



## Human Pose Estimation: Extension and Application

Digvijay Singh CVIT, IIIT Hyderabad Guide: Professor C V Jawahar



#### Content



- 1. Problem Definition and Challenges
- 2. Previous work on Human Pose Estimation: Classic to New
- 3. Extension: Human Pose Estimation in Videos
- 4. Application: Parsing Clothes in Images







### HPE: Problem Definition

OR



Locating Body Parts





IIT Hyc eral ad

#### Motivation



Most studied subject : HUMAN

Where is the human located? How many humans? What activity is the human performing? What gesture humans are making?



Activity: People toasting. 8 humans: 4 male, 4 female. Background: Club Interaction: Holding glasses



Activity: Boys playing football. 8 humans: 8 boys Background: Field Interaction: Kicking football



IIT Tyd erat ad

### Motivation for HPE



#### Detecting Cloth Segments<sup>1</sup>

## Understanding Human-Object<sup>2</sup> interaction and performed Activity





Yamaguchi et al, Parsing Clothing in fashion photographs.
 Yao et al, Modeling Mutual Context of Object and Human Pose in Human-Object Interaction Activites.



## Challenges for HPE



Clothing Variations Different clothing induces different silhouette. Different clothing covers/reveals different body regions.



Illumination Variation Outdoor illumination is highly variable.



#### Background Clutter

Background comprises of wide range of color distribution.





## Challenges for HPE





Occlusion: Self and External

Natural settings can have occlusions induced by different entities present.

#### Body part Foreshortening

Body parts projecting with a higher angle with respect to image plane have foreshortened projection on image plane.

Motion Blur (Video data specific) Quick moving body parts cause blur while capturing.





#### Our Contribution



#### **Extending Human Pose Estimation for Videos**

- Uses off-the-shelf pose estimator for images.
- Learns video-specific features using semi-supervised algorithm.
- Iterative self-training propagates correctness across sequence which is more robust than using contemporary tracking strategies.

#### **Applying HPE for Cloth Parsing in Images**

- By incorporating more robust pose information, the model is shown to perform more reliably in unrestricted settings.
- A peek at relevant information extraction from obtained results like obtaining cloth type/color patterns usually worn.



#### Content



- 1. Problem Definition and Challenges
- 2. Previous work on Human Pose Estimation: Classic to New
- 3. Extension: Human Pose Estimation in Videos
- 4. Application: Parsing Clothes in Images





#### 1. Histogram of Gradients (HoG)

Image divided into dense grid of uniformly spaced cells.

Gradient orientations are accumulated in all the grid cells.



#### HoG and Human Detection

Simple linear SVM trained using HoG features gives reliable pedestrian detector.



\* Dalal et al. Histograms of Oriented Gradients for Human Detection.



IIT Tyc eral



### Previous work on HPE: Classic

# 2. Poselets: Body Part Detectors trained with Clustering in configuration space

Keypoint based similarity distance:

$$D(a,b) = w_a \cdot \left\| |x_a - x_b| \right\|_2^2 \cdot (1 + v(a,b))$$

 $x_a, x_b$ : 3D coordinates of sample a, b; v(a, b): Penalizes visual dissimilarity



\* Bourdev et al, Poselets: Body Part Detectors trained with Clustering in configuration space.





#### Poselets (contd.)

**Training**: Poselet classifiers are linear SVM trained using poselet examples as positives and non-person patches as negatives. Features used are HoG.





**Inference**: All poselet classifiers used to get probability maps for different body parts. Probability maps are combined using optimal max-margin framework to get detection box.

**Pros:** Each part classifier independent of the presence/absence of other parts.

**Cons:** Fails to exploit holistic representation of human pose.



IIT

-lyc era



\* Fischler et al, The Representation and Matching of Pictorial Structures.





#### Pictorial Structure for Human Pose : YR Model



- *I*: Image
- $L_i$ : Location of part *i*

Slide Credits: Yang et al, Articulated Human Detection with Flexible Mixtures-of-Parts.





#### Pictorial Structure for Human Pose : YR Model



- $\alpha_i$ : Unary template for part *i*
- $\phi(I, l_i)$  : Local image features at location i





#### Pictorial Structure for Human Pose : YR Model



- $\psi(l_i, l_j)$  : Spatial features between  $l_i$  and  $l_j$
- $\beta_{ij}$ : Pairwise springs between parts *i* and *j*





#### Previous work on HPE: New

#### 4. DeepPose: Human Pose Estimation via Deep Neural Networks

For k keypoints, the architecture is designed to regress a vector of size 2k

LossValue(L<sub>2</sub>) = argmin 
$$\sum_{\theta \in I} ||y_i - \varphi(x_i, \theta)||_2^2$$



Parameters  $\varphi$  are fine-tuned using stochastic gradient descent while back-propagation.

