

# Pictorial Structures Revisited: People Detection and Articulated Pose Estimation



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT



Mykhaylo Andriluka



Stefan Roth



Bernt Schiele

Department of Computer Science  
TU Darmstadt

# Generic model for human detection and pose estimation



## Human pose estimation

[Felzenszwalb&Huttenlocher, ICCV'05], [Ren et al., ICCV'05], [Sigal&Black, CVPR'06], [Zhang et al., CVPR'06], [Jiang&Marin, CVPR'08], [Ramanan, NIPS'06], [Ferrari et al., CVPR'08], [Ferrari et al., CVPR'09]



often rather simple appearance model  
focus on finding optimal assembly of parts



## People Detection

[Viola et al., ICCV'03], [Dalal&Triggs, CVPR'05], [Leibe et al., CVPR'05], [Andriluka et al., CVPR'08]



complex appearance model  
no pose model or limited to walking motion

# Generic model for human detection and pose estimation



## Human pose estimation

[Felzenszwalb&Huttenlocher, ICCV'05], [Ren et al., ICCV'05], [Sigal&Black, CVPR'06], [Zhang et al., CVPR'06], [Jiang&Marin, CVPR'08], [Ramanan, NIPS'06], [Ferrari et al., CVPR'08], [Ferrari et al., CVPR'09]



often rather simple appearance model  
focus on finding optimal assembly of parts



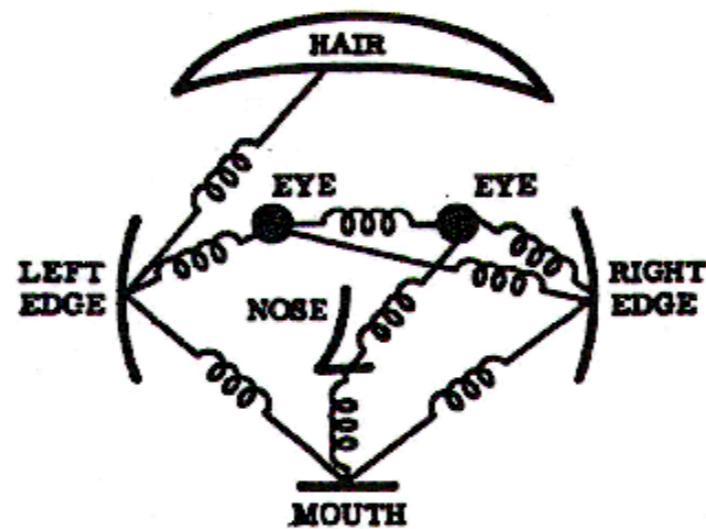
## People Detection

[Viola et al., ICCV'03], [Dalal&Triggs, CVPR'05], [Leibe et al., CVPR'05], [Andriluka et al., CVPR'08]



complex appearance model  
no pose model or limited to walking motion

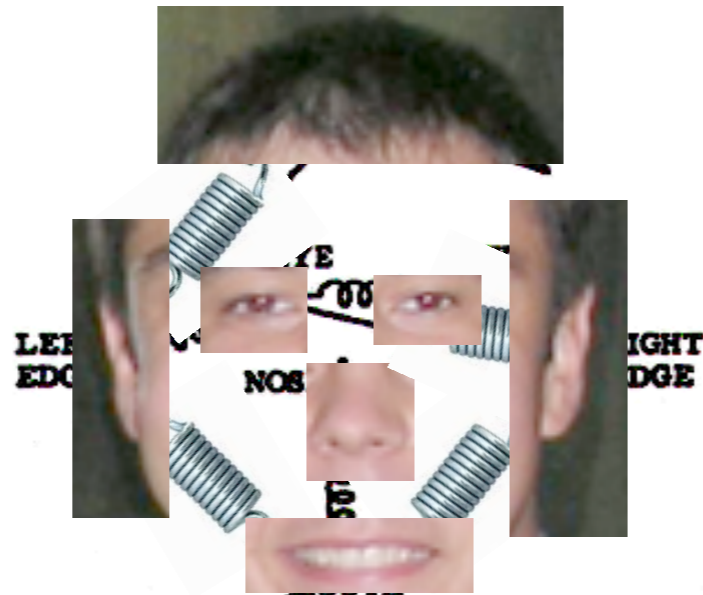
# Can we make pictorial structures model effective for these tasks?



[Fischler&Elschlager, 1973]



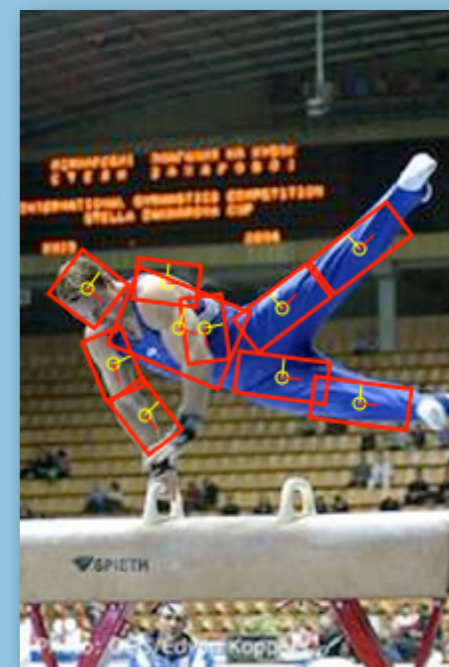
# Can we make pictorial structures model effective for these tasks?



Yes... if the model components are chosen right.

# Pictorial Structures Model

- Body is represented as **flexible configuration of body parts**

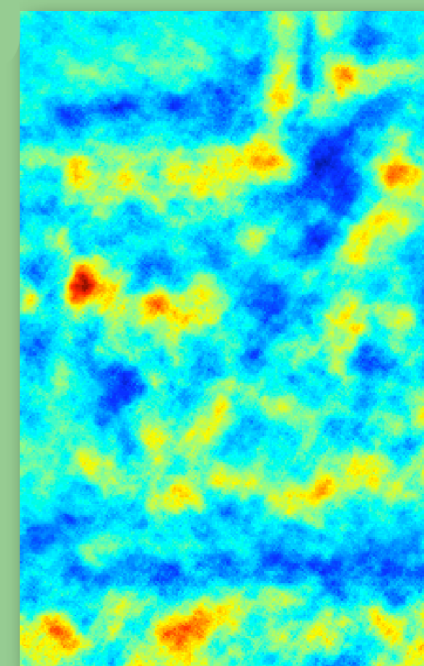
 $L$  $L = \{\mathbf{l}_0, \mathbf{l}_1, \dots, \mathbf{l}_N\}$  - configuration of parts $D = \{\mathbf{d}_0, \mathbf{d}_1, \dots, \mathbf{d}_N\}$  - part evidence

posterior over body poses

$$p(L|D) \propto p(D|L)p(L)$$

likelihood of observations

prior on body poses

 $\mathbf{d}_i$ 

# Pictorial Structures Model

Pictorial structures allow  
**exact** and **efficient** inference.

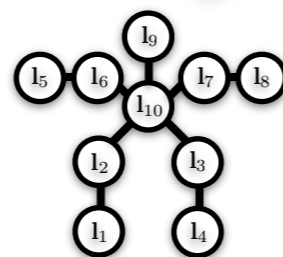


- tree-structured prior
- independent part appearance model
- discretized part locations

