



**Master in Artificial Intelligence (UPC-URV-UB)**

# Tri-modal Human Body Segmentation

Master of Science Thesis

Cristina Palmero Cantariño

Advisor: Sergio Escalera Guerrero



February 6, 2014





































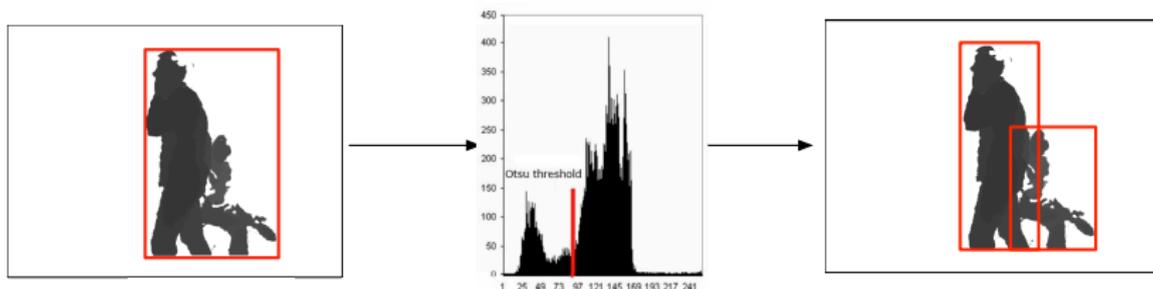


# Extraction of regions of interest

## Bounding box generation from regions of interest

People overlap:

- Bimodal disparity distribution.
- Otsu's threshold to split regions.



# Extraction of regions of interest

Bounding box transformation and correspondence to other modalities

All modalities must have the same number of bounding boxes, corresponding to the same regions of interest.

Tasks:

- 1 Find correspondence between rgb/depth and thermal regions of interest.
- 2 Compute the corresponding bounding boxes in thermal modality generated after applying Otsu's threshold in depth modality.

# Extraction of regions of interest

Bounding box transformation and correspondence to other modalities

- 1 Find correspondence between rgb/depth and thermal regions of interest.
  - Iterative search among depth and thermal modalities.
  - Takes into account deviation among them.
  - Best match: bounding box coordinates, amount of overlap and area similarity.
  - Correspondence function:

$$b_{iq}^{\text{thermal}} = \beta(b_{ij}^{\text{depth}}) \quad (1)$$

where  $b_{ij}$  is the  $j$ -th bounding box in frame  $i$ .

# Extraction of regions of interest

Bounding box transformation and correspondence to other modalities

2 Compute the corresponding bounding boxes in thermal modality generated after applying Otsu's threshold in depth modality.

- Assuming bounding boxes of both rgb/depth and thermal modalities are proportional, find the equivalence ratio to create the split bounding boxes in thermal.
- Ratio  $k$ :

$$k_h = \frac{h_{b_{ij}^{\text{depth}}}}{h_{b_{iq}^{\text{thermal}}}}, k_w = \frac{w_{b_{ij}^{\text{depth}}}}{w_{b_{iq}^{\text{thermal}}}} \quad (2)$$

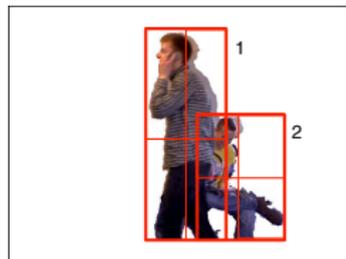
where  $h$  and  $w$  are the size of a given bounding box.

# Extraction of regions of interest

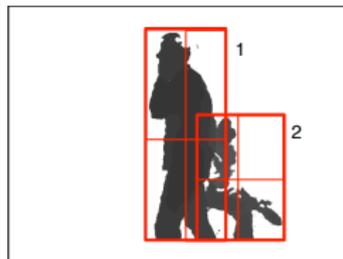
Bounding box transformation and correspondence to other modalities

Result:

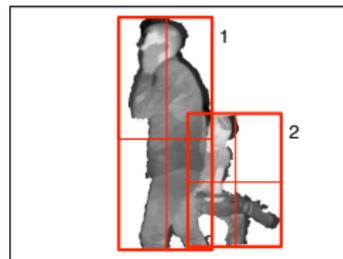
- Correspondence of regions of interest among modalities.
- Grid partitioning  $2 \times 2$  cells per bounding box.