\* Toshev et al, DeepPose: Human Pose Estimation via Deep Neural Networks.



era



#### Previous work on HPE: New

5. Articulated Pose Estimation by a Graphical Model with Image Dependent Pairwise Relations (IDPR)

**Assumption**: Local image features around a part can reliably detect the part as well as help in determining relative positions of all its neighbors.





IIT Tyd erat ad

#### Classic vs New: Comparison

11			
	Method	PCK(%)	Examples
FLIC dataset	YR	70.00	
	IDPR	91.14	
	Method	PCK(%)	Examples
LSP dataset	YR	55.81	
	IDPR	67.91	



#### Content



- 1. Problem Definition and Challenges
- 2. Previous work on Human Pose Estimation: Classic to New
- 3. Extension: Human Pose Estimation in Videos
- 4. Application: Parsing Clothes in Images







## Drawbacks of existing work

• Pairwise over-rides unary response to make false predictions.



IIT Hyd erab Id



## Drawbacks of existing work

- Pairwise over-rides unary response to make false predictions.
- Very high number of possible configurations for body parts that needs higher parameterization.
- Contemporary tracking methods are restricted by the amount of object movement and outlook change it can handle.





## Problem Formulation

• Video data has consistent appearance and slow changing pose configuration.







## Problem Formulation

- Video data has consistent appearance and slow changing pose configuration.
- While temporal consistency can be learnt from supervised data, modelling unseen video-specific appearance is difficult.
- We propose to augment off-the-shelf HPE models with additional parameters to encode local appearance.
- These additional parameters are learnt using confident pose estimations for a particular video.
- Final model is run on all the frames in the video.





Digvijay et al. Fine-Tuning Human Pose Estimations in Videos, In WACV16



IT lyd erab id







### Proposed Method

#### Modified Self-training Pipeline:







#### Exemplar-SVM

While training, each exemplar generates a classifier with exemplar as only positive and large amount of negatives not belonging to exemplar's class.



$$\Omega_E(\mathbf{w}, b) = ||\mathbf{w}||^2 + C_1 h(\mathbf{w}^T \mathbf{x}_E + b) + C_2 \sum_{\mathbf{x} \in \mathcal{N}_E} h(-\mathbf{w}^T \mathbf{x} - b)$$
$$h(\mathbf{x}) = \max(1-\mathbf{x}, 0) \text{ "hinge-loss"}$$

The method exploits:

- Effectiveness of discriminative detectors
- Explicit nearest-neighbor correspondence for instances from similar class.

\* Malisiewicz et al, Ensemble of Exemplar-SVMs for Object Detection and Beyond.





### Proposed Method

#### Modified Self-training Pipeline:







IIT

lyd

erab

ıd



### Proposed Method









### Proposed Method

Modified Self-training Pipeline:





IIT

Hyd erat id



## Pose Quality Ranking SVM

Based on the characteristics of feasible pose geometric configurations in current frame as well as local sequence of frames.





IIT



## Pose Quality Ranking SVM

Per frame criteria:

- Left-Right shoulder
- Left-Right hip
- Left-Right torso length

 $\begin{array}{c} x(3) < x(15) \\ x(10) < x(22) \\ t1 < d(3,10)/d(15,22) < t2 \end{array}$ 

Temporal neighborhood consistency criteria:

- Pose scale
- Left shoulder movement  $t1 \le d(x(3),nbd(x(3)) \le t2$
- Right hip movement

t1 < sc(i)/nbd(sc(i)) < t2t1 < d(x(3),nbd(x(3)) < t2

t1 < d(x(22),nbd(x(22)) < t2







#### Coarse-to-Fine Strategy for Exemplar Selection

To ensure samples picked in each iteration of self-training belong to different temporal segments.

Strategy:

- Temporal regions with mean pose quality score higher than mean pose quality score of full sequence are chosen.
   Coarsely k regions with highest mean output scores are picked.
  - Coarsely, k regions with highest mean output scores are picked.
- From each selected temporal region, maximum *j* exemplars with highest pose quality score are finally picked.





#### Ensemble of Exemplar SVMs for Semi-Supervised Self-Training



Exemplars picked using our strategy are converted into Exemplar-SVMs.

For full-pose, 26 parts synthesize 26 E-SVMs, with body-part as the only positive and random patches from non-human image as negatives.