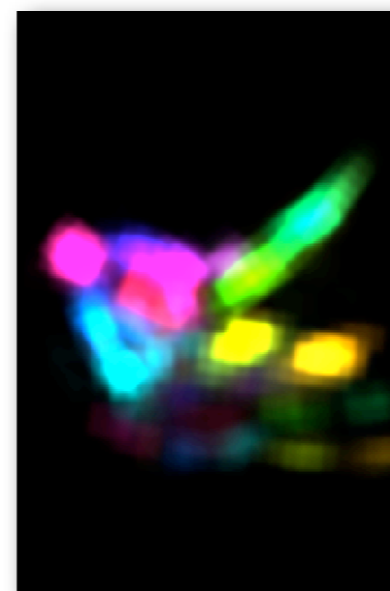
- Gaussian pairwise part relationships



sum-product BP

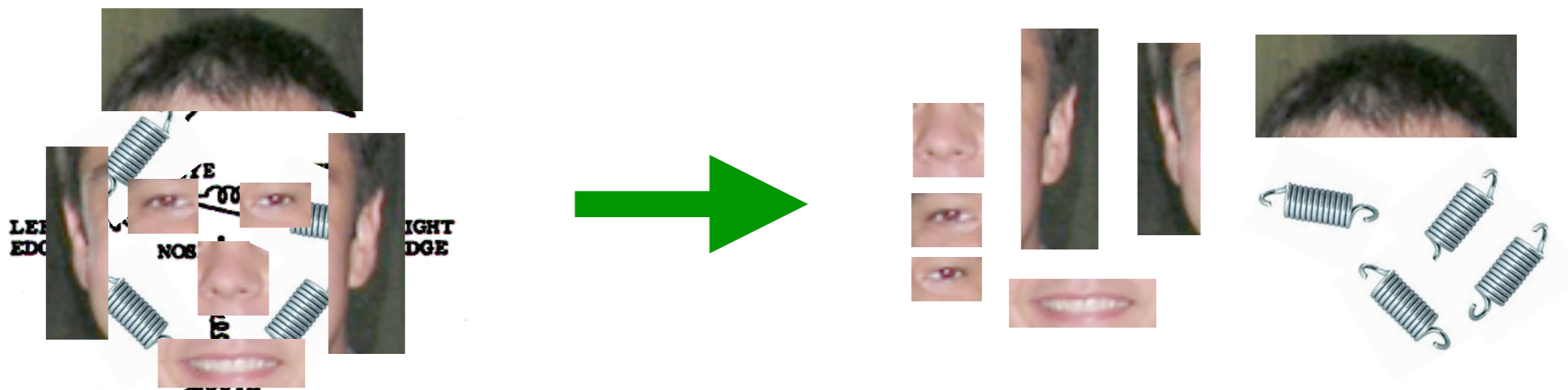


posterior marginals



$$p(\mathbf{l}_i | D) \propto \sum_{L \setminus \mathbf{l}_i} p(L | D)$$

# Can we make pictorial structures model effective for these tasks?

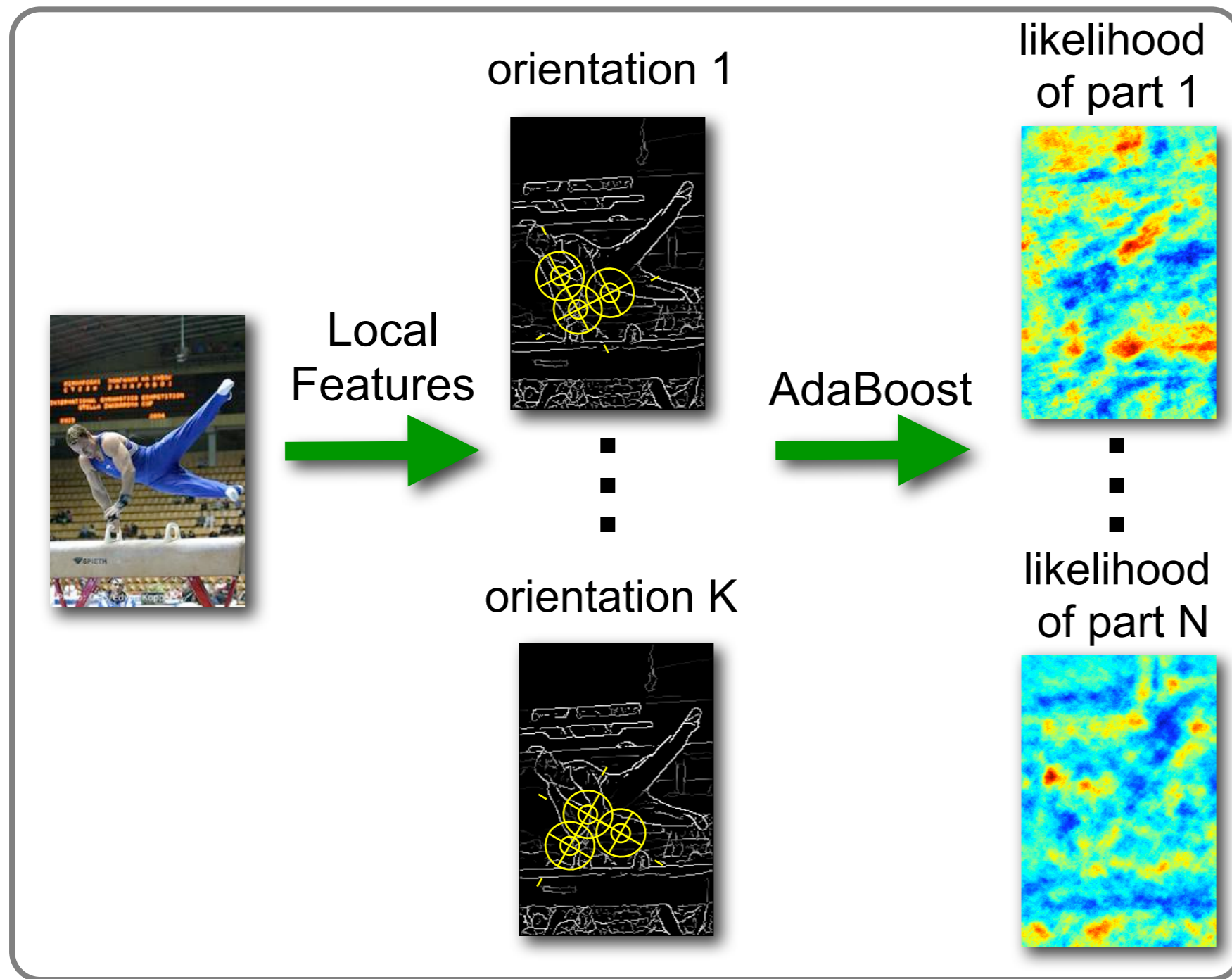


So... what are the right components?

