# Detecting Actions, Poses, and

# Objects wi

#### Chaitanya D









# Detecting Actions, Poses, and Objects with Relational Phraselets

Chaitanya Desai <u>UC Irvine</u>

Deva Ramanan



# K-way action classification



walking riding-bike jumping phoning taking-picture using computer

# K-way action classification



Requires bounding-box annotation on test images

#### What's wrong with K-way classification?



Detecting and localizing people performing actions is challenging

#### What's wrong with K-way classification?



Ignores the complexity of the ways in which people + objects can interact











Localize person (+ interacting object)







Localize person (+ interacting object)







Localize person (+ interacting object)
Classify action of each detected instance
Estimate pose of person (+ interacting object)

#### Challenge 1: human pose estimation



variation in appearance



variation in pose, viewpoint

### Challenge 2: person-object occlusions



Occluded person leg

# (Revised) action understanding



Localize
Estimate pose
Classify action
Estimate occlusions

## Related work: PASCAL Action Classification Challenge



Everingham et al 2011 Yao et al ICCV 11 Maji et al CVPR11

Few previous entries appear to output an explicit human skeleton Exceptions: next talk, Yang et al CVPR 10

# Our approach

Articulated pose estimation



Visual composites

Geometric parts

## Articulated pose estimation



Pictorial structures



Yang & Ramanan 11

Ioffe & Forsyth 01 Felzenswalb & Huttenlocher 05 Ferrari et al.08 Andruikula et al. 09 Johnson & Everingham 11







Models assume local appeara













Models assume local appeara Problem: T







#### Visual Phrases

Sadeghi and Fahardi, CVPR 11



Occluded leg not present in template



Person on horse

#### Visual Phrases

#### Sadeghi and Fahardi, CVPR 11







Person on jumping horse

Person on horse

Person standing next to horse

Problem: one may need lots of large composite templates

## Geometric parts (poselets)

Bourdev & Malik ICCV09 Maji et al CVPR11



## Geometric parts (poselets)

Bourdev & Malik ICCV09 Maji et al CVPR11



**Problem:** difficult to ensure that a globally-consistent arrangement of poselets will fire on a detection

# Approach

Articulated pose estimation



Visual composites

Geometric parts

#### Articulated models + visual composites

1. Define articulated model for person+object composite





Articulated models + visual composites + geometric parts

1. Define articulated model for person+object composite

2. Use local part mixtures ("phraselets") to capture different occlusion states



# Learning phraselets

#### Define phraselets as commonly-occuring geometric configurations

"Poselet-like" clusters



Given training data (with annotated landmarks), find clusters of landmark configurations relative to each joint

#### Clusters



#### Model occlusions with separate clusters



#### Visible left elbow

Occluded left elbow

Mixture label corresponds to visible/occlusion state

#### Local mixtures of phraselets



#### Local mixtures of phraselets



#### Local mixtures of phraselets



## Geometry-dependent parts





Part appearance (local mixture, denoted by color) depends on the location and appearance of other parts

## Inference & Learning



**Inference:** Infer part locations + mixtures with dynamic programming on trees

Learning: Tune linear parameters (including occlusion constraints) with SVM solver











An occluded mixture template may learn all 0 weights; let the learning algorithm decide!









An occluded mixture template may learn all 0 weights; let the learning algorithm decide!



2. Small patches are not as discriminative as larger templates (visual phrases / poselets)





An occluded mixture template may learn all 0 weights; let the learning algorithm decide!



2. Small patches are not as discriminative as larger templates (visual phrases / poselets)

Any connected set of phraselets can behave like a larger template (rigid s



## **Experimental Results**



We use 2010 & 2011 PASCAL Action Recognition benchmark We augment training and validation set with landmark annotations





#### Often produce meaningful poses



#### False-positive "phoning" detections

(taking a picture and scratching head have similar poses)

Articulated pose estimation

Action detection/ localization

Action classification

# Articulated pose estimation

#### Action detection/ localization

Action classification

On par with Poselets, comparable to state-of-the-art

# Articulated pose estimation

Considerably outperform pictorial structures

#### Action detection/ localization

Action classification

On par with Poselets, comparable to state-of-the-art

# Articulated pose estimation

Considerably outperform pictorial structures

#### Action detection/ localization

Action classification

On par with Poselets, comparable to state-of-the-art

## Action detection

Treat each action as an "object" and evaluate standard criteria (AP)



Significantly outperform state-of-art for detection





#### Action understanding

(detection, classification, & pose estimation)



# A look



#### Action understanding (detection, classification, & pose estimation)

#### Geometric part models

(interdependence of geometry and appearance)





#### **Riding Horse**