#### Ensemble of Exemplar SVMs for Semi-Supervised Self-Training



Unary potential term is redefined as:

$$\sum_{i \in V} (\eta w_i^{t_i} + (1 - \eta) \hat{w}_i^{t_i}) \cdot \phi(I, p_i)$$

where  $\hat{w}_i^{t_i}$  are the (normalized) weights learnt from the exemplar SVM for the  $i^{th}$  part tuned for type  $t_i$ , and  $\eta$  is a parameter that controls the weight given to the exemplars' contribution to the unary score.



![](_page_36_Picture_0.jpeg)

![](_page_36_Picture_1.jpeg)

![](_page_36_Picture_2.jpeg)

#### Percentage of Correct Keypoints (PCK):

Estimated keypoints are correct that lie within a threshold distance  $D_t$  from its GT counterpart.

 $D_t = \beta . \max(h, w)$ 

#### **Datasets Used**

Dataset	Avg Frame length	Activities	Complexity	
VideoPose	30	TV shows	Easy	
Poses in the Wild	30	Movies	Medium	
CVIT-Sports-Videos	131	Sports	High	

![](_page_37_Picture_0.jpeg)

#### Datasets Used

![](_page_37_Picture_2.jpeg)

![](_page_37_Picture_3.jpeg)

![](_page_37_Picture_4.jpeg)

![](_page_37_Picture_5.jpeg)

![](_page_37_Picture_6.jpeg)

![](_page_37_Picture_7.jpeg)

#### Poses in the Wild

![](_page_37_Picture_9.jpeg)

![](_page_37_Picture_10.jpeg)

#### VideoPose

![](_page_38_Picture_0.jpeg)

![](_page_38_Picture_1.jpeg)

## **CVIT-Sports-Videos**

- 11 sports videos retrieved from YouTube.
- Activities included: cricket-bowling, cricket-batting and football-kicking.
- Full human pose (26 keypoints) labelling
- Total 1446 frames averaging to 131 frames per video.

![](_page_38_Picture_7.jpeg)

Half-body self occlusion

![](_page_38_Picture_9.jpeg)

Extreme body deformation

![](_page_39_Picture_0.jpeg)

![](_page_39_Picture_1.jpeg)

#### Quantitative Results

Method	PCK(%)		
YR	72.41		
FT @ phase2	74.31		
FT-Full	74.19		

![](_page_39_Figure_4.jpeg)

![](_page_39_Picture_5.jpeg)

Final phase 5 estimations

![](_page_40_Picture_0.jpeg)

![](_page_40_Picture_1.jpeg)

#### Quantitative Results

Method	PCK(%)		
YR	72.41		
FT @ phase2	74.31		
FT-Full	74.19		
FT-Full + PP	74.74		

PP, Post Processing: Neighborhood Interpolation  $est_i = mean(est_{i-1}, est_{i+1})$ 

![](_page_41_Picture_0.jpeg)

![](_page_41_Picture_1.jpeg)

#### Quantitative Results

Method	PCK(%)	
YR	72.41	
FT @ phase2	74.31	Full-body evaluation
FT-Full	74.19	
FT-Full + PP	74.74	

Upper-body evaluation	Method	<b>CVIT-Sports-Videos</b>	PIW	VP
	YR	64.87	70.74	63.35
	[1]	42.12	71.43	71.95
	YR + SiftFlow	31.75	48.04	60.01
	FT @ phase 2	65.20	72.78	69.11
	FT-Full	64.90	72.44	68.65
	FT-Full + PP	64.50	73.26	68.96

[1] Cherian et al, Mixing body-part sequences for Human Pose Estimation.

![](_page_42_Picture_0.jpeg)

#### Qualitative Results

![](_page_42_Picture_2.jpeg)

Top Row: YR Bottom Row: FT-Full

![](_page_42_Picture_4.jpeg)

![](_page_42_Picture_5.jpeg)

CVIT-Sports-Videos

![](_page_42_Picture_7.jpeg)

#### Poses in the Wild

![](_page_42_Picture_9.jpeg)

![](_page_43_Picture_0.jpeg)

### Discussions

![](_page_43_Picture_2.jpeg)

- Proposed self-training methodology that empowers unary response and captures more intricate pose configurations.
- Introduced new CVIT-Sports-Videos dataset having 11 videos from sports domain.
- Novel pose quality scoring criteria that helps in selecting instances for each iteration of self-training.
- Quantitative and qualitative results show we surpass previous state-of-the-art based on pictorial structures.
- **Parameter Selection:** The value of biasing parameter  $\eta$  that controls the contribution from base model and exemplars is determined automatically using pose quality score in a manner determined by cross validation.