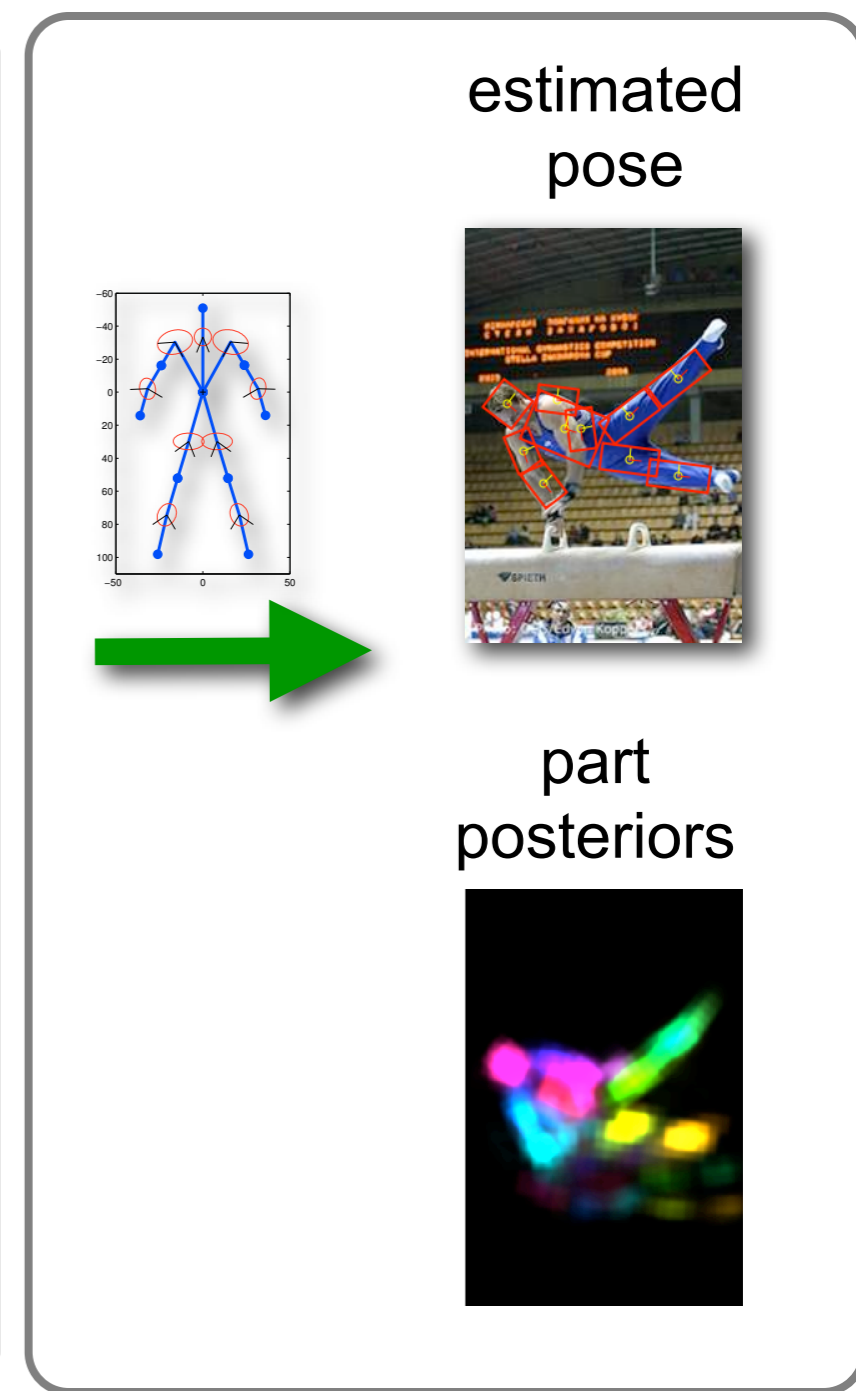


# Model Components

## Appearance Model:

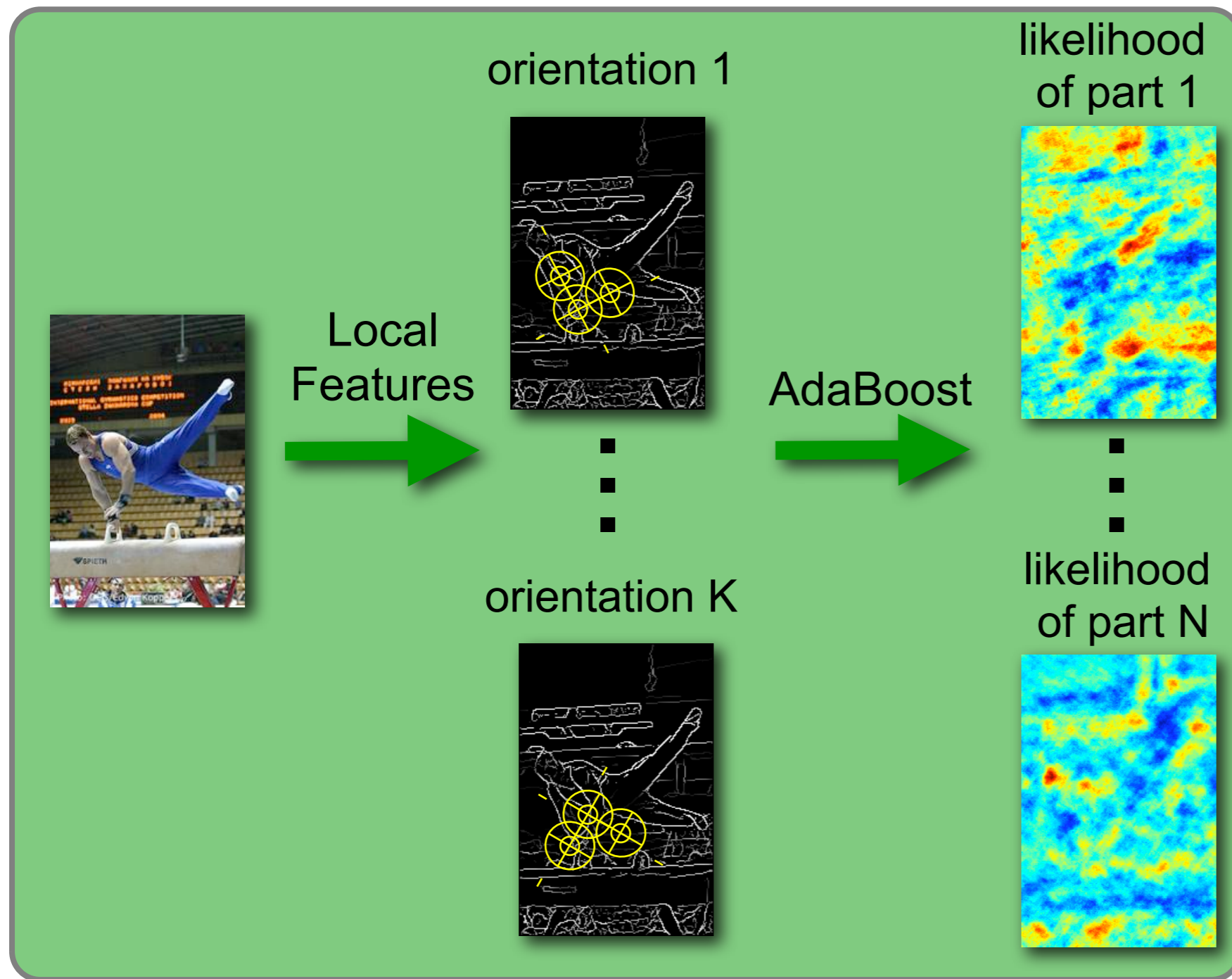


## Prior and Inference:

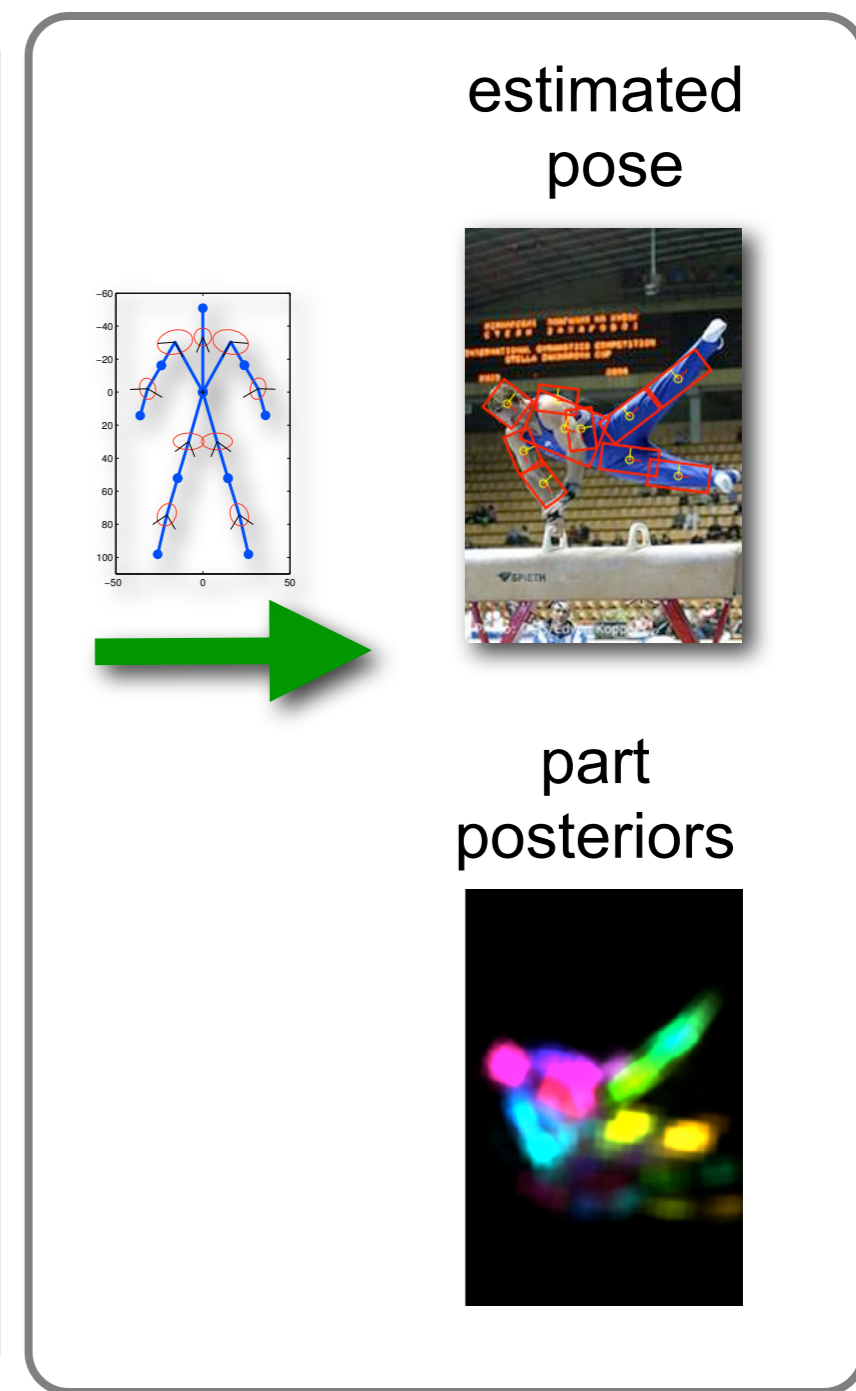


# Model Components

## Appearance Model:

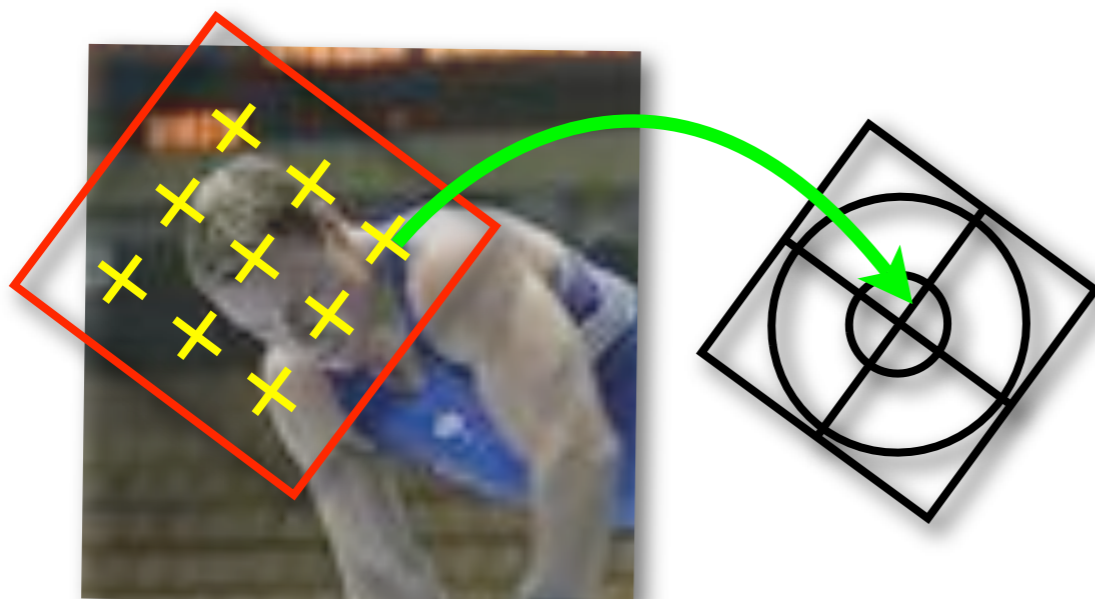


## Prior and Inference:



# Likelihood Model

- Build on recent advances in object detection:
  - ▶ state-of-the-art image descriptor: **Shape Context**  
[Belongie et al., PAMI'02; Mikolajczyk&Schmid, PAMI'05]
  - ▶ **dense representation**
  - ▶ discriminative model: **AdaBoost** classifier for each body part



- Shape Context: 96 dimensions  
(4 angular, 3 radial, 8 gradient orientations)
- Feature Vector: concatenate the descriptors inside part bounding box
- head: 4032 dimensions
- torso: 8448 dimensions

# Likelihood Model

- Part likelihood derived from the boosting score:

decision stump weight      decision stump output

$$\tilde{p}(\mathbf{d}_i | \mathbf{l}_i) = \max \left( \frac{\sum_t \alpha_{i,t} h_t(\mathbf{x}(\mathbf{l}_i))}{\sum_t \alpha_{i,t}}, \varepsilon_0 \right)$$

part location

small constant to deal with part occlusions



# Likelihood Model

Input image

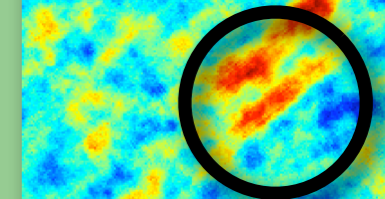
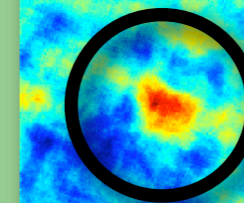
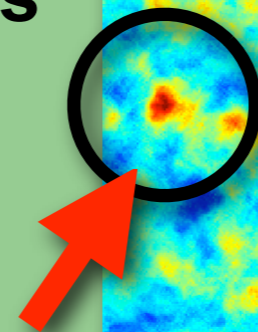


Head

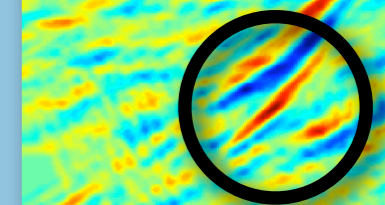
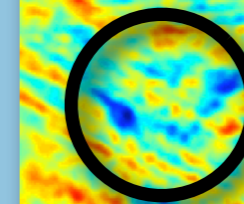
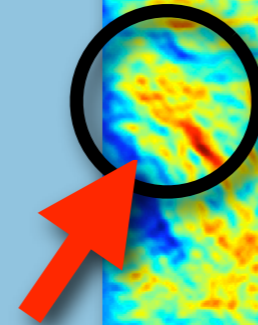
Torso

Upper leg

Our part  
likelihoods



[Ramanan,  
NIPS'06]





# Likelihood Model

Input image

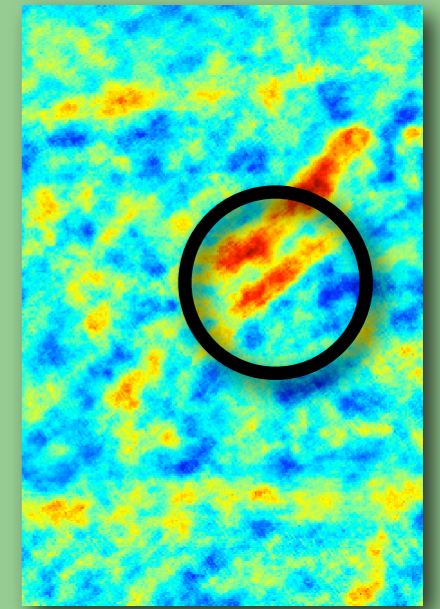
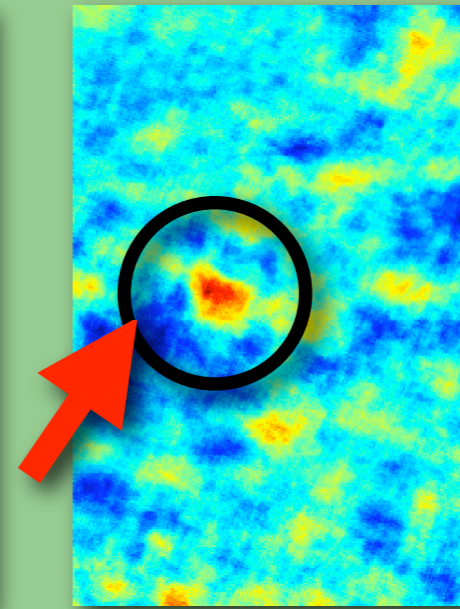
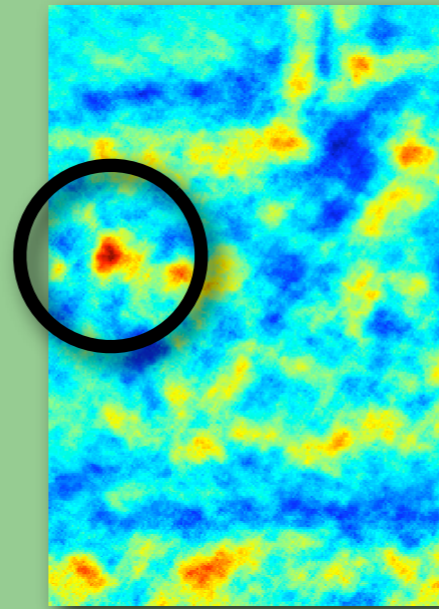