RGB



Depth



Thermal

# Extraction of regions of interest

## Bounding box ground truth generation

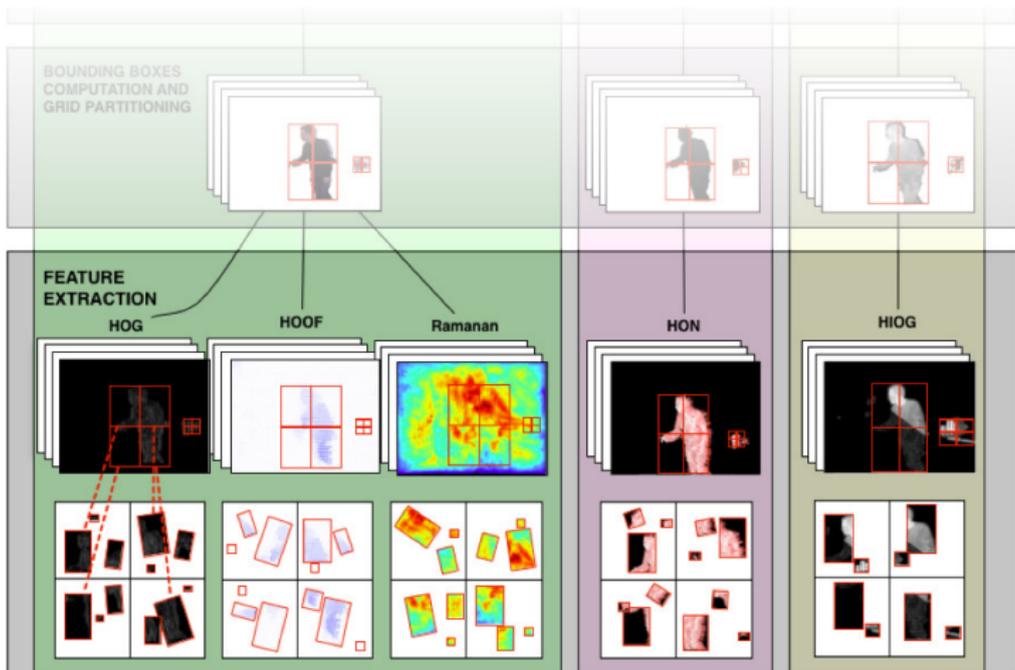
Comparing overlap between:

- Bounding boxes extracted from Ground Truth Masks
- Bounding boxes extracted from Background Subtraction Masks

Label:

$$t_r^d = \begin{cases} 0 & \text{(Object)} & \text{if overlap} \leq 0.1 \\ -1 & \text{(Unknown)} & \text{if } 0.1 < \text{overlap} < 0.6 \\ 1 & \text{(Subject)} & \text{if overlap} \geq 0.6 \end{cases}$$

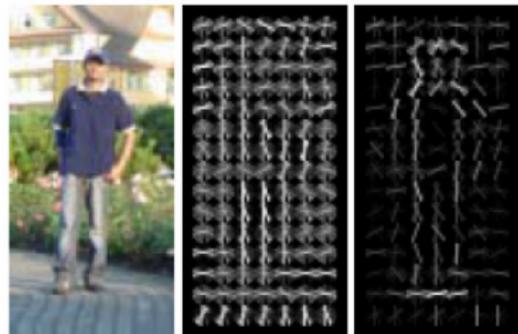
# Feature extraction



# Feature extraction: Color

## Histogram of Oriented Gradients (HOG)

- Unsigned gradients (0 - 180 degrees).
- 9-bin histogram.
- Contribution to the histogram given by the vector magnitude.
- No block overlap applied.
- Final vector of 288 values per cell.



