

# *Human Behavior understanding in cars and around*

Prof. Ing. Rita Cucchiara

Imagelab,  
Dipartimento di Ingegneria «Enzo Ferrari»,  
Università di Modena e Reggio Emilia,

Director of the  
Research Center in ICT, Softech-ICT  
Modena, Italy

Invited talk  
Venezia , 4/11/2016  
@ECLT Ca'Foscari



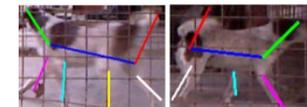
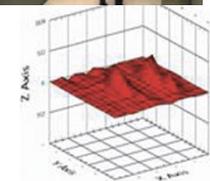
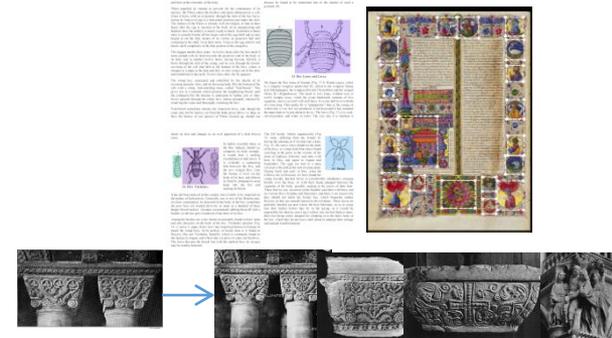
- UNIMORE Università di Modena e Reggio Emilia, Italy
- funded in 1175, in 2 cities:
- Modena (11 Departments, 14.100 students)
- Reggio Emilia (3 Departments, 5.400 students)



- Dipartimento di Ingegneria «Enzo Ferrari», Modena
- 6.000 students, c.a.100 faculties, 7 curricula
- Mechanical Engineering, Vehicle Engineering, Material engineering
- Computer Engineering, Electronic Engineering,
- Civil and Environmental Engineering
  - 2 International Phd Curricula (High Mechanics, **ICT**)
  - 2 High Technology Network Centres: Intermech and Softech-ICT
  - 2 Master in ICT: Vision Learning and Multimedia Technology (MUMET 2017), Cybersecurity Academy

- Pattern recognition and Image processing
- Medical Imaging
- Digitalized Document analysis
- **Multimedia**
- Multimedia big data annotation
- Video captioning
- **2D, 3D, wearable Computer vision**
- Augmented experiences in culture and museums
- Experience with Wearable devices, floors and IoT
- **Computer vision for Behavior analysis**
- Children and people behavior analysis
- Surveillance (in crowd)
- Automotive driver behavior understanding

[www.imagelab.unimore.it](http://www.imagelab.unimore.it)

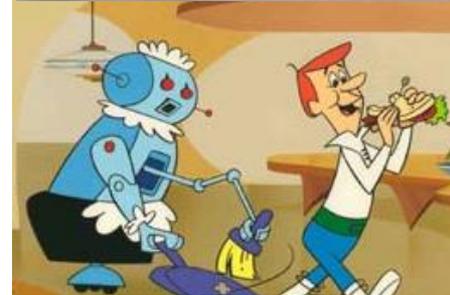




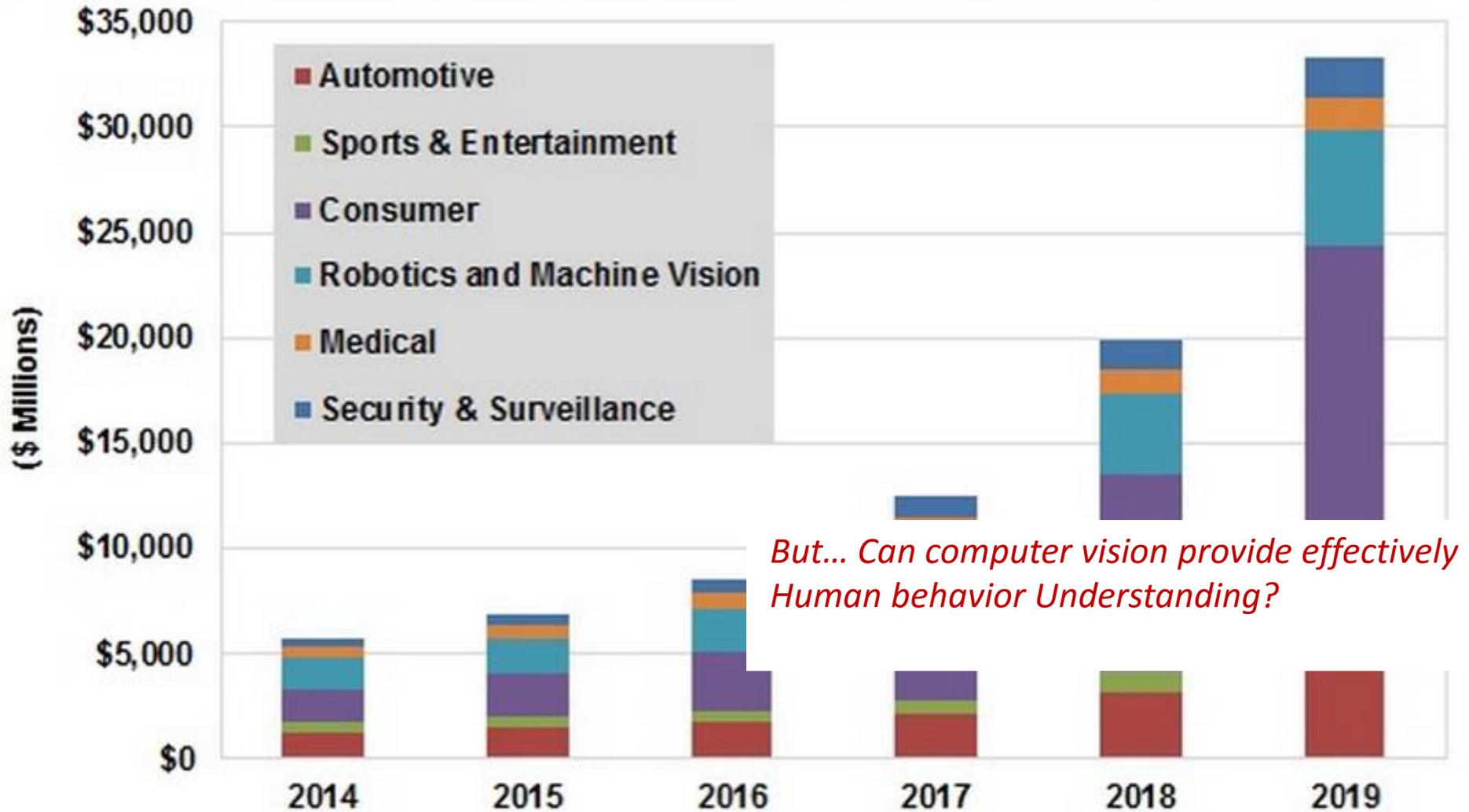
Computer Vision and Human Behavior Understanding  
*in cars and around*

- **Computer Vision** is the scientific discipline studying **how to perceive and understand the world through visual data by computers.**
- **Pattern Recognition** is the ensemble of theories, models, techniques **to recognize patterns in unknown or unordered data, generally images and multimedia content.**
- **Machine Learning** is the science of **getting computers to act without being explicitly programmed**

**Computer Vision & Machine Learning** *represent now the most advanced frontier of Artificial Intelligence studies*



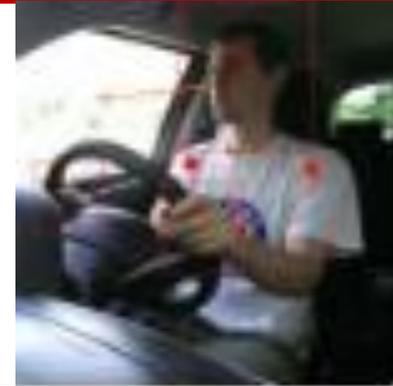
## Computer Vision Revenue by Vertical Market, World Markets: 2014-2019



*But... Can computer vision provide effectively Human behavior Understanding?*

- What he is doing?
- What are they doing?
- Single and collective behavior
- Collaborative or not collaborative behaviors

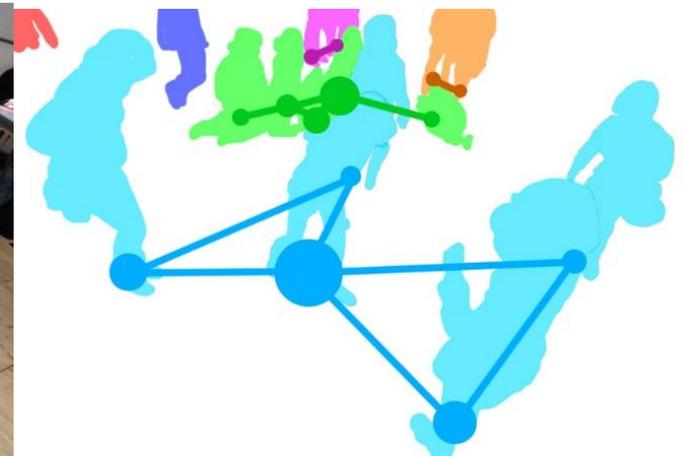
MIUR-FAR project  
DriveAttention  
2016-2018



PRIN project  
"2106-2018



project "VAEX 2015-2018



MIUR project "the educating city" cluster smart city 2015-2018



- Why:
  - To support sociologists' and psychologists' work (e.g. education, social interaction..) to understand humans
  - To support computers' work in on-line or off-line knowledge extraction about humans for **a huge number of applications, services and systems** (surveillance, HCI, automotive, augmented experiences..)