Head

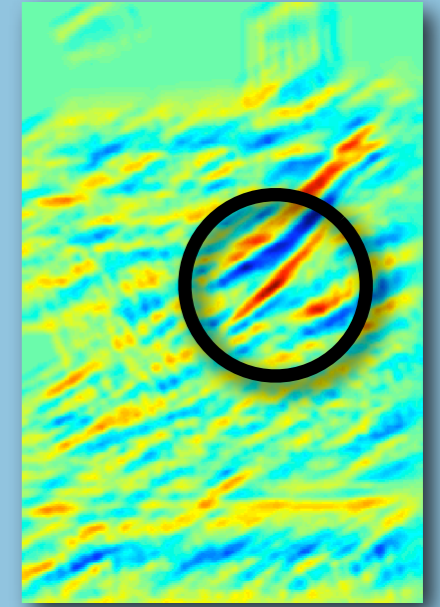
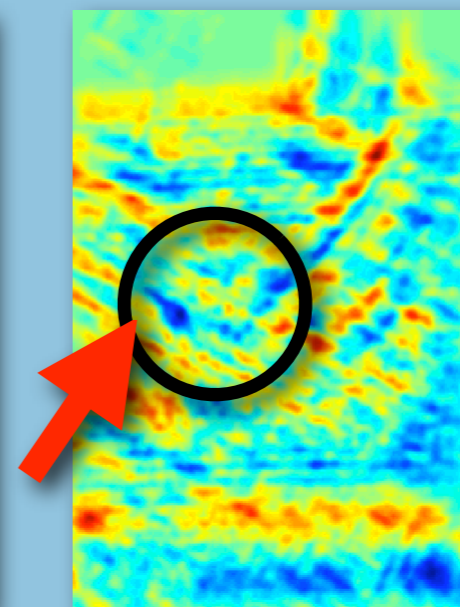
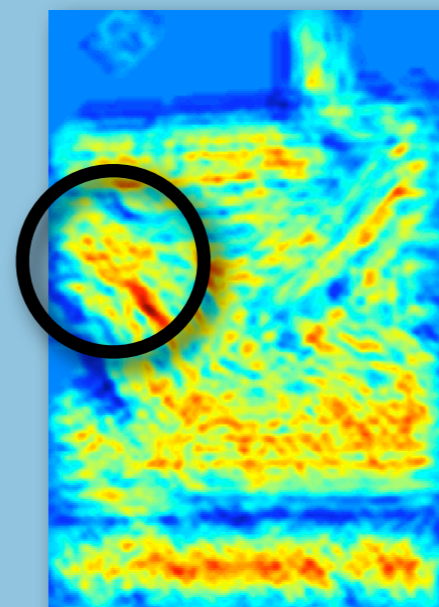
Torso

Upper leg

Our part  
likelihoods



[Ramanan,  
NIPS'06]





# Likelihood Model

Input image

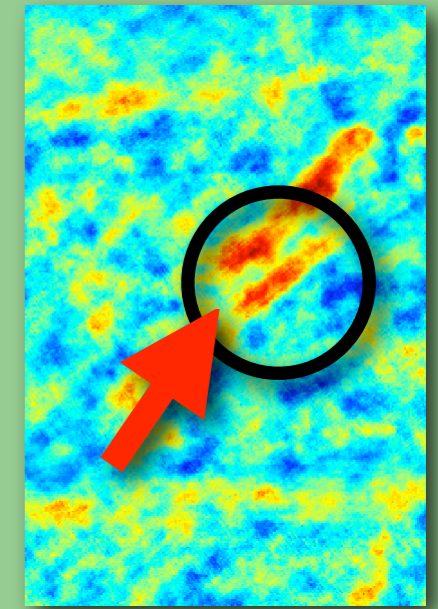
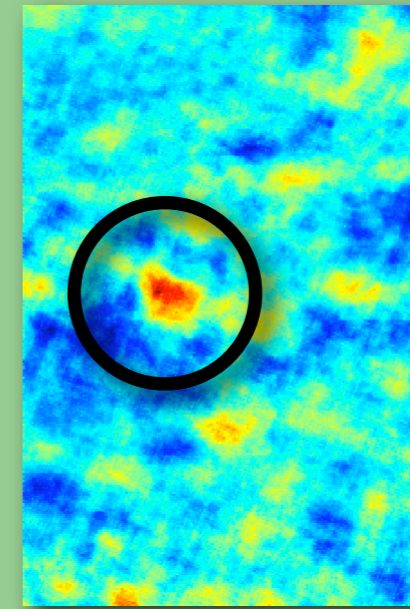
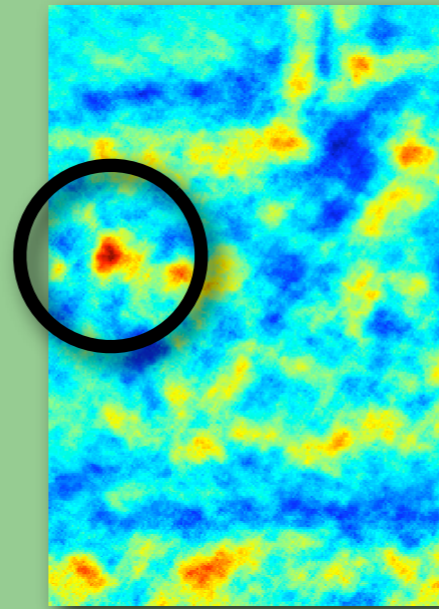


Head

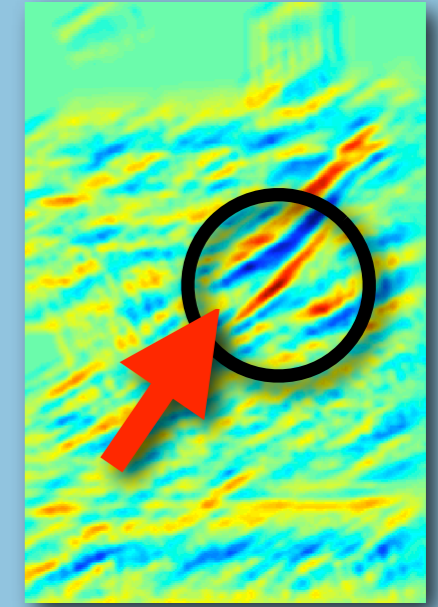
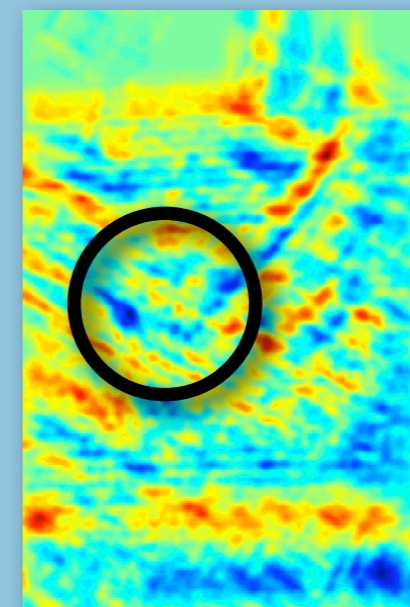
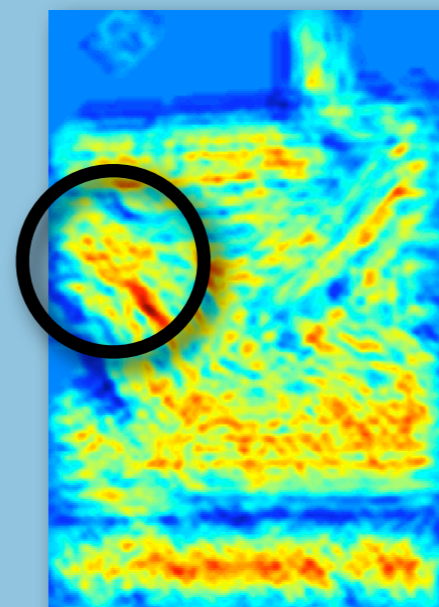
Torso

Upper leg

Our part  
likelihoods

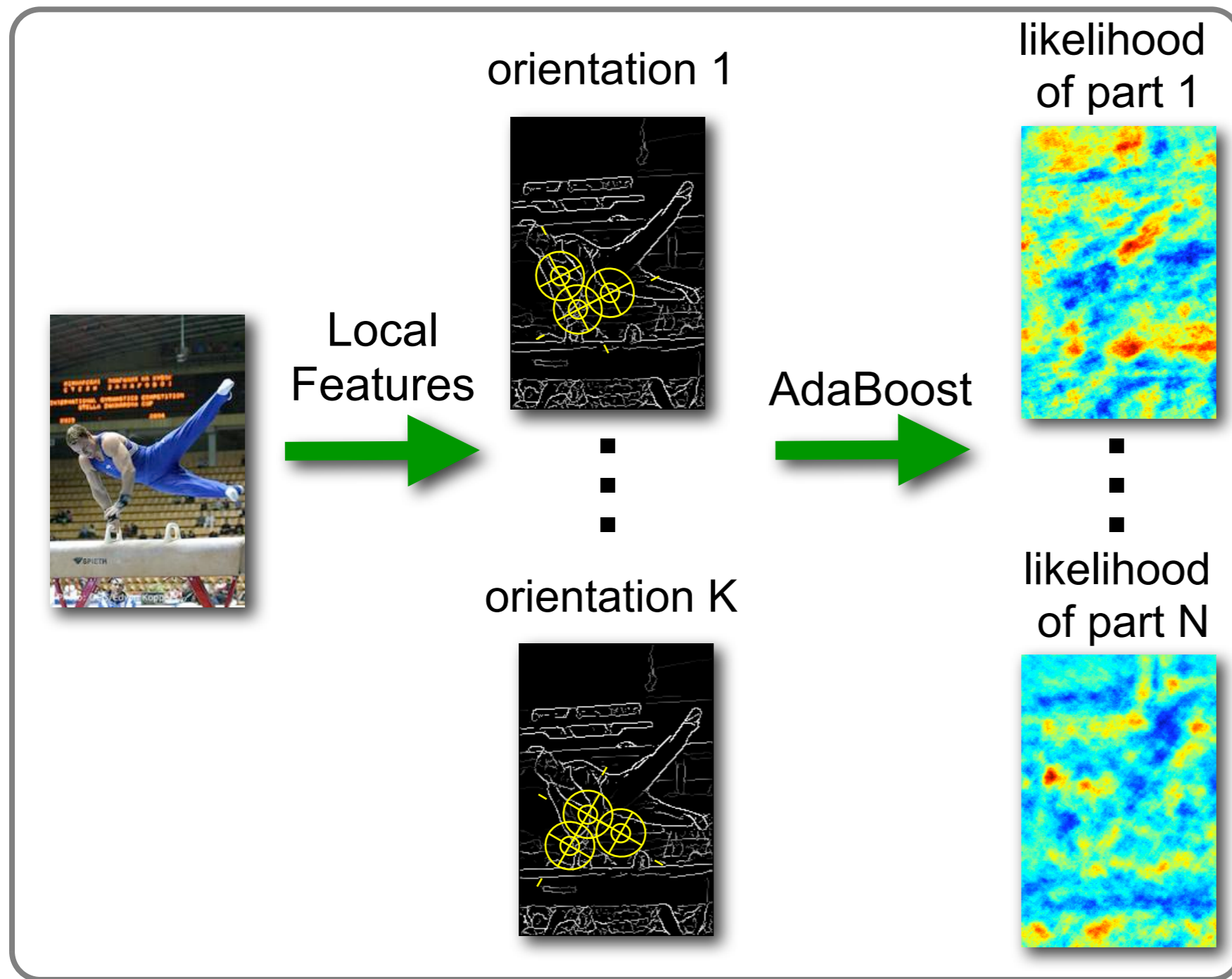


[Ramanan,  
NIPS'06]