HOG

---

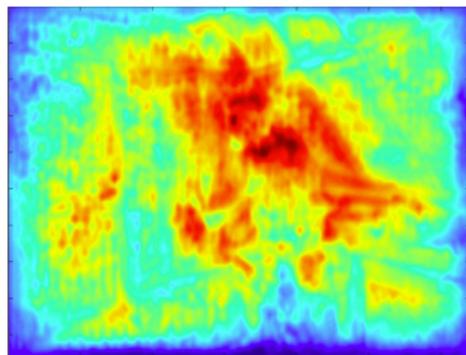
Navneet Dalal and Bill Triggs. "Histograms of oriented gradients for human detection". In: *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*. Vol. 1. IEEE, 2005, pp. 886–893.



# Feature extraction: Color

## Score Maps (SM)

- Score map based on Gabor filters.
- $C = 6$  component filters per body part.
- $M = 26$  body parts.
- $L$  scales per image.



Score maps from Ramanan et al

$$score(p_l) = \frac{1}{C} \frac{1}{M} \sum_{c \in C} \sum_{m \in M} score(p_l)_c^m \quad (3)$$

$$score(p) = \frac{1}{L} \sum_{l \in L} score(p_l)' \quad (4)$$

---

Yi Yang and Deva Ramanan. "Articulated pose estimation with flexible mixtures-of-parts". In: *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*. IEEE. 2011, pp. 1385–1392. [▶](#) [◀](#) [⏪](#) [⏩](#) [☰](#) [🔍](#) [🔄](#)





# Classifiers overview

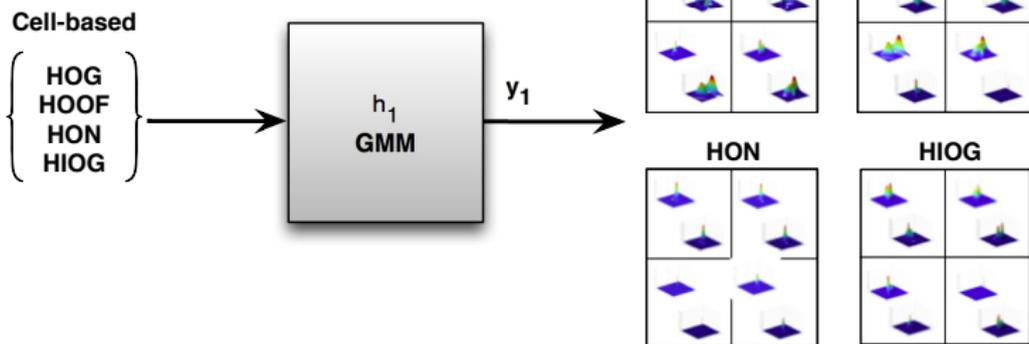
- 1 Statistical Learning:
  - Gaussian Mixture Models (GMM)
  - Subject and object probabilities
  - Individual prediction
- 2 Multi-modal fusion approaches:
  - Naive approach
  - Discriminative classifiers
  - Stacked learning fashion