- 1997-2000 MIT Alex Pentland: PFINDER projects and understanding interactions
- 2006- datasets for action analysis (Weizmann ICCV2005), action understanding now is popular\*
- 2010- 7 workshops on **HBU** (from 2010: IAPR, AMI, IROS, ACM MM, ECCV, ACMMM2016)
- **2011- Chalearn** workshops 2011- 2016; CVPR 2016 challenge “Looking at People”
- Many many datasets...

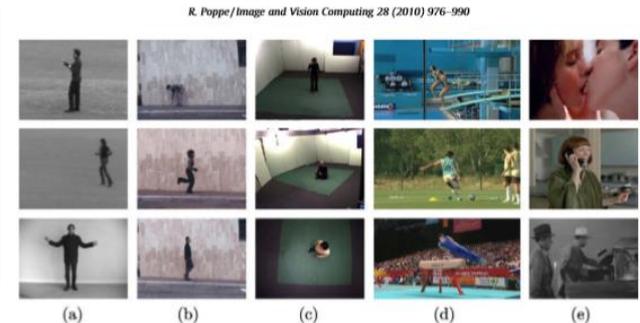
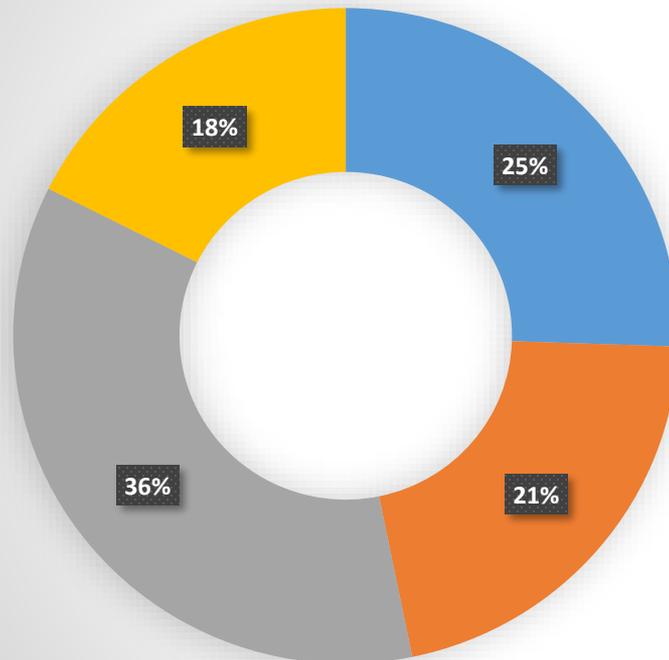


Fig. 1. Example frames of (a) KTH dataset, (b) Weizmann dataset, (c) Inria XMAS dataset, (d) UCF sports action dataset and (e) Hollywood human action dataset.

- One for all: \* R. Poppe “A survey on vision based action recognition” Image and vision computing 2010

- HBU by vision:
- More than 1000 papers from 2010 to 2015
- In 2016? Surely more

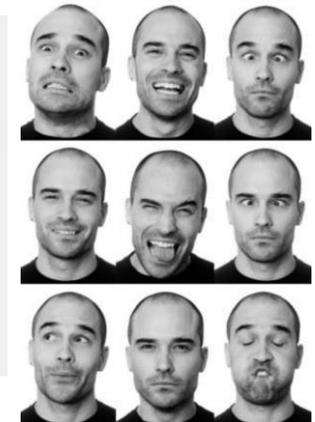
1000 scientific papers from google scholar 2015



- HBU & surveillance
- HBU & multimedia
- HBU & interaction
- HBU & health care

## Movements

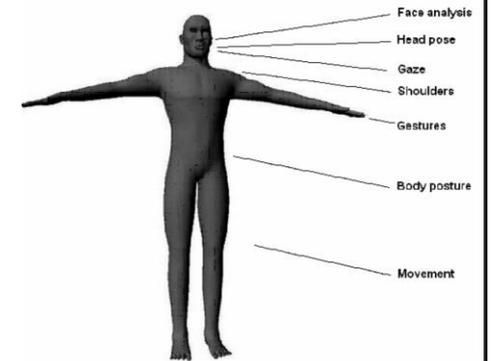
- Body movements
- Gestures
- Poses
- Gaze and Expression



## Actions

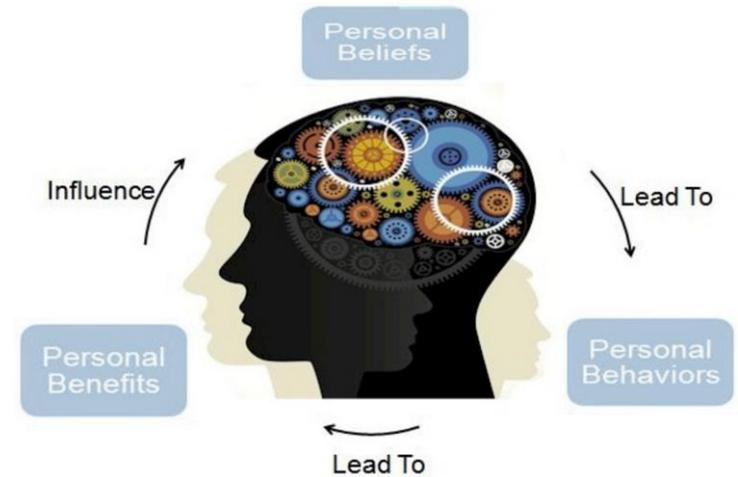


## Activities



## Behavior?

- Behavior
- {*movements, actions, activities*} +  
{*environment, objects, people*} +  
{*purposes, beliefs, habits*}

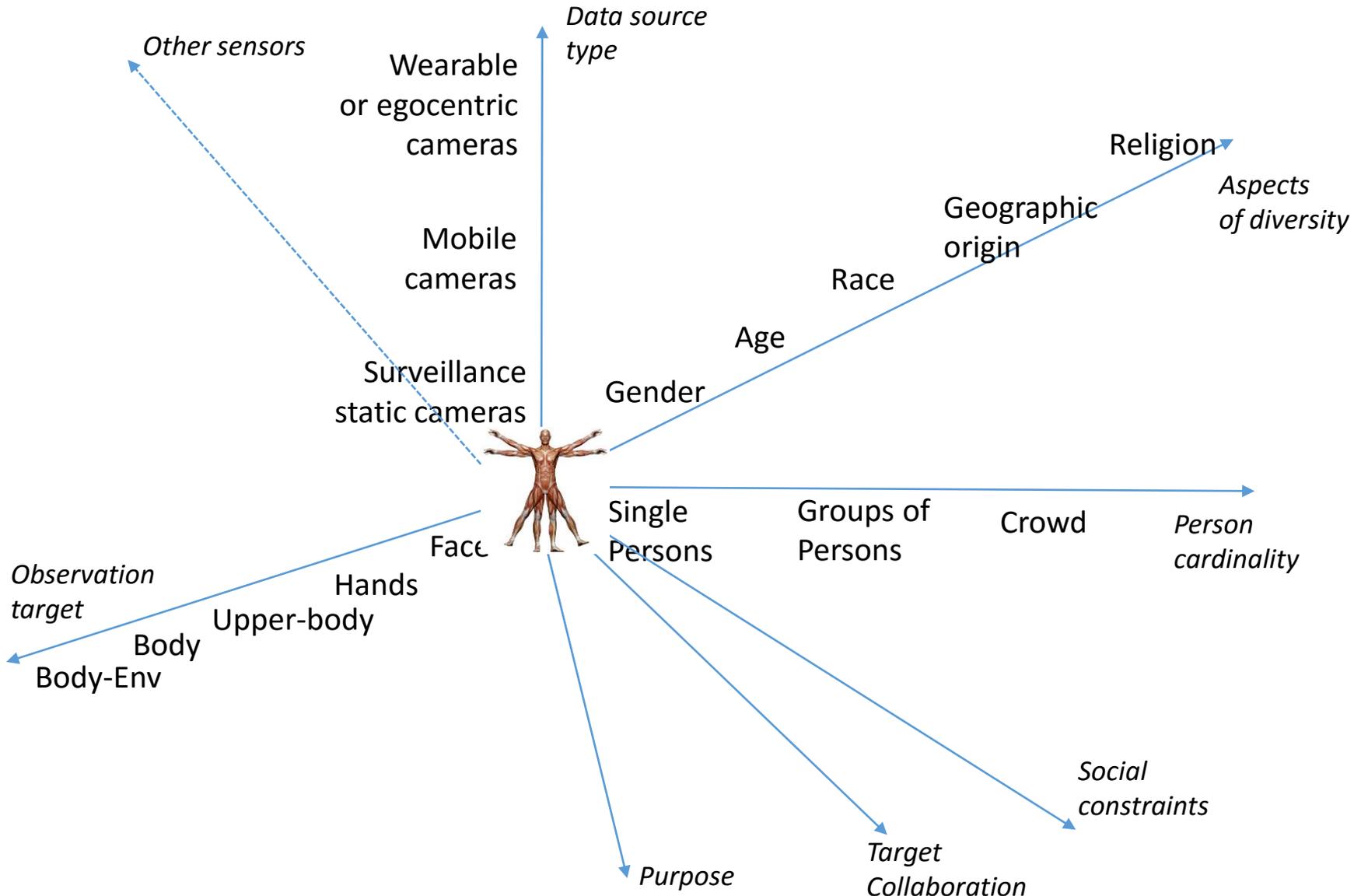


- Self-behavior, Personal behavior, Social behavior, ...