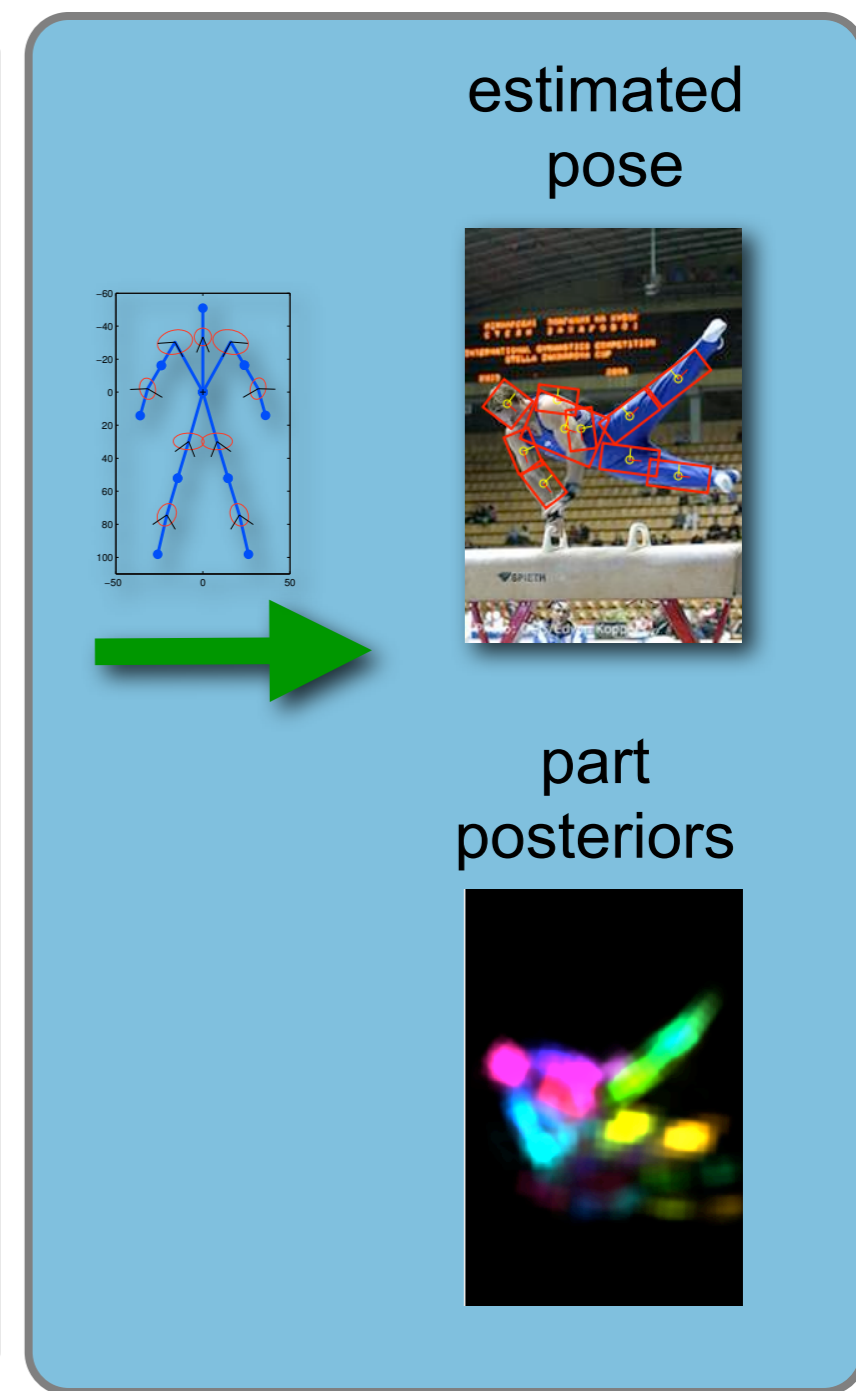


# Model Components

## Appearance Model:



## Prior and Inference:



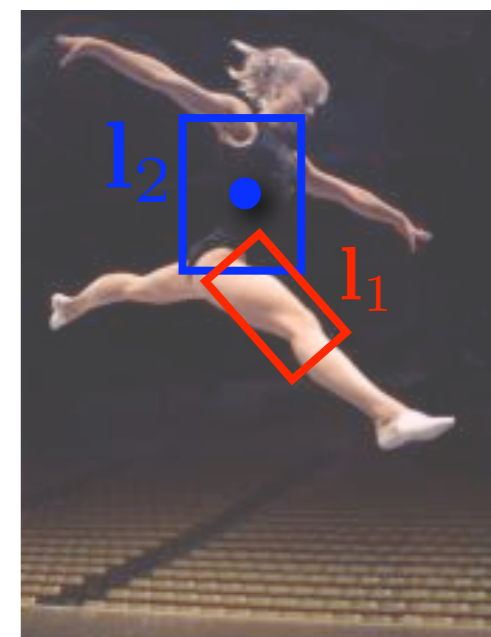


# Kinematic Tree Prior

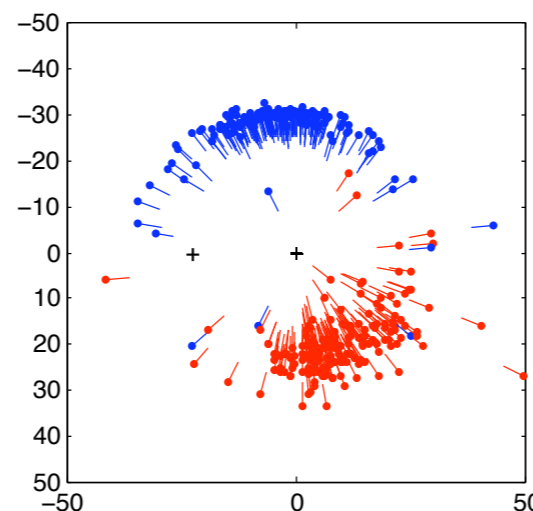
- Represent pairwise part relations  
[Felzenszwalb & Huttenlocher, IJCV'05]

$$p(L) = p(\mathbf{l}_0) \prod_{(i,j) \in E} p(\mathbf{l}_i | \mathbf{l}_j),$$

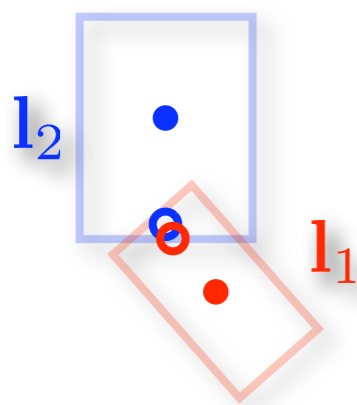
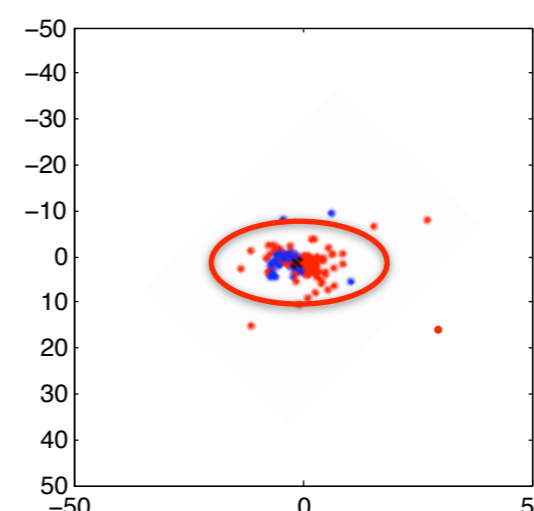
$$p(\mathbf{l}_2 | \mathbf{l}_1) = \mathcal{N}(T_{12}(\mathbf{l}_2) | T_{21}(\mathbf{l}_1), \Sigma^{12})$$