![](_page_44_Picture_0.jpeg)

#### Content

![](_page_44_Picture_2.jpeg)

- 1. Problem Definition and Challenges
- 2. Previous work on Human Pose Estimation: Classic to New
- 3. Extension: Human Pose Estimation in Videos
- 4. Application: Parsing Clothes in Images

![](_page_44_Picture_7.jpeg)

![](_page_45_Picture_0.jpeg)

IIT Hyd

erat Id

### **Application: Clothing Parsing**

![](_page_45_Picture_2.jpeg)

Simple Image

![](_page_45_Picture_4.jpeg)

Complex Occlusion

![](_page_45_Picture_6.jpeg)

![](_page_45_Picture_7.jpeg)

![](_page_45_Picture_8.jpeg)

![](_page_45_Picture_9.jpeg)

Patterns

![](_page_45_Picture_10.jpeg)

![](_page_45_Picture_11.jpeg)

![](_page_45_Picture_12.jpeg)

![](_page_45_Picture_13.jpeg)

![](_page_45_Picture_14.jpeg)

![](_page_45_Picture_15.jpeg)

![](_page_45_Picture_16.jpeg)

![](_page_45_Figure_17.jpeg)

![](_page_46_Picture_0.jpeg)

## Application: Clothing Parsing

![](_page_46_Picture_2.jpeg)

Understanding the dependence of obtaining cloth segments and determining cloth types on **pose information**. Settings are divided into two categories:

![](_page_46_Picture_4.jpeg)

#### Restricted

- Front-facing single human present in the center of the image with similar scale and width.
- Pose configurations are very generic and comprehensible.
- Almost all human body parts are clearly visible
- Dataset: Fashionista

![](_page_46_Picture_10.jpeg)

#### Unrestricted

- No restrictions on orientation, scale or number of humans present in the image.
- Pose configuration composition can vary from very simple to very intricate.
- No restriction on missing human body parts.
- Dataset: H3D

![](_page_47_Picture_0.jpeg)

## Application: Clothing Parsing

![](_page_47_Picture_2.jpeg)

Understanding the dependence of obtaining cloth segments and determining cloth types on **pose information**. Settings are divided into two categories:

![](_page_47_Picture_4.jpeg)

Restricted

More assumptions.

Feasibility of a more regularized model.

Pictorial Structure based models like YR model work reliably well under such assumptions.

Tackled by Yamaguchi et al.

![](_page_47_Picture_10.jpeg)

Unrestricted

Fewer assumptions.

More random variables to keep account of.

Poselets (body part detector) works better in such settings.

Our problem.

![](_page_48_Picture_0.jpeg)

#### **Robust Pose Information**

![](_page_48_Picture_2.jpeg)

![](_page_48_Picture_3.jpeg)

![](_page_48_Picture_4.jpeg)

Pictorial Structure based models underperform in natural settings. Poselets manage to approximate visible body-parts.

#### Given prior information, certain assumptions can be made

Cloth & body part correlation

![](_page_48_Picture_8.jpeg)

![](_page_48_Picture_9.jpeg)

Pants-Legs

Cloth & appearance correlation

![](_page_48_Picture_12.jpeg)

Skin

Cloth & position correlation

![](_page_48_Picture_15.jpeg)

IIT Ιус eral ιd

![](_page_49_Picture_0.jpeg)

### Application: Parsing Clothes in Unrestricted Images

![](_page_49_Picture_2.jpeg)

#### Features

(1) RGB Histogram (2) CIE L\*a\*b Histogram (3) Gabor filter response (4) Normalized x-y coordinates histogram
(5) Normalized x-y coordinates w.r.t body joints.

![](_page_49_Picture_5.jpeg)

$$L^* = \min_{L} \sum_{i \in V} \Phi(l_i \mid Z, \phi(s_j, Z)) + \lambda_1 \sum_{(i,j) \in E} \psi_1(l_i, l_j) + \lambda_2 \sum_{(i,j) \in E} \psi_2(l_i, l_j \mid \psi(s_i, s_j, Z))$$

where  ${\bf Z}$  is the human pose information,  ${\bf I}$  is the image and

#### Segmentation

![](_page_49_Picture_9.jpeg)

$$\psi(s_i, s_j, Z) \equiv [(\phi(s_i, Z) + \phi(s_j, Z))/2, |(\phi(s_i, Z) - \phi(s_j, Z))/2|]$$