# HBU by Vision, a multidimensional space



- Levels of details



Human parts



Lev1 self Behavior  
expression,  
gesture, pose

Human Full body



Lev2 whole person  
Behavior  
Action and  
interaction

Human(s)  
In the environment



Lev3 Social  
Behavior

Humans  
In crowd



## Understanding humans for New natural Human Computer interaction systems \*



What are they doing?



## Deaf Sign language



Maja Pantic, Alex Pentland, Anton Nijholt and Thomas Huang *Human Computing and Machine Understanding of Human Behavior: A Survey* ICMI 2006

S. Rautaray, A Agrawal *Vision based hand gesture recognition for human computer interaction: a survey* [Artificial Intelligence Review](#) January 2015,

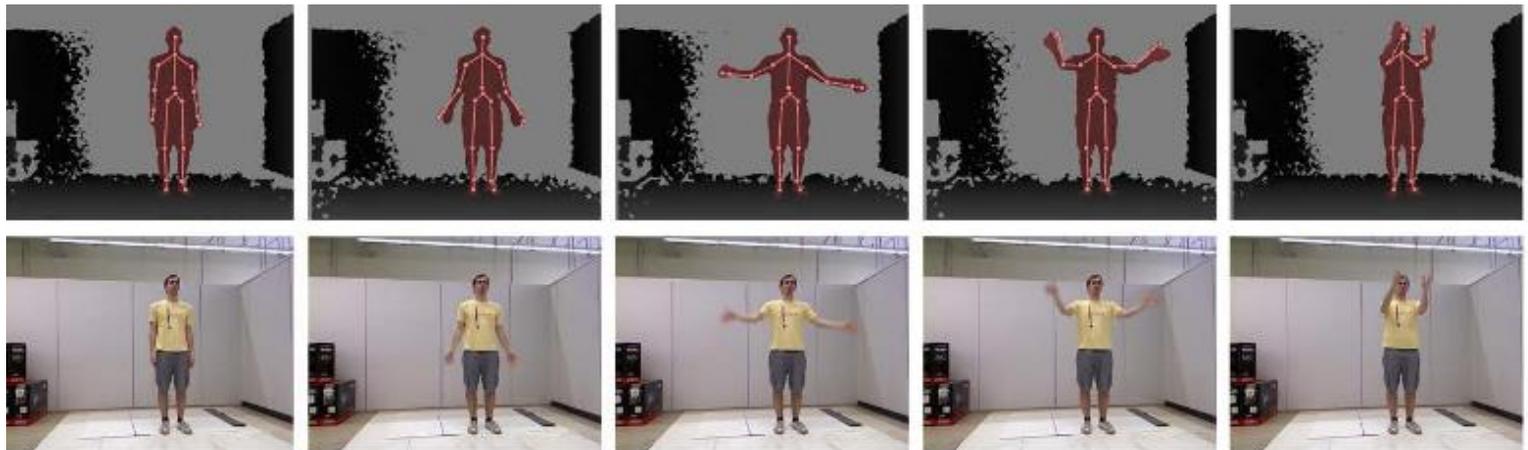
- Human attention in the car



## Lev2: Body actions



Fig. 7. Sample input frame of the Weizmann dataset



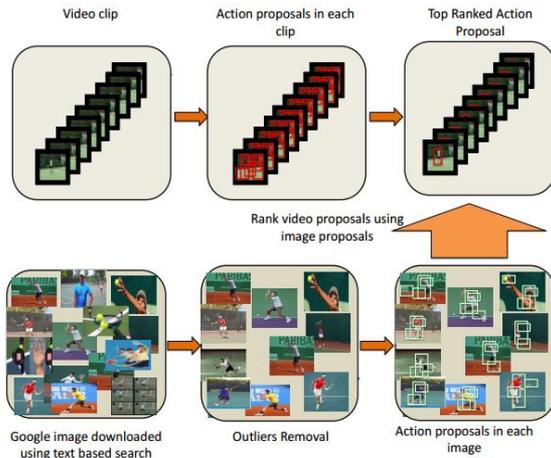
R.Vezzani, D.Baltieri, R.Cucchiara HMM Based Action Recognition with Projection Histogram Features ICPRW2010 supported by EU THIS Project

G.Borghini, R. Vezzani, R.Cucchiara; ["Fast gesture recognition with Multiple Stream Discrete HMMs on 3D Skeletons"](#) *Proceedings of the 23rd International Conference on Pattern Recognition*, Cancun, Dec 4-8, 2016, 2016

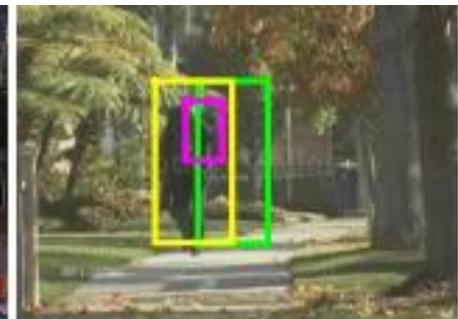
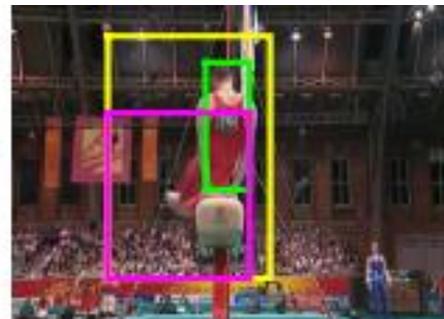
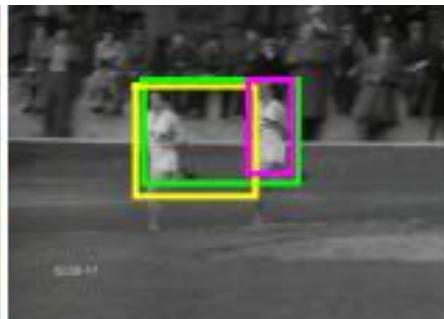
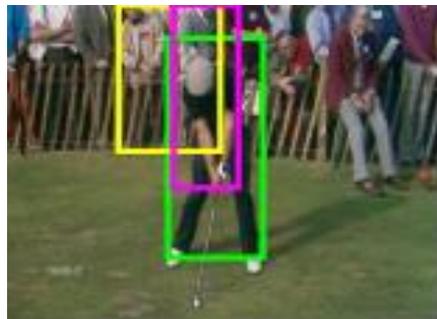
- CVPR2016

## What if we do not have multiple videos of the same action? — Video Action Localization Using Web Images

Waqas Sultani, Mubarak Shah



- Graph representation on DL features
- Graph based optimization
- Probabilistic Hough Matching for proposals
- Optimization in superpixel
- Image and video proposal