# Cell classification

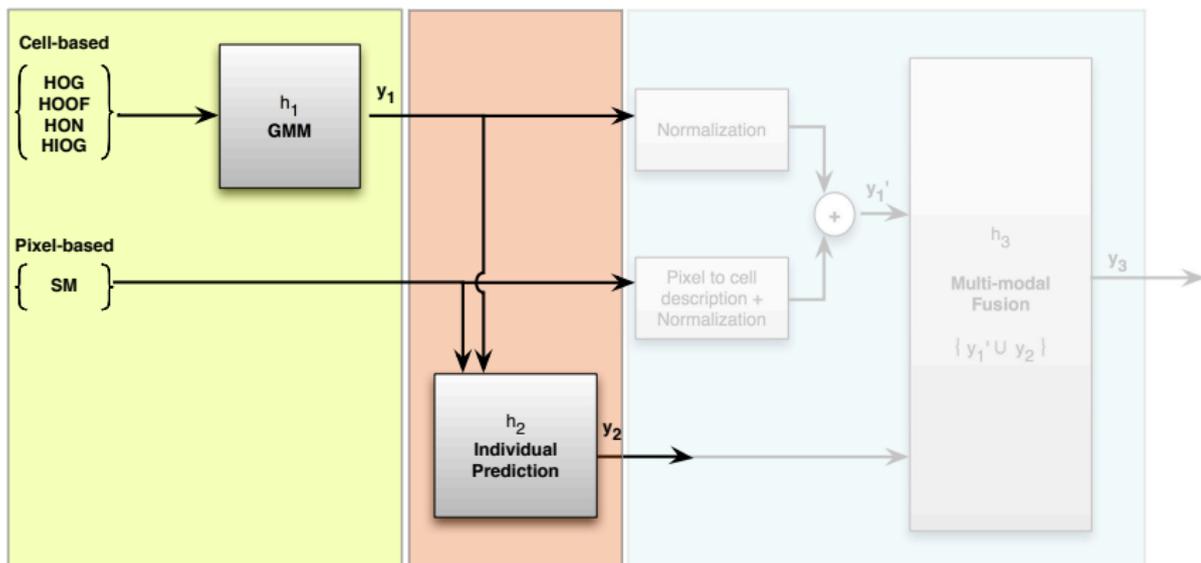
## Gaussian Mixture Models

- Unsupervised learning method for fitting multiple Gaussians to a set of multi-dimensional data points to obtain a likelihood  $\mathcal{L}$ .
- Trained using Expectation Maximization algorithm.

$$\mathcal{L} = \prod_{\mathbf{x} \in \mathbf{X}} \prod_{k=1}^K p(\mathbf{x}|k) P(k) \quad (5)$$



# Individual Prediction



# Individual Prediction

## Cell-based

Predict if a region corresponds to subject or object, for each cell-based descriptor individually. Grid cell voting  $v$ :

$$v = \sum_{i,j} \mathbb{1}\{\mathcal{L}_{ij}^{d,\text{sub}} > \mathcal{L}_{ij}^{d,\text{obj}}\} \quad (6)$$

Based on a threshold  $v_{thr}$  that defines the minimum number of positive votes needed to assign the subject label to the given region:

$$v_{thr} = \frac{v_{\text{grid}} h_{\text{grid}}}{2} \quad (7)$$

Final decision  $\hat{t}_r^d$ :

$$\hat{t}_r^d = \mathbb{1}\left\{v > v_{thr}\right\} \vee \left\{\mathbb{1}\{v = v_{thr}\} \cdot \mathbb{1}\left\{\sum_{i,j} (\hat{\mathcal{L}}_{ij}^{d,\text{sub}} - \hat{\mathcal{L}}_{ij}^{d,\text{obj}}) > 0\right\}\right\} \quad (8)$$





# Multi-modal fusion

## Naive approach

- 1 Voting among all descriptors using individual predictions  $\hat{t}_r^d$ .
- 2 If there is a strong agreement between descriptions, those descriptions that differ are not taking into account in the third step.
- 3 Cell level fusion:

$$\bar{\mathcal{L}}_{ij}^{d,\text{sub}} = \sum_{d \in \mathcal{D}'} \hat{\mathcal{L}}_{ij}^{d,\text{sub}}, \quad \bar{\mathcal{L}}_{ij}^{d,\text{obj}} = \sum_{d \in \mathcal{D}'} \hat{\mathcal{L}}_{ij}^{d,\text{obj}} \quad (10)$$

- 4 Predict  $\hat{t}_r$  following the same procedure as in individual prediction.

# Multi-modal fusion

## SVM-based approach

- Discriminative supervised binary classifier that learns a model which represents the instances as points in space, mapped in such a way that instances of different classes are separated by a hyperplane in a high dimensional space.
- Approaches:
  - Simple:  $\{\hat{\mathcal{L}}_{ij}^{d,\text{sub}}, \hat{\mathcal{L}}_{ij}^{d,\text{obj}}\}$
  - Stacked:  $\{\hat{\mathcal{L}}_{ij}^{d,\text{sub}}, \hat{\mathcal{L}}_{ij}^{d,\text{obj}}, \hat{t}_r^d\}$



