- $\phi$  is the feature vector
- Ψ<sub>1</sub> is the pairwise potential modeling the label cooccurrence

Learning: Logistic regression

- $\Phi$  is the unary potential modeling the appearance
- $\Psi_2 \quad \begin{array}{l} \text{is the pairwise potential} \\ \text{modeling the appearance} \\ \text{dependent co-occurrence} \end{array}$

Inference: Loopy belief propagation

![](_page_50_Picture_0.jpeg)

#### Datasets Used

![](_page_50_Picture_2.jpeg)

![](_page_50_Picture_3.jpeg)

Fashionista

H3D

IIT Iyd erab id

![](_page_51_Picture_0.jpeg)

## Qualitative Results

![](_page_51_Picture_2.jpeg)

![](_page_51_Picture_3.jpeg)

![](_page_51_Picture_4.jpeg)

![](_page_51_Picture_5.jpeg)

![](_page_51_Picture_6.jpeg)

![](_page_51_Picture_7.jpeg)

![](_page_51_Picture_8.jpeg)

![](_page_51_Picture_9.jpeg)

#### Fashionista

![](_page_51_Picture_11.jpeg)

![](_page_52_Picture_0.jpeg)

#### Quantitative Results

![](_page_52_Picture_2.jpeg)

Method	F	Full-a		Full-m		Unary	
	Pixel acc	mAGR	Pixel acc	mAGR	Pixel acc	mAGR	
[1]	$61.0 \pm 5.0$	$34.9\pm3.9$	$49.9 \pm 4.5$	$39.9 \pm 4.8$	$49.5\pm4.3$	$39.6 \pm 4.4$	
Ours	$74.5 \pm 4.7$	$49.7 \pm 3.8$	$68.5 \pm 5.4$	$55.2 \pm 4.5$	$68.4 \pm 5.4$	$54.8 \pm 4.3$	
Ours+Noc	$2 \parallel 77.4 \pm 4.0$	$57.0 \pm 3.3$	$70.2 \pm 5.2$	$63.1\pm4.5$	$70.1 \pm 5.2$	$62.4 \pm 4.3$	
	H3D						
Method	Ful	Full-a		Full-m		Unary	
	Pixel acc	mAGR	Pixel acc	mAGR	Pixel acc	mAGR	
[1]	$89.0\pm0.8$	$63.4 \pm 1.5$	$88.3\pm0.8$	$69.6 \pm 1.7$	$88.2\pm0.8$	$69.8 \pm 1.8$	
Ours	$87.7\pm0.8$	$62.7 \pm 1.8$	$86.0\pm1.0$	$70.2\pm2.0$	$85.9 \pm 1.1$	$70.4\pm2.0$	
Fashionista							

Fasmonista

![](_page_53_Picture_0.jpeg)

IIT Hyc

![](_page_53_Picture_1.jpeg)

## Clothing Pattern Mining

Cloth Co-occurrence

- Upper-body and Lower-body cloth occurrences considered.
- Outermost cloth as the representative cloth label.
- Top results: Skirt-top, Shorts-top, Blouse-skirt, Tights-dress and Cardigan-dress

![](_page_53_Picture_7.jpeg)

![](_page_54_Picture_0.jpeg)

IIT Hyc

![](_page_54_Picture_1.jpeg)

## Clothing Pattern Mining

Color Co-occurrence

- Upper-body and Lower-body cloth color occurrences considered.
- Vector quantized dominant color of the representative cloth is color label.
- Top results (upper:lower): blue:blue, red:red, red:blue, blue:red, white:blue

![](_page_54_Picture_7.jpeg)

![](_page_55_Picture_0.jpeg)

#### Discussions

![](_page_55_Picture_2.jpeg)

- Cloth parsing is formulated by incorporating more robust pose information for unrestricted settings.
- Proposed work outperforms previous state-of-the-art that uses Pictorial Structure based model.
- Qualitative observation shows that results are not reliable.
- The problem needs more breaking down. Even very accurate pose estimations can aid to an extent. More abstract information need to be parameterized like fine-grained cloth attributes, shape and layering order.

![](_page_55_Picture_7.jpeg)

![](_page_56_Picture_0.jpeg)

### **Related Publications**

- Digvijay Singh, Vineeth Balasubramanian, and C V Jawahar. Fine-Tuning Human Pose Estimations in Videos. In WACV, 2016.
- Nataraj Jammalamadaka, Ayush Minocha, Digvijay Singh, and C V Jawahar. Parsing Clothes in Unrestricted Images. In BMVC, 2013.