- CVPR2016

## VLAD<sup>3</sup>: Encoding Dynamics of Deep Features for Action Recognition

Yingwei Li<sup>†</sup>

Weixin Li<sup>†</sup>

Vijay Mahadevan<sup>§</sup>

Nuno Vasconcelos<sup>†</sup>

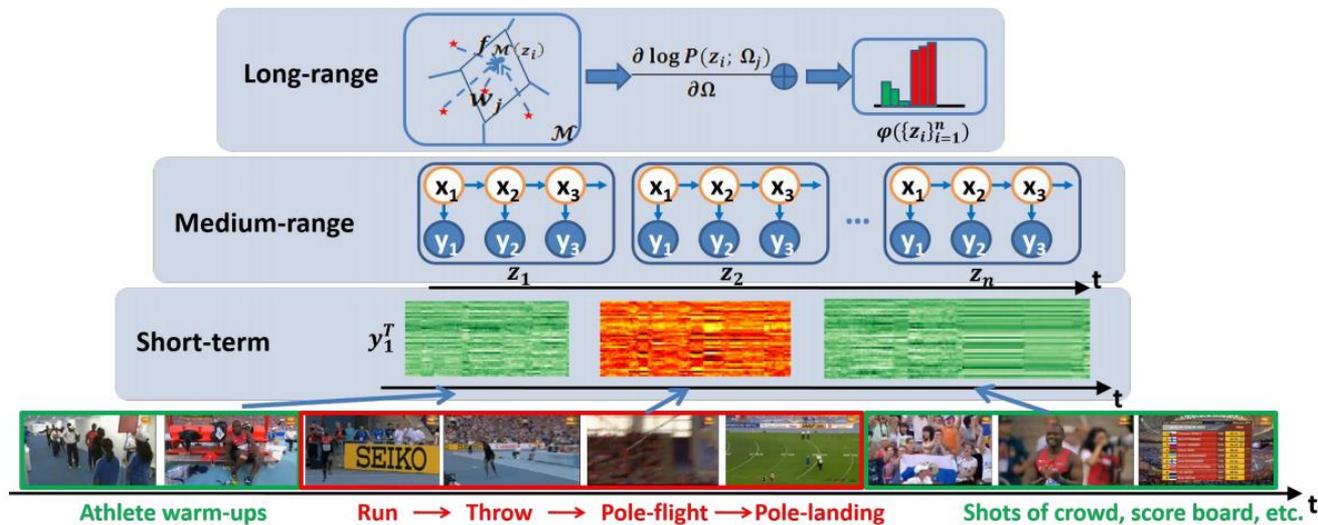
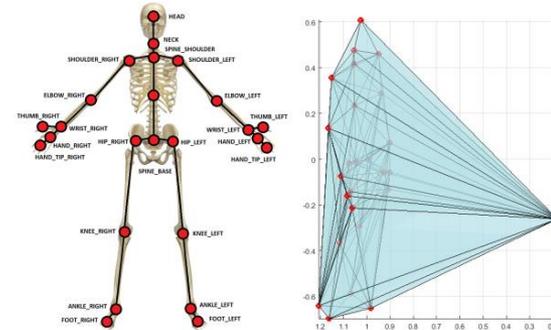
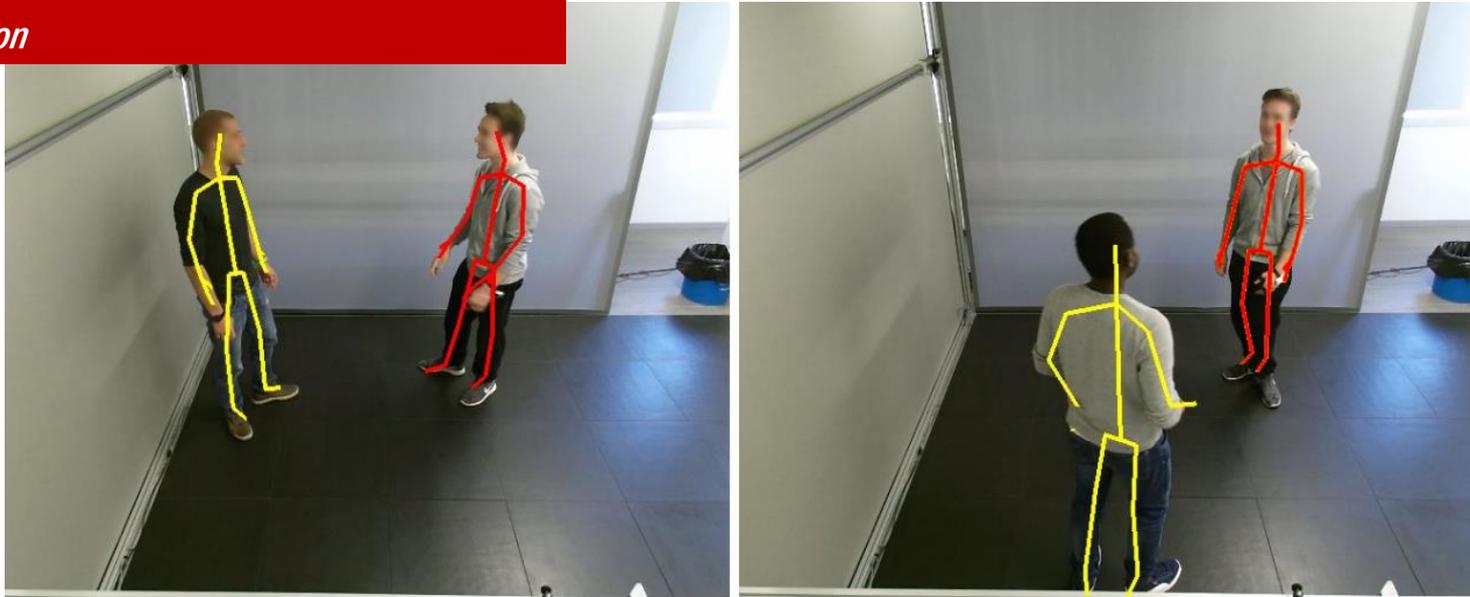


Figure 1: The VLAD<sup>3</sup> is inspired by the hierarchical structure of video dynamics. A **short-term** stage captures short-term appearance and motion patterns with deep features. A **medium-range** stage models the dynamics of segments of deep features, using an LDS. Finally, a **long-range** stage computes and pools a VLAD descriptor, derived from the LDS.



## Measuring human interaction



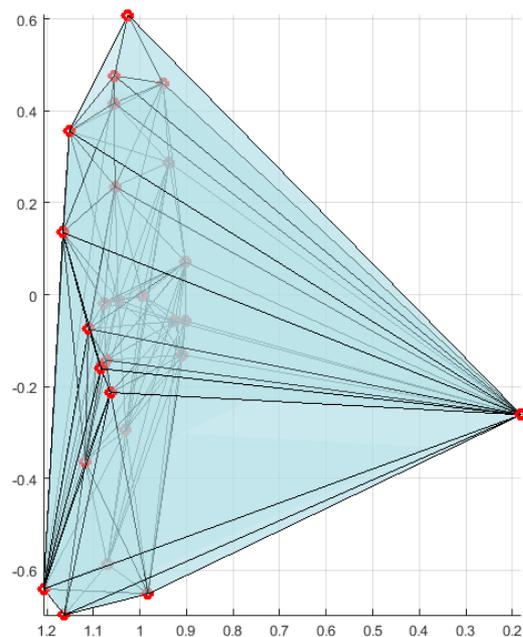
Palazzi,A.; Calderara,S.; Biccocchi,N; Vezzali,L.; di Bernardo,G.; Zambonelli,F; Cucchiara, R. "[Spotting prejudice with nonverbal behaviours](#)" ACM International Joint Conference on Pervasive and Ubiquitous Computing, Heidelberg, Settembre 2016,

- Simple but effective CV features

## Volume distance

$$F_{vol}(f) = Vol(DT(P(f)))$$

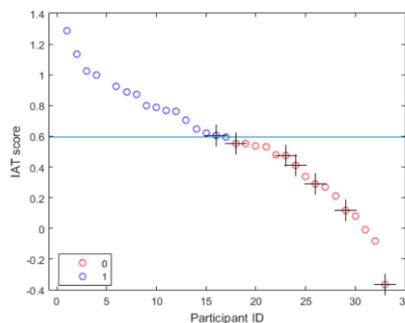
$DT(P(f))$  = Delaunay Triangulation



## Mutual distance

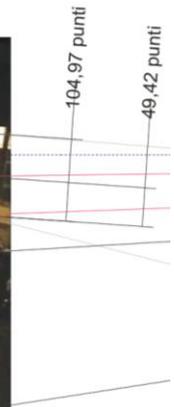
$$D_m(f) = dist(C^p, C^c)$$

$$C^p(f) = \frac{1}{m} \sum_{i=1}^m joint_i^p(f) \quad C^c(f) = \frac{1}{m} \sum_{i=1}^m joint_i^c(f)$$



- What are they doing?

*Real-time surveillance*



ENVI-VISION



EGO-VISION

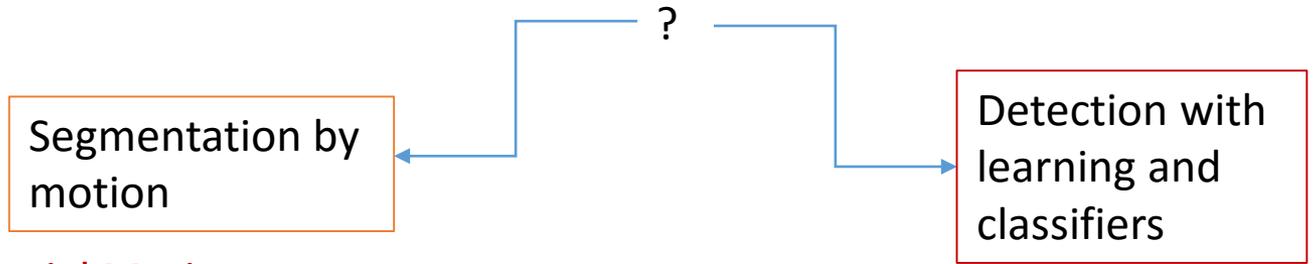


We need

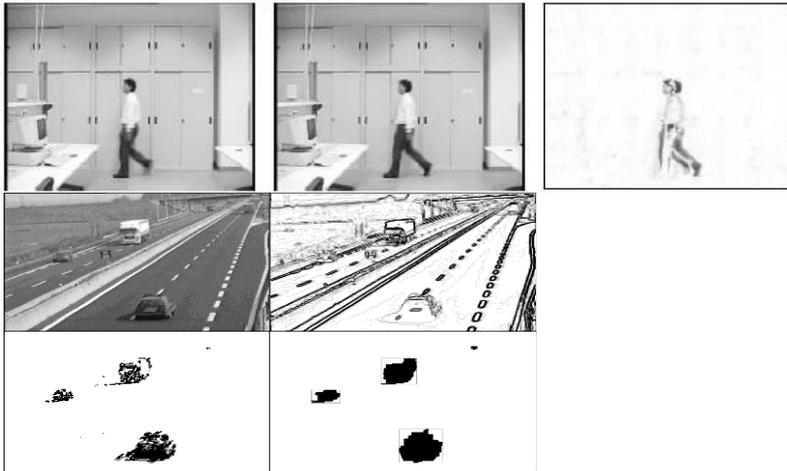
- Detection
- Tracking ( single , MTT, MTT in Multiple cameras..)
- Understanding movements, actions, activity, behaviors and social relation



# Segmentation/detection in surveillance



## Differential Motion



## HOG, part based, CNNs People detectors



## Background suppression



R. Cucchiara, C. Grana, M. Piccardi, A. Prati "Detecting Moving Objects, Ghosts and Shadows in Video Streams", *IEEE Trans on PAMI*, 2003

End-To-End People Detection in Crowded Scenes  
Russell Stewart, Mykhaylo Andriluka, Andrew Y. Ng *CVPR2016*

## Tracking (few) people

- Tracking few people in a constrained environment: «solved problem» 😊



Tracking by detection: using people detection for initialize ROI-based tracking (eg particle filter)

In semi-constrained world  
Tracking is possible



- **Is tracking a solved problem?**



- We tried to answer this questions in an “**experimental evaluation**”

- Even in case of single target tracking\*

- - a very large dataset

- of 14 categories of challenges

- - a large set of performance measures

- - a large experimentation

- (with code available over 3 clusters in 3 labs)

315 video ALOV++  
<http://www.alov300.org>  
<http://imagelab.ing.unimo.it/dsm>

MOTA; OTA; Deviaton....  
F-Measure  
SURVIVAL CURVES..