part locations relative  
to the joint



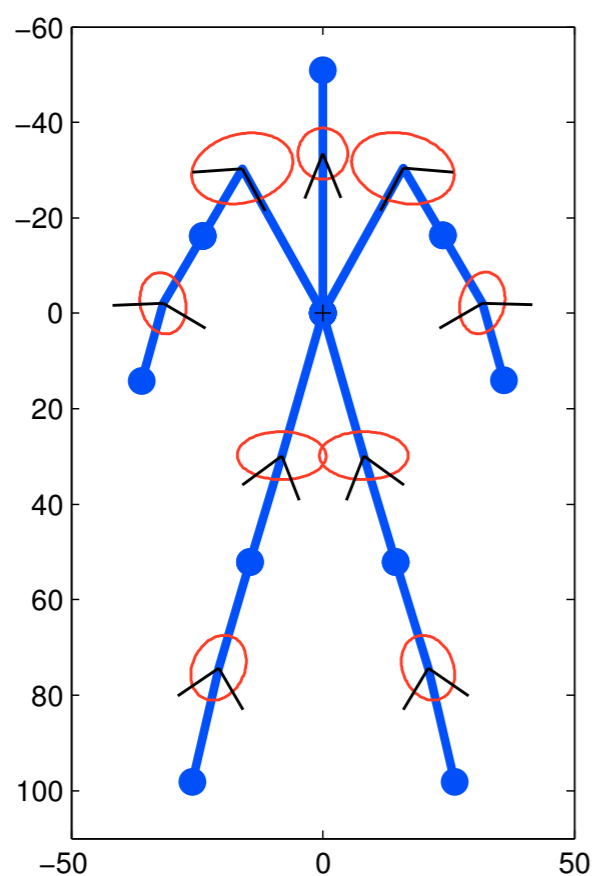
transformed  
part locations



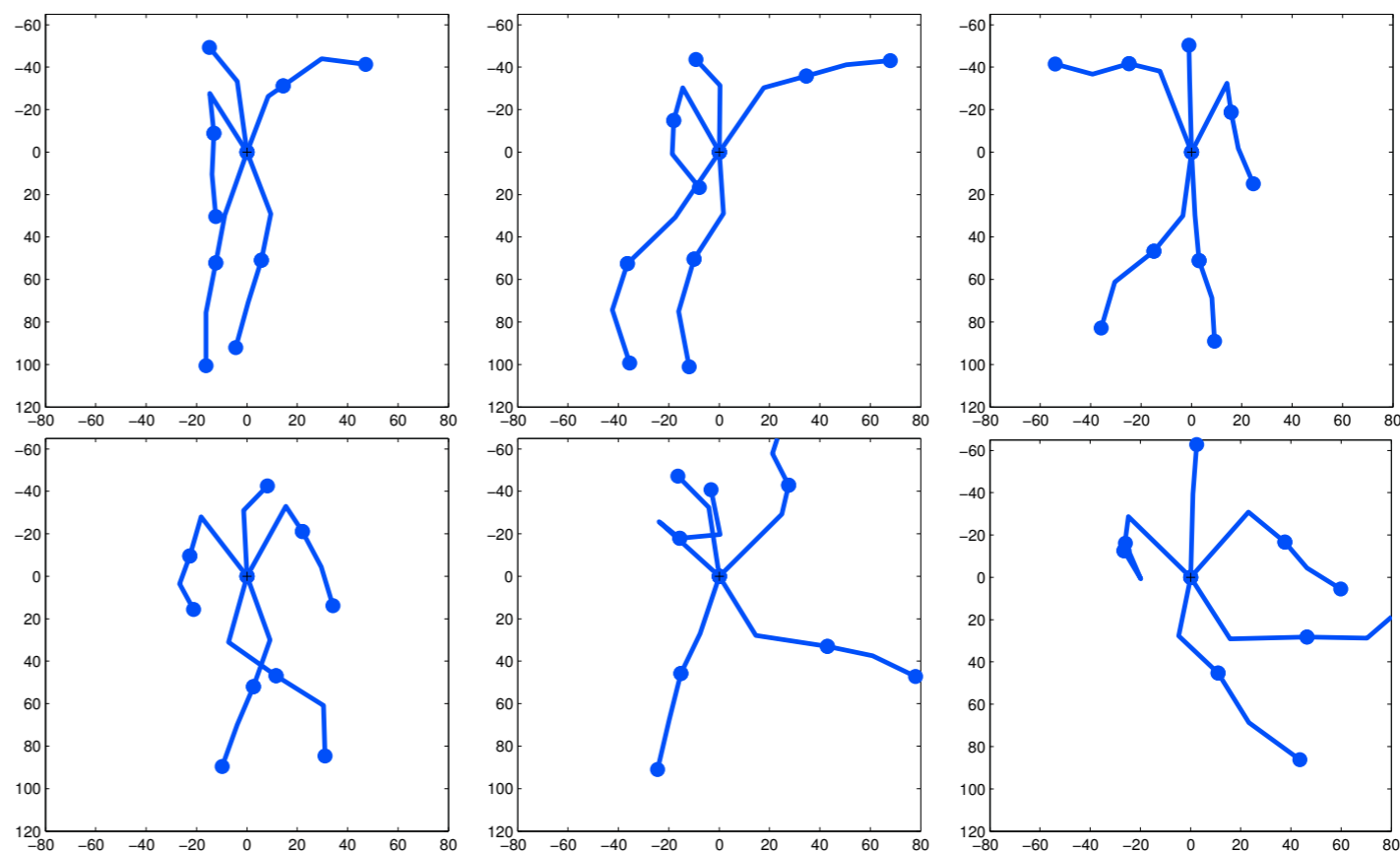
# Kinematic Tree Prior

- Prior parameters:  $\{T_{ij}, \Sigma^{ij}\}$
- Parameters of the prior are estimated with **maximum likelihood**

mean pose



several independent samples



# Evaluation Scenarios

1. Human Pose Estimation  
“People” dataset  
[Ramanan, NIPS’06]



2. Upper-body Pose Estimation  
“Buffy” dataset  
[Ferrari et al., CVPR’08]



3. Pedestrian Detection  
“TUD Pedestrians” dataset  
[Andriluka et al., CVPR’08]





# Evaluation Scenarios

1. Human Pose Estimation  
“People” dataset  
[Ramanan, NIPS’06]



2. Upper-body Pose Estimation  
“Buffy” dataset  
[Ferrari et al., CVPR’08]



3. Pedestrian Detection  
“TUD Pedestrians” dataset  
[Andriluka et al., CVPR’08]

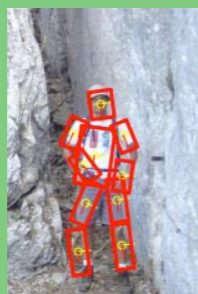




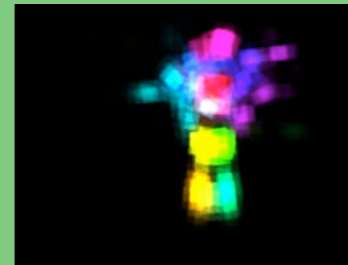
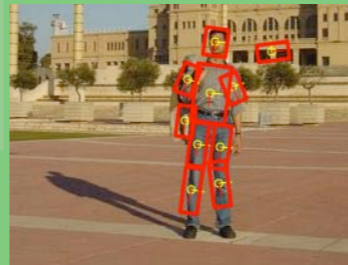
# Scenario 1: Qualitative Results

Our model

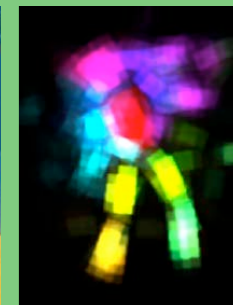
8/10



8/10

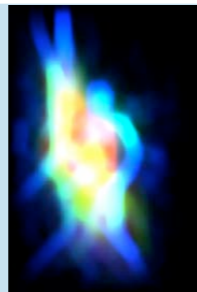


7/10

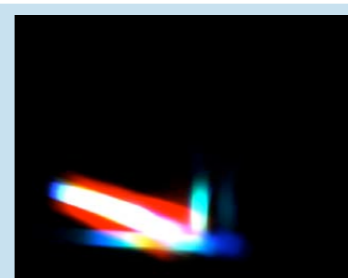
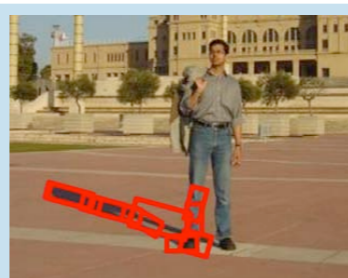


[Ramanan,  
NIPS'06]

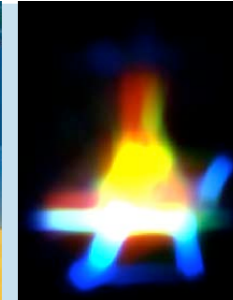
7/10



0/10



3/10



Our model

8/10



6/10

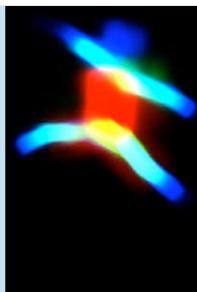


7/10

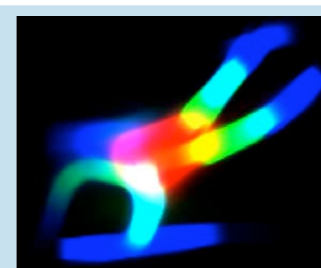


[Ramanan,  
NIPS'06]

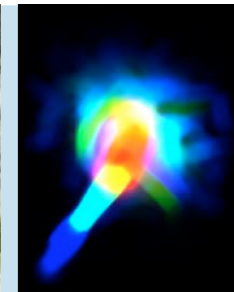
4/10



3/10



3/10





# Scenario 1: Quantitative Results

Method	Torso	Upper legs	Lower legs	Upper arm	Forearm	Head	Total
[Ramanan, NIPS'06] 2nd parse	52	30	29	17	13	37	27
Our inference, edge features from [Ramanan, NIPS'06]	63	48	37	26	20	45	37
Our part detectors (SC)	29	12	18	3	4	40	14
Our prior, our part detectors (SC)	<b>81</b>	<b>63</b>	<b>55</b>	<b>47</b>	<b>31</b>	<b>75</b>	<b>55</b>
Our prior, our part detectors (SIFT)	78	58	54	44	31	66	52



# Scenario 1: Quantitative Results

Method	Torso	Upper legs	Lower legs	Upper arm	Forearm	Head	Total
[Ramanan, NIPS'06] 2nd parse	52	30	29	17	13	37	27
Our prior, edge features from [Ramanan, NIPS'06]	63	48	37	26	20	45	37
Our part detectors (SC)	29	12	18	3	4	40	14
Our prior, our part detectors (SC)	81	63	55	47	31	75	55
Our prior, our part detectors (SIFT)	78	58	54	44	31	66	52





# Scenario 1: Quantitative Results

Method	Torso	Upper legs	Lower legs	Upper arm	Forearm	Head	Total
[Ramanan, NIPS'06] 2nd parse	52	30	29	17	13	37	27
Our inference, edge features from [Ramanan, NIPS'06]	63	48	37	26	20	45	37
<b>Our part detectors (SC)</b>	29	12	18	3	4	<b>40</b>	14
Our prior, our part detectors (SC)	<b>81</b>	<b>63</b>	<b>55</b>	<b>47</b>	<b>31</b>	<b>75</b>	<b>55</b>
Our prior, our part detectors (SIFT)	78	58	54	44	31	66	52

SC = Shape Context

# Scenario 1: Quantitative Results

Method	Torso	Upper legs	Lower legs	Upper arm	Forearm	Head	Total
[Ramanan, NIPS'06] 2nd parse	52	30	29	17	13	37	27
Our inference, edge features from [Ramanan, NIPS'06]	63	48	37	26	20	45	37
Our part detectors (SC)	29	12	18	3	4	40	14
Our prior, our part detectors (SC)	<b>81</b>	<b>63</b>	<b>55</b>	<b>47</b>	<b>31</b>	<b>75</b>	<b>55</b>
Our prior, our part detectors (SIFT)	78	58	54	44	31	66	52

SC = Shape Context



# Scenario 1: Quantitative Results

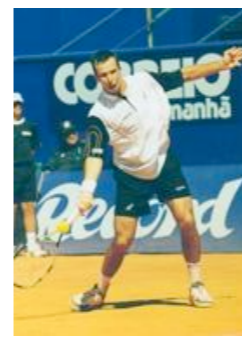
Method	Torso	Upper legs	Lower legs	Upper arm	Forearm	Head	Total
[Ramanan, NIPS'06] 2nd parse	52	30	29	17	13	37	27
Our inference, edge features from [Ramanan, NIPS'06]	63	48	37	26	20	45	37
Our part detectors (SC)	29	12	18	3	4	40	14
<b>Our prior, our part detectors (SC)</b>	<b>81</b>	<b>63</b>	<b>55</b>	<b>47</b>	<b>31</b>	<b>75</b>	<b>55</b>
Our prior, our part detectors (SIFT)	78	58	54	44	31	66	52

