = (x - t_x, y - t_y, L - t_L \theta - t_\theta)$
- Child type selection: $f_c(h_i, I) = \max_{s_c} f_{c_s c}(T(h_i, t_{i_c}^{s, i, c}), I) + b_{i_c}^{s, i, c}$
- Aggregation: $f_{i_s}(h_i, I) = f_{i_s}^A(h_i, I) + \sum_c f_c(h_i; I)$

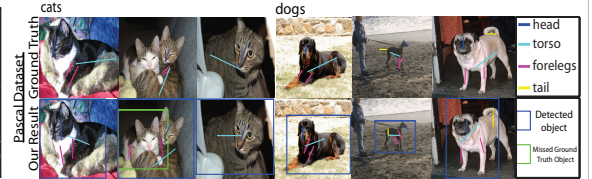
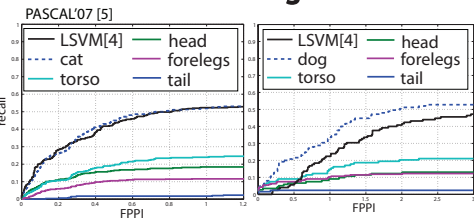
4. Results on Human Dataset



PASCAL stickmen	torso	head	upper arm	lower arm	obj	Recall
Det. Pose methods	0.311/48.43%	0.318/49.52%	0.122/19.06%			
Eichner et al [2]	0.525/81.80%	0.231/35.92%	0.316/49.27%	0.155/24.15%	0.642	
APM (ours)	0.550/85.57%	0.439/68.33%	0.326/50.73%	0.151/23.54%		



5. Results on Cats & Dogs Dataset



6. Conclusion

- Improvement in both object detection and pose estimation: recursive coarse-to-fine and multiple part-type representation
- Novel performance measure: the part recall vs. FPP1 curve

Acknowledgments

ONR grant
N000141110389

[1] B. Sapp, A. Toshev, and B. Taskar. Cascaded models for articulated pose estimation. ECCV, 2010.
 [2] M. Eichner and V. Ferrari. Better appearance models for pictorial structures. BMVC, 2009.
 [3] J. Bourdev, S. Maji, T. Brox, and J. Malik. Detecting people using mutually consistent poselet activations. In ECCV, 2010.
 [4] F. F. Fitzgibbon, B. G. S. Coombes, D. McArtor, and D. Kananian. Object detection with discriminatively trained part-based models. TPAMI, 2011.
 [5] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL VOC2007 Results.