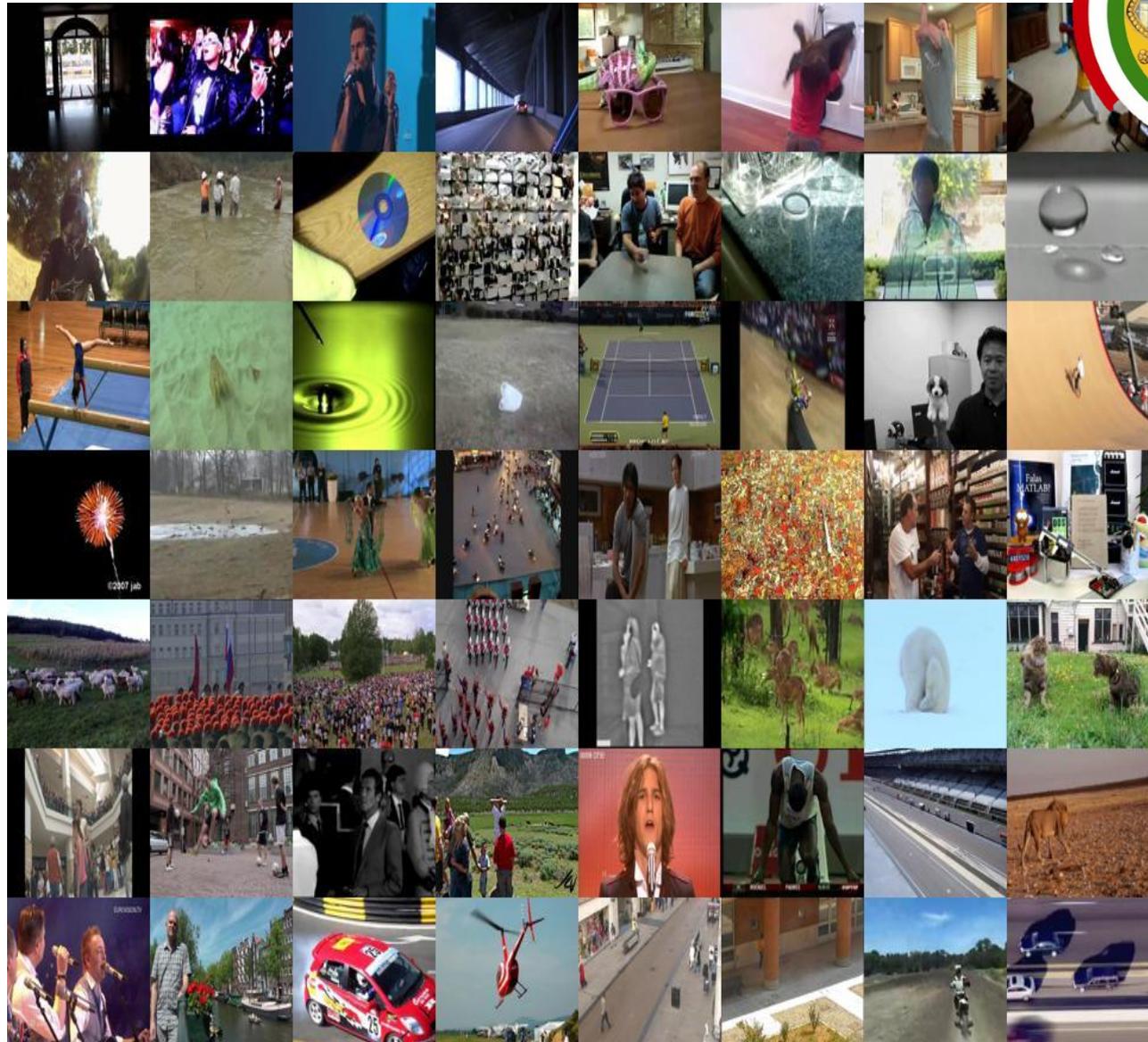
19 trackers  
BASELINES  
STATE OF THE ART

\* *D.Chu, A.Smeulders, S.Calderara, R.Cucchiara, A. Dehghan, M.Shah* **Visual Tracking: an Experimental Survey**  
*Transactions on PAMI 2013*

# 14 tracking challenges in 313 videos



- [01-LIGHT](#)
- [02-SURFACECOVER](#)
- [03-SPECULARITY](#)
- [04-TRANSPARENCY](#)
- [05-SHAPE](#)
- [06-MOTIONSMOOTHNESS](#)
- [07-MOTIONCOHERENCE](#)
- [08-CLUTTER](#)
- [09-CONFUSION](#)
- [10-LOWCONTRAST](#)
- [11-OCCLUSION](#)
- [12-MOVINGCAMERA](#)
- [13-ZOOMINGCAMERA](#)
- [14-LONGDURATION](#)

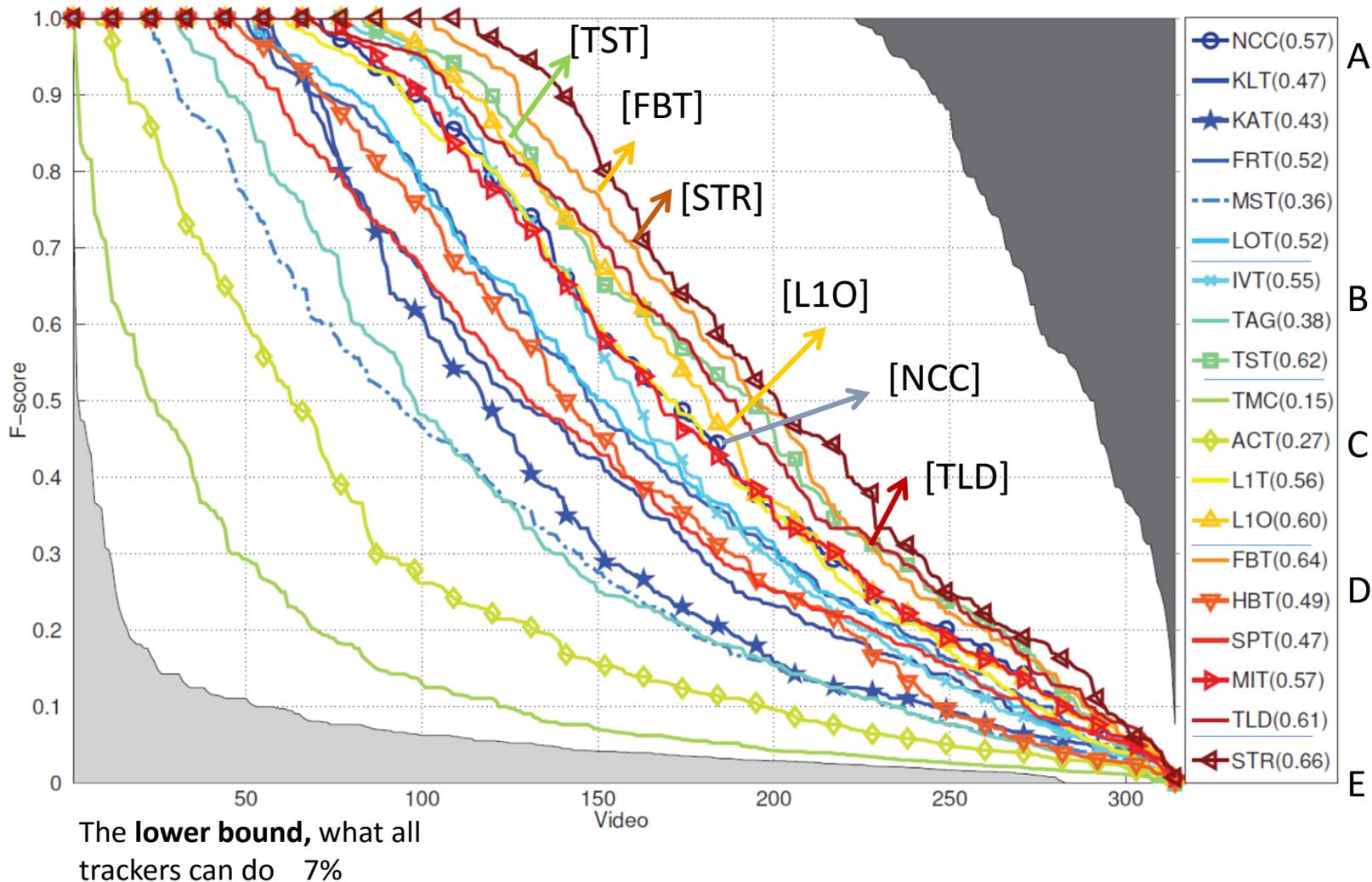




# A comprehensive view Survival curve

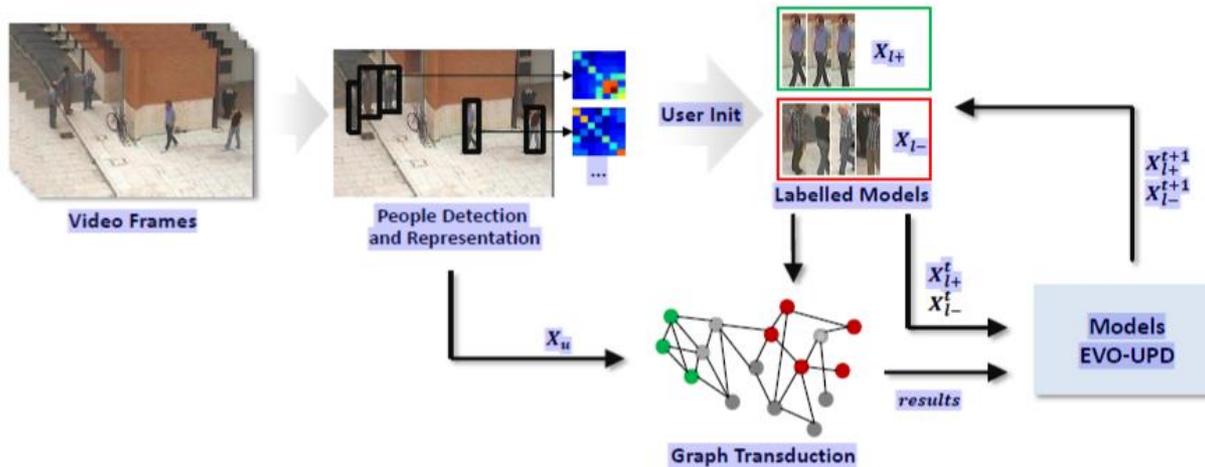
About the 30%, correctly tracked only

The **upper bound**, taking the best of all trackers at each frame 10%

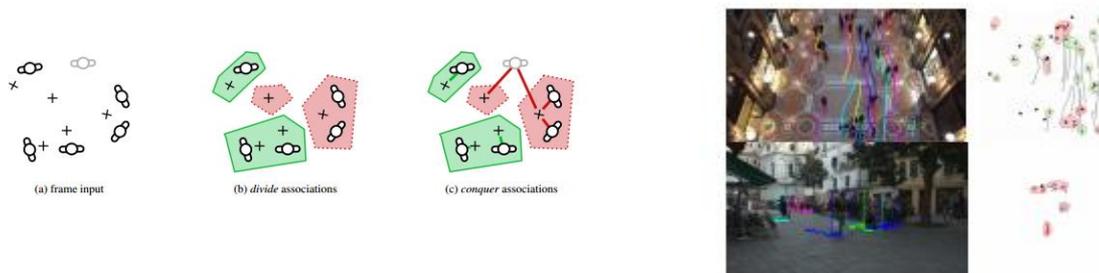


The **lower bound**, what all trackers can do 7%

## Tracking by detection with transductive learning \*

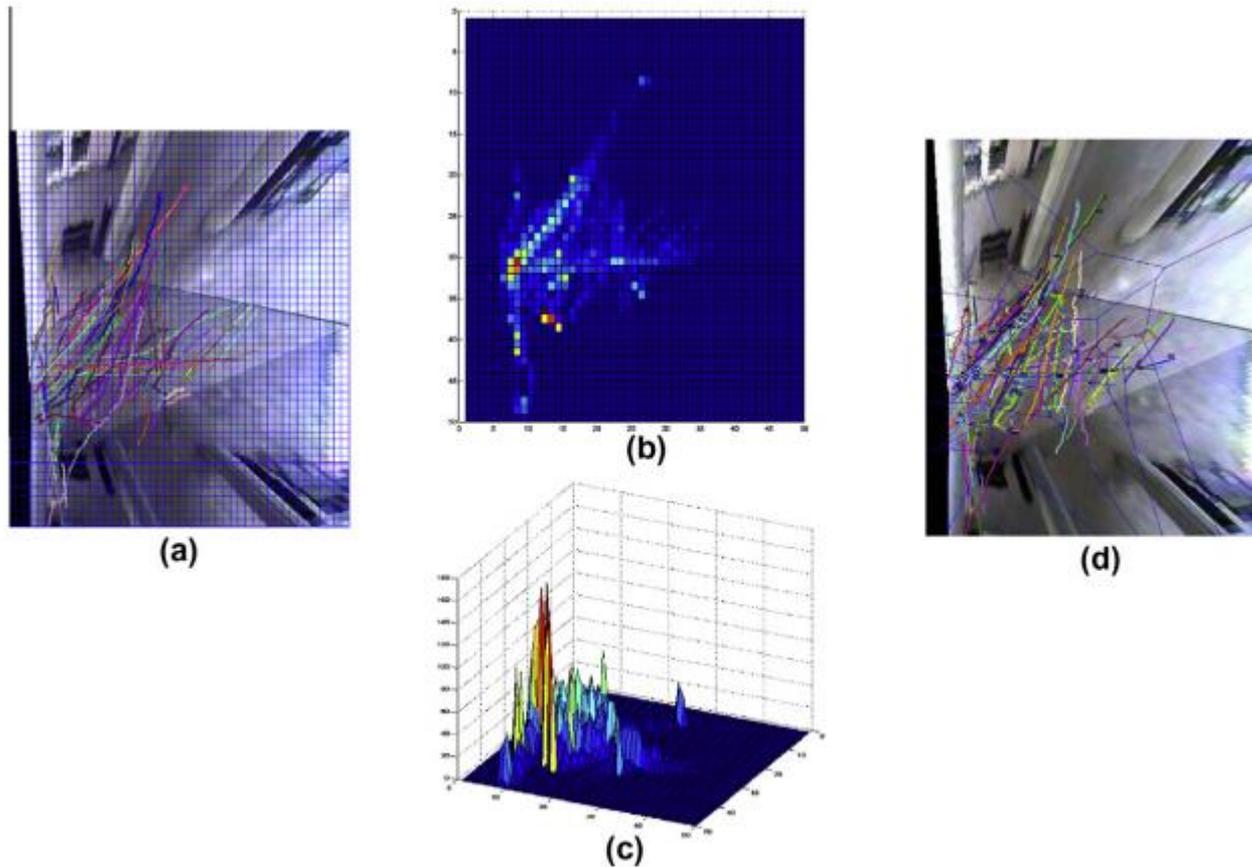


## Tracking by detection with structural SVM \*\*

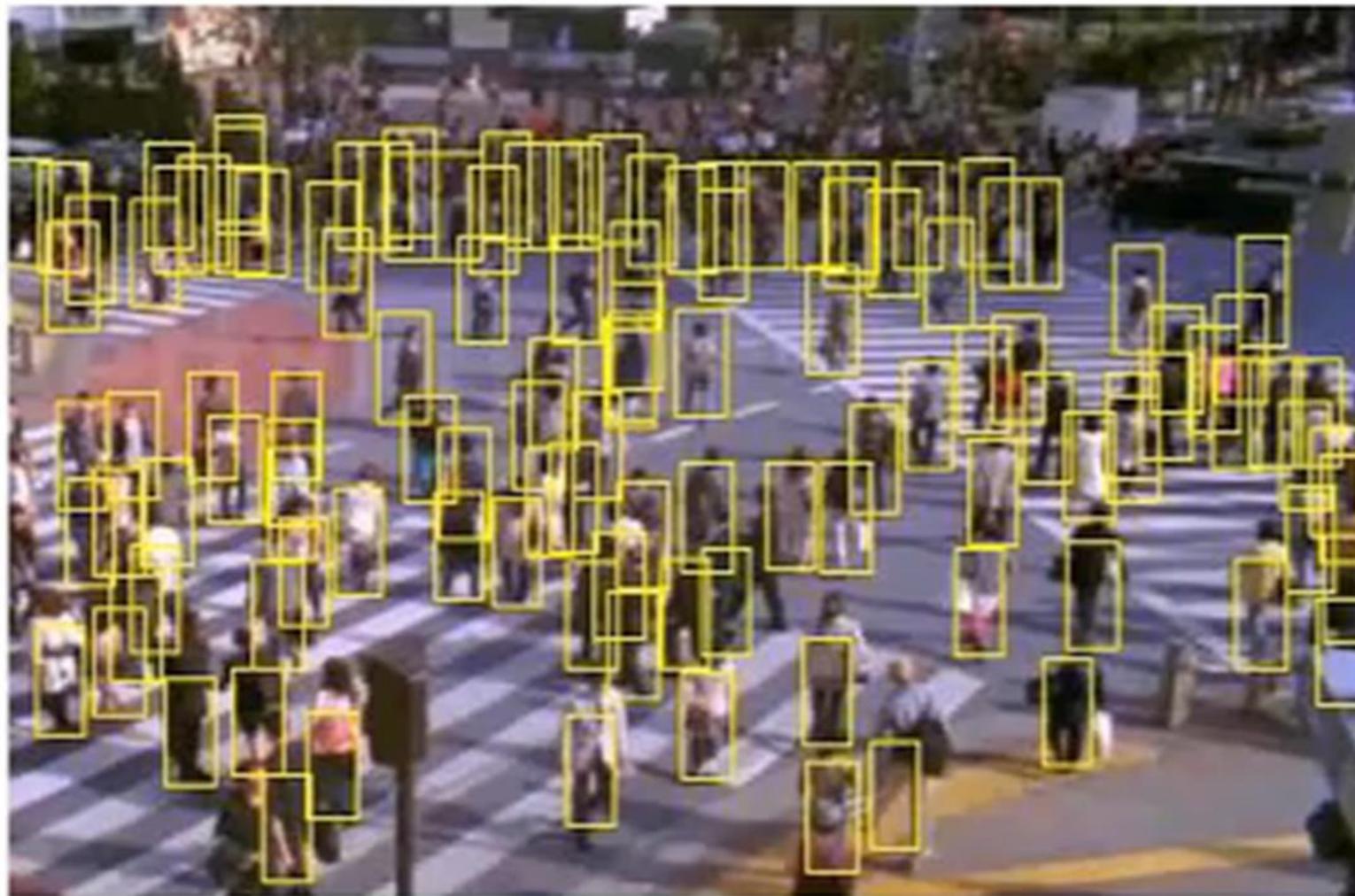


\*D. Coppi, S. Calderara, R. Cucchiara "Transductive People Tracking in Unconstrained Surveillance"  
Transactions on CSVT 2015

\*\* Francesco, Solera; Simone, Calderara; Rita, Cucchiara "[Learning to Divide and Conquer for Online Multi-Target Tracking](#)" Proceedings of ICCV 2015



**Fig. 2.** Irregular partitioning of the image area through Voronoi diagrams: (a) Reports the first regular division of the image ( $50 \times 50 = 2500$  cells in this example); (b) shows the top view of the 2D histogram, while (c) shows a side view; and (d) shows the resulting Voronoi diagram with 50 cells.





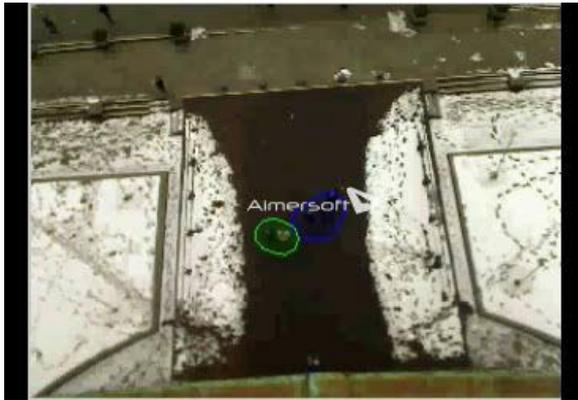
- If tracking were solved...

If the trajectories of every pedestrian in the scene (more or less) were available..  
**would we be able to discern the behaviour of groups?**

Features: Proxemics and Granger causality

Structure function: pair-wise correlation clustering

Group detection: Structured SVM [groups]



- UNIMORE and Duke University
- Duke dataset

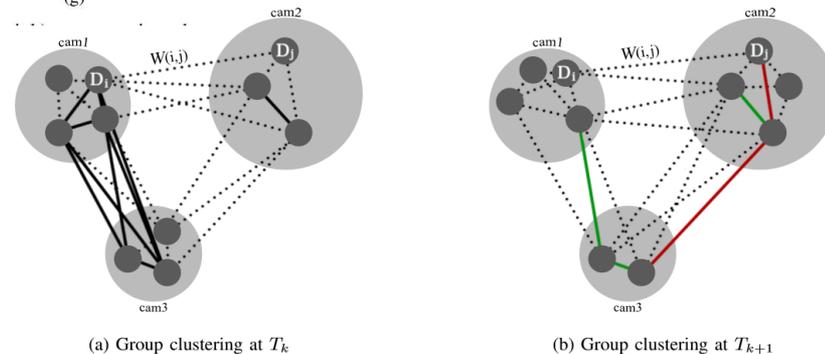
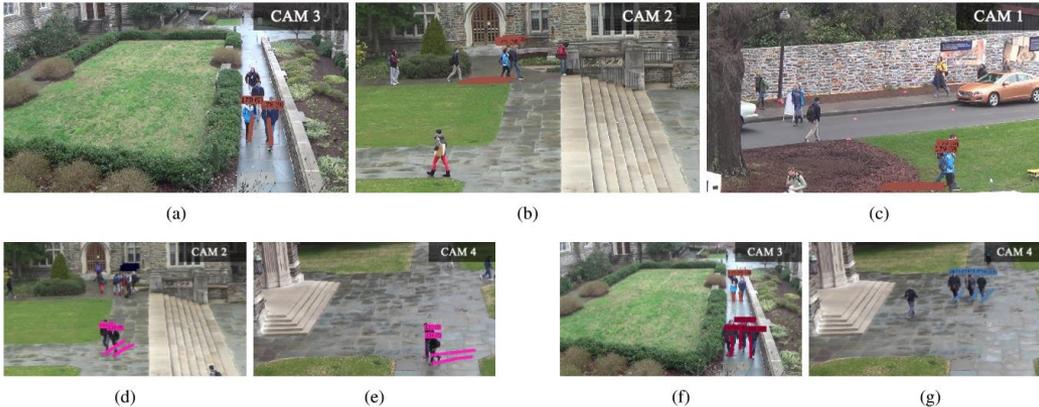


Fig. 2: The problem of tracking groups is cast as correlation clustering. In (a) all the detected groups  $D_i$  observed in the time window  $T_k$  are taken into account, all the pairwise correlation  $W(i, j)$  are computed (dashed lines) and a solution to tracking is found (solid lines). In (b), since time windows overlap in time,  $T_{k+1}$  will include group associations that were already solved in  $T_k$ . The new clustering is thus constrained by the previous solution forcing some observations to join (green lines) and others to remain separated (red lines), inducing consistent results across different time windows.

## Tracking Social Groups Within and Across Cameras

Francesco Solera, Simone Calderara, Ergys Ristani, Carlo Tomasi, Rita Cucchiara

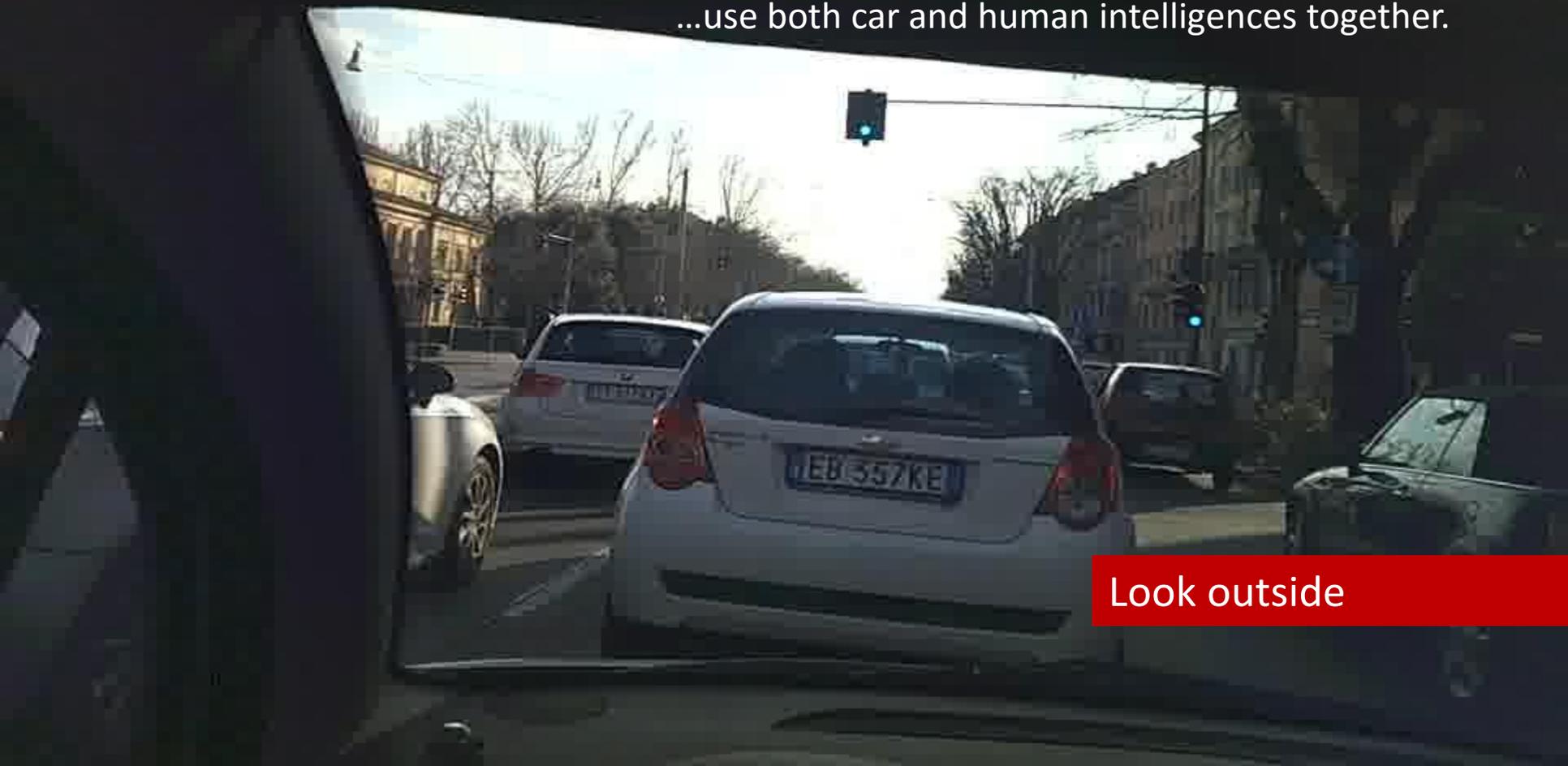
IEEE Trans. On Circuit and Systems for Video technology 2016

Let's come back to self-behavior understanding



- Understanding the human pose/gaze
- Understanding the human attention/distraction
- Learn the human driving behavior

...use both car and human intelligences together.



Look outside

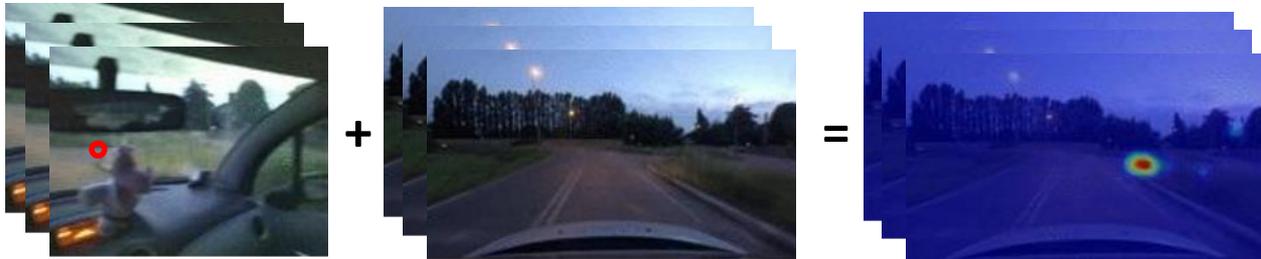
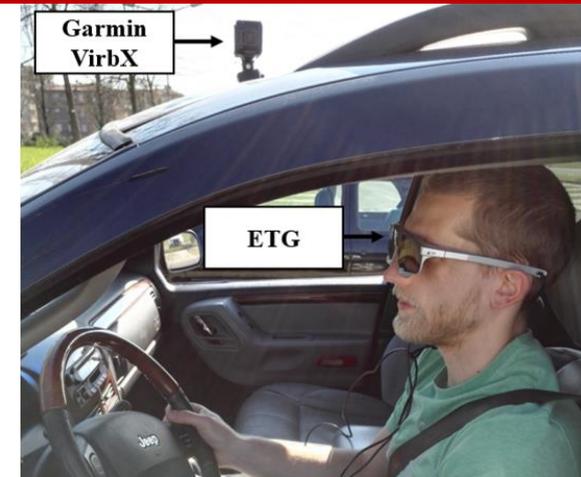


Look inside

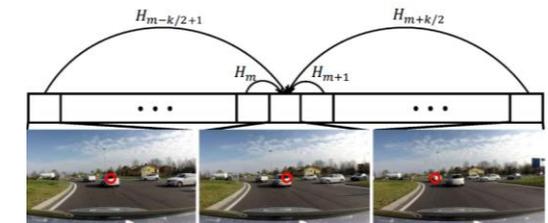


## Acquisition rig:

- Car mounted camera: Garmin VirbX 1080p/25fps, embedded GPS
- Eye tracker POV: SMI ETG HD camera 720p/30fps



*Gaze position projected on the video of the roof-mounted camera*



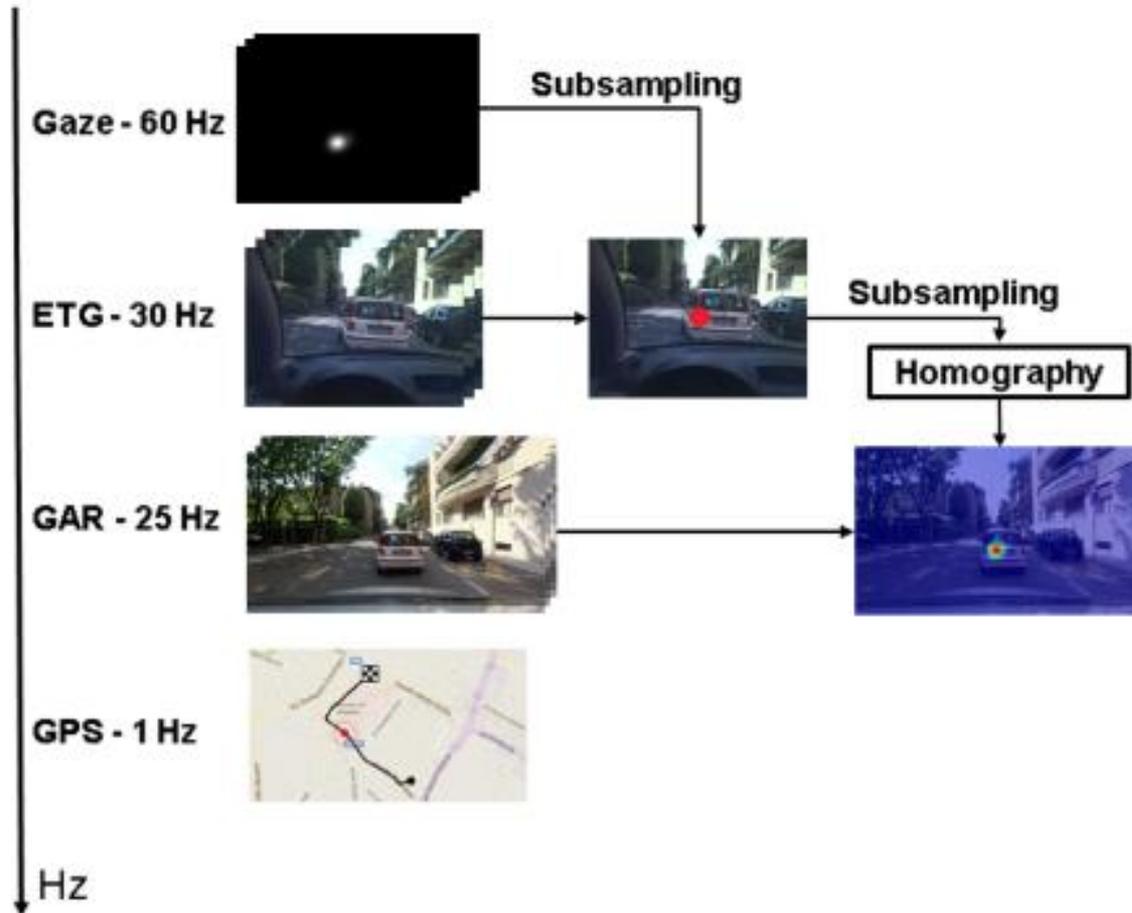
## Ground Truth:

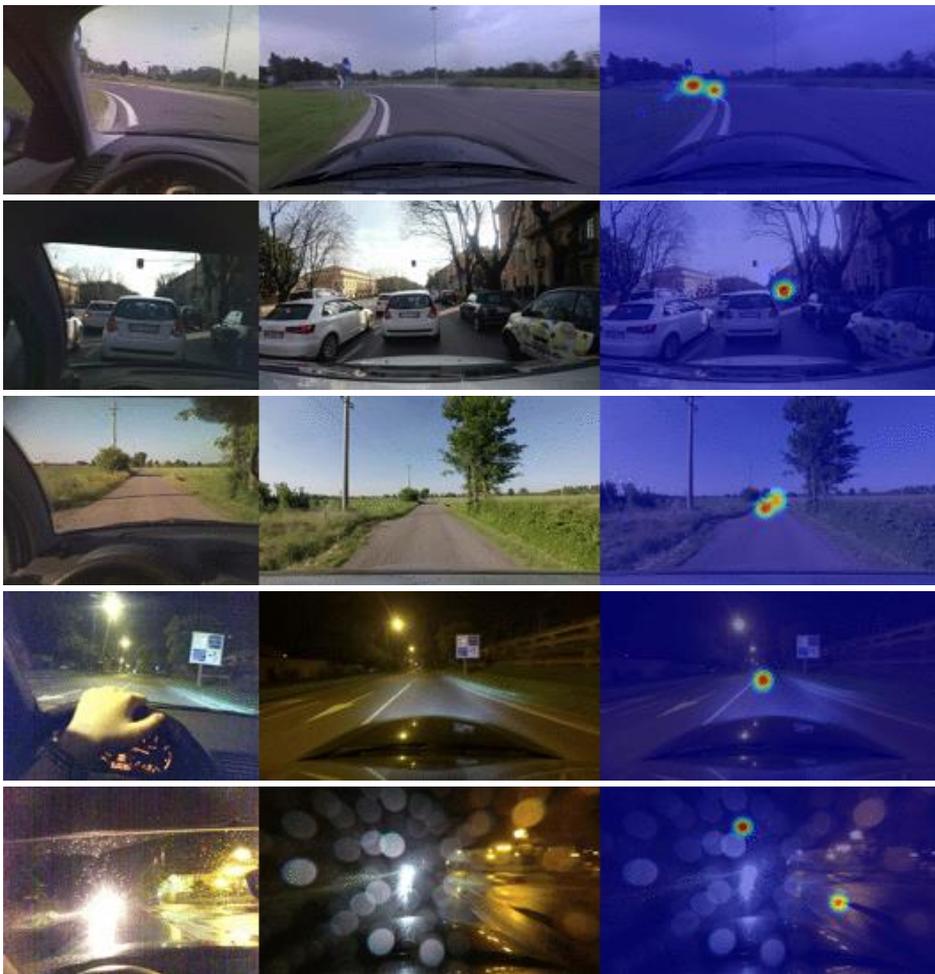
- Attentional map integrated over 25 frames (1 sec)
- Speed/GPS and driving course information

- Image registration and synchronization



- synchronization