SC = Shape Context



# Evaluation Scenarios

1. Human Pose Estimation  
“People” dataset  
[Ramanan, NIPS’06]



2. Upper-body Pose Estimation  
“Buffy” dataset  
[Ferrari et al., CVPR’08]

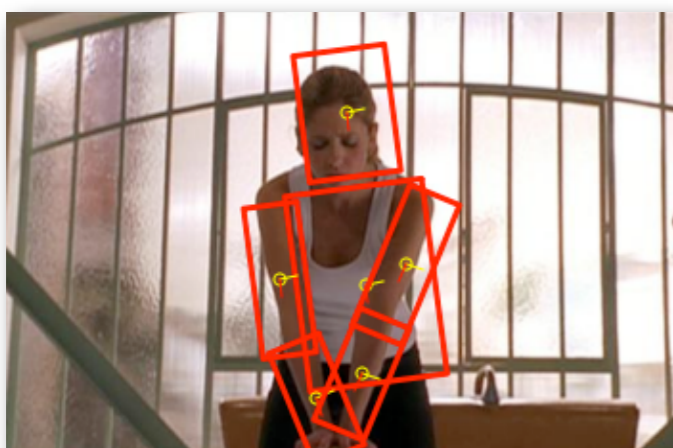
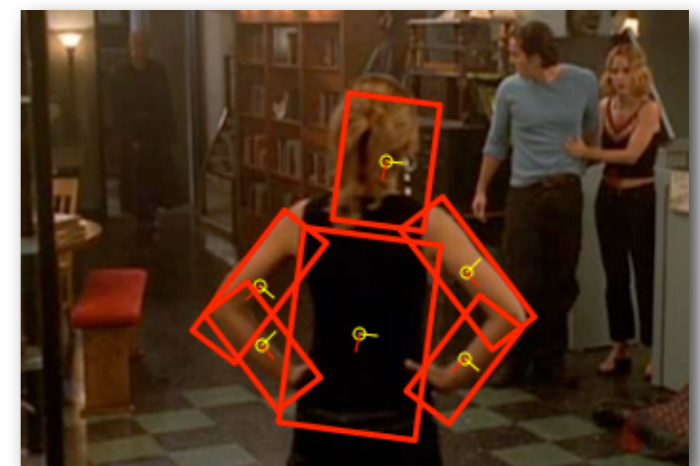
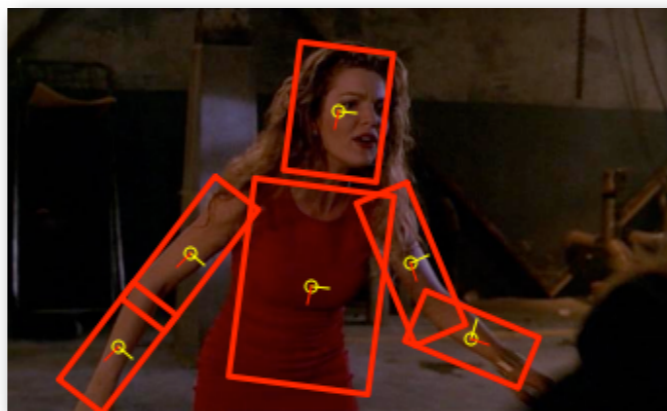
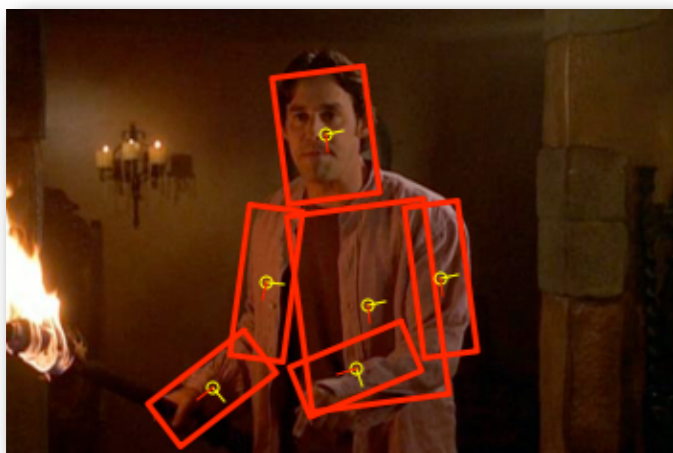
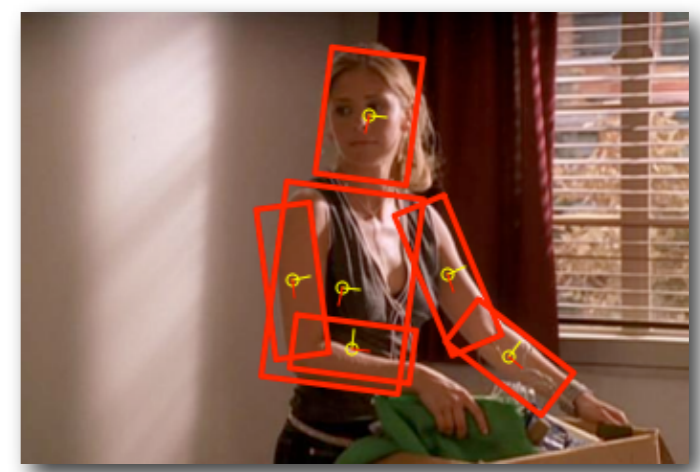
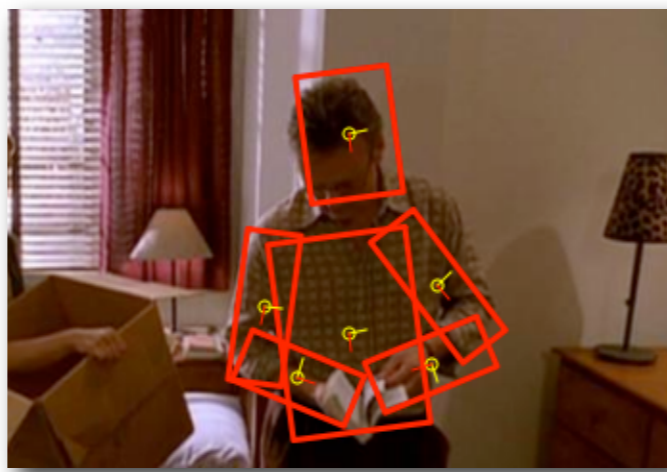
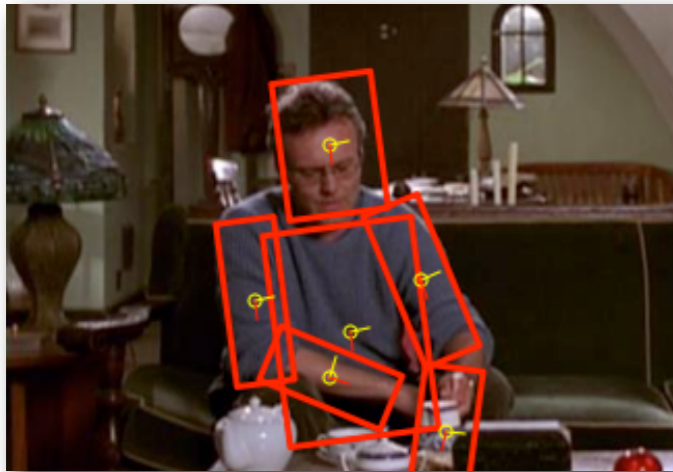


3. Pedestrian Detection  
“TUD Pedestrians” dataset  
[Andriluka et al., CVPR’08]



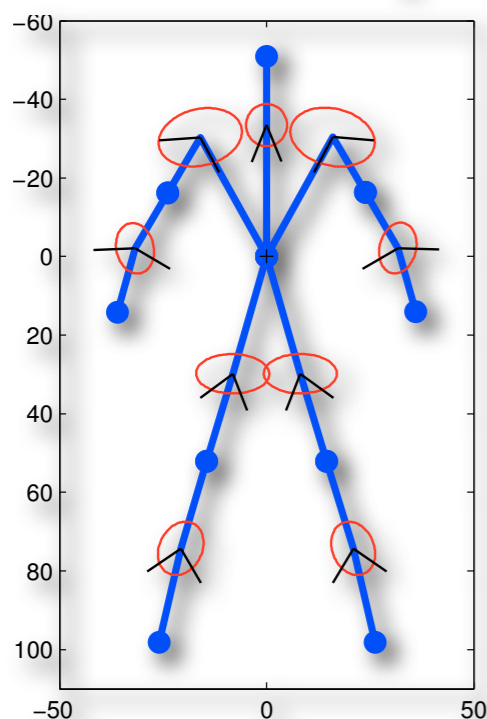


# Estimated upper-body poses



# Quantitative Results

Method	Torso	Upper arm	Lower arm	Head	Total
[Ferrari et al. CVPR'08]	-	-	-	-	57.9
detectors only	18.9	6.8	3.1	47.2	14.3
<b>full model</b>	<b>90.7</b>	79.3	41.2	<b>95.9</b>	71.3

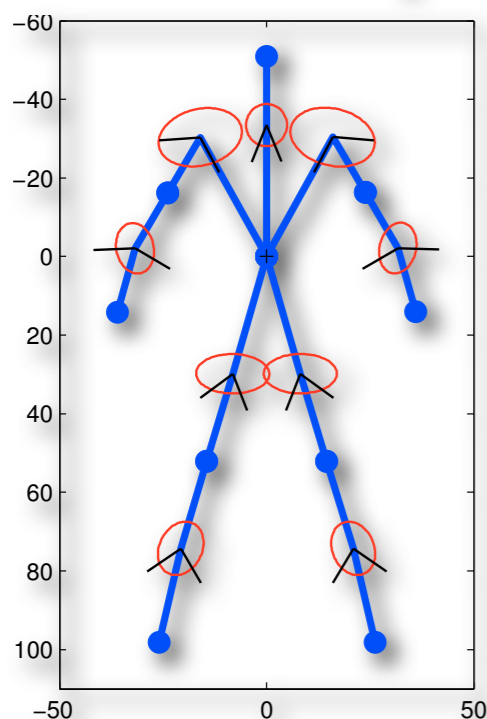


- generic model
- prior and appearance learned on the “People” dataset



# Quantitative Results

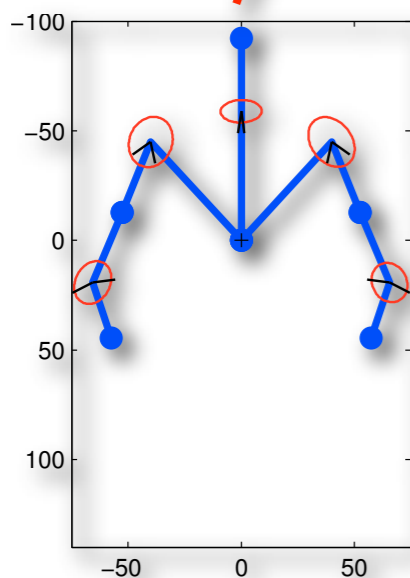
Method	Torso	Upper arm	Lower arm	Head	Total
[Ferrari et al. CVPR'08]	-	-	-	-	57.9
detectors only	18.9	6.8	3.1	47.2	14.3
full model	<b>90.7</b>	79.3	41.2	<b>95.9</b>	71.3
[Ferrari et al. CVPR'09]	-	-	-	-	<b>72.2</b>



- generic model
- prior and appearance learned on the "People" dataset

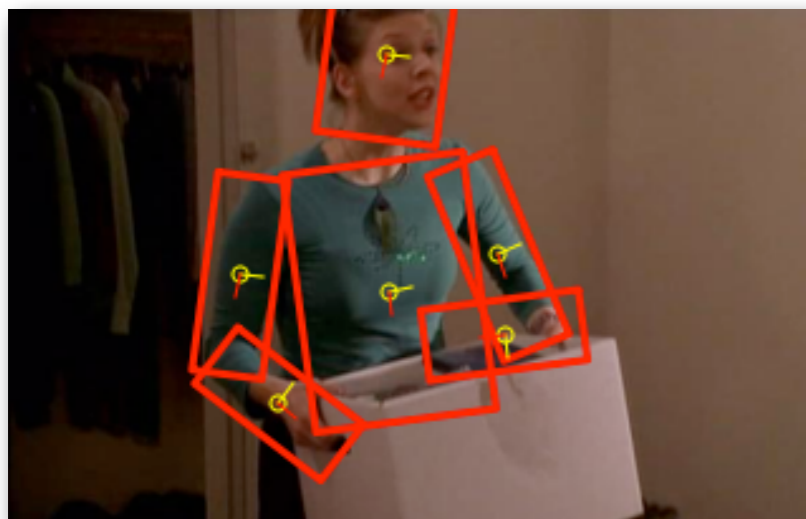
# Quantitative Results

Method	Torso	Upper arm	Lower arm	Head	Total
[Ferrari et al. CVPR'08]	-	-	-	-	57.9
detectors only	18.9	6.8	3.1	47.2	14.3
full model	<b>90.7</b>	79.3	41.2	<b>95.9</b>	71.3
[Ferrari et al. CVPR'09]	-	-	-	-	72.2
full model, Buffy pose prior	<b>90.7</b>	<b>81.35</b>	<b>46.5</b>	95.5	<b>73.5</b>

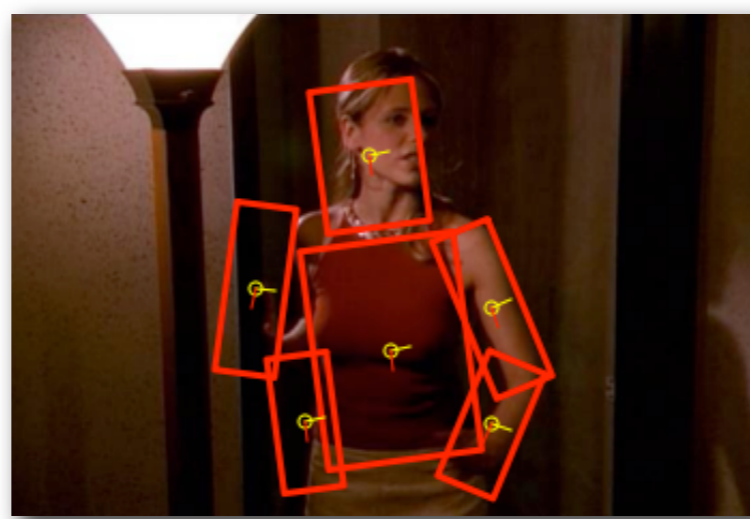


- specialized upper body prior
- appearance learned on the "People" dataset

# Typical Failure Cases



Foreshortening



Part occlusion



Detections on other  
body parts



# Evaluation Scenarios

1. Human Pose Estimation  
“People” dataset  
[Ramanan, NIPS’06]



2. Upper-body Pose Estimation  
“Buffy” dataset  
[Ferrari et al., CVPR’08]

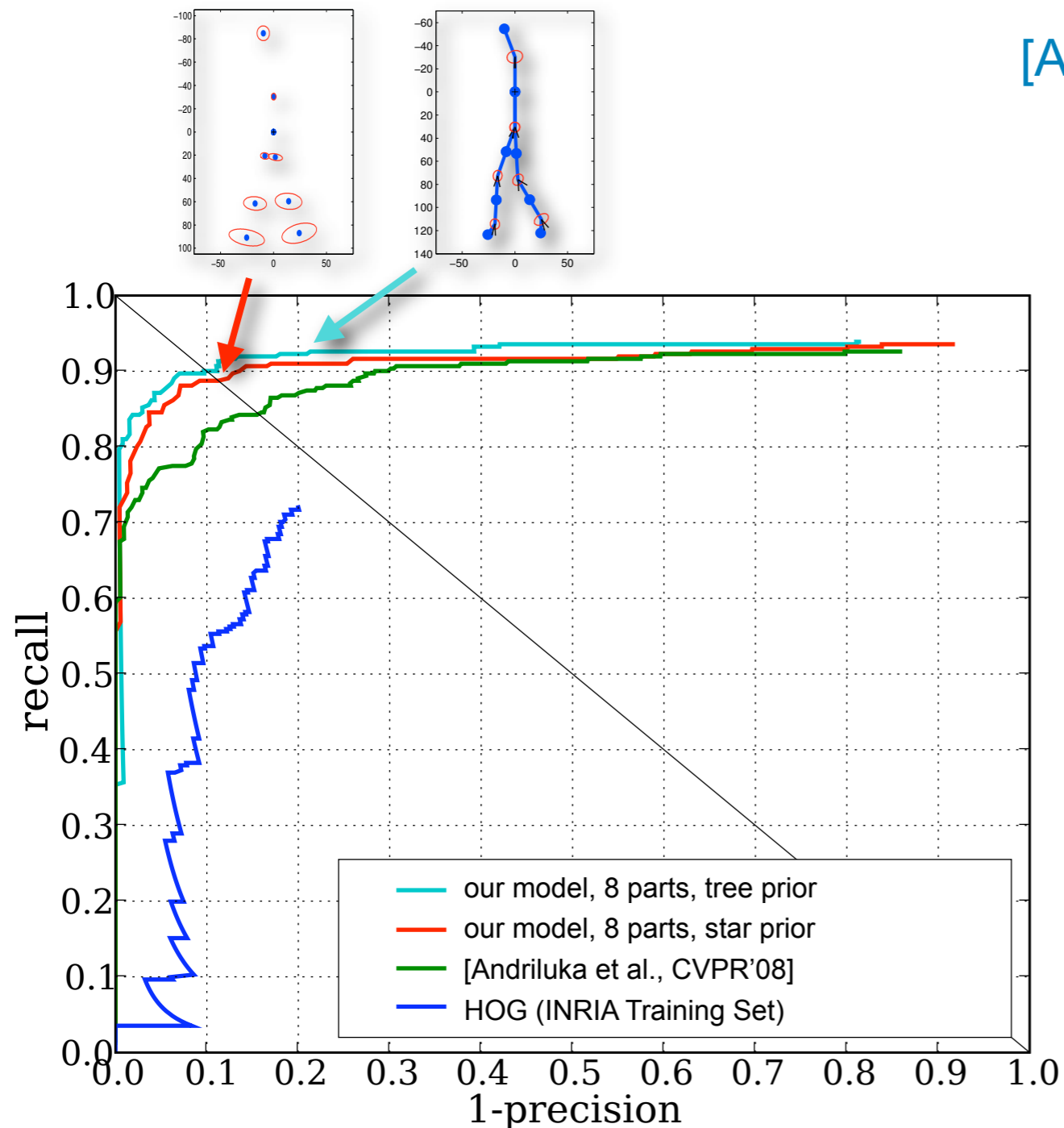


3. Pedestrian Detection  
“TUD Pedestrians” dataset  
[Andriluka et al., CVPR’08]



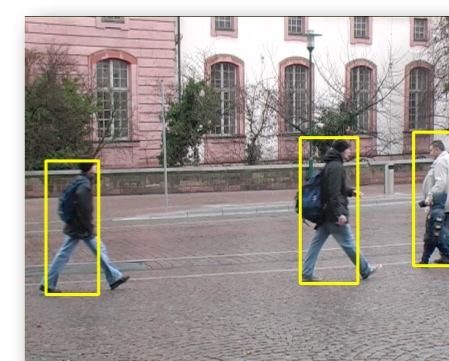
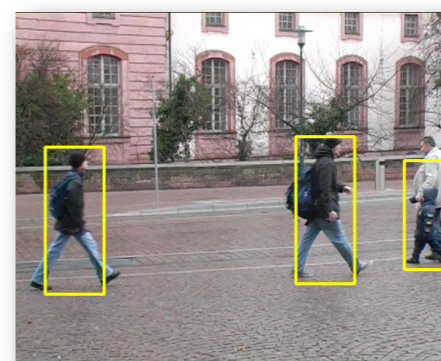
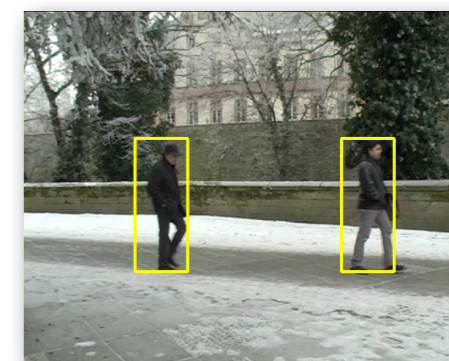
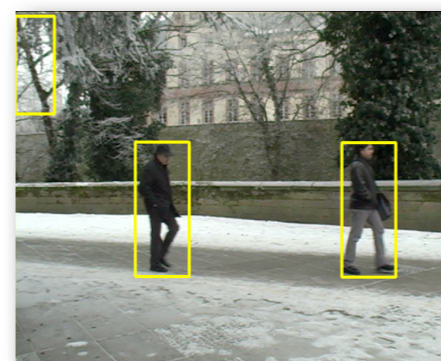
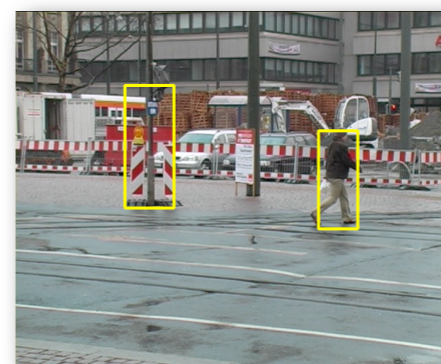
# People Detection: Results

- Comparison with state-of-the-art in people detection



[Andriluka et al., CVPR'08]

This work





# Conclusion & Future Work

- Success of **pose estimation by “body part” detection**
  - ▶ use well understood pose estimation framework (Pictorial Structures)
  - ▶ use appropriate representation for kinematic dependencies
  - ▶ use state of the art appearance representation (SIFT, SC) and classification (AdaBoost)
- Next steps:
  - ▶ estimate poses in 3D
  - ▶ part occlusions
  - ▶ appearance constraints between parts





# Thanks!

- **Acknowledgements:**
  - ▶ Thanks to Krystian Mikolajczyk for image descriptors code
  - ▶ Thanks to Christian Wojek for AdaBoost code and helpful suggestion
  - ▶ Thanks to Deva Ramanan and Vittorio Ferrari for making their code and datasets publicly available
  - ▶ This work is partially funded by German Research Foundation (DFG) through GRK 1362.
- Code and pre-trained models will be available at:
  - ▶ <http://www.mis.informatik.tu-darmstadt.de/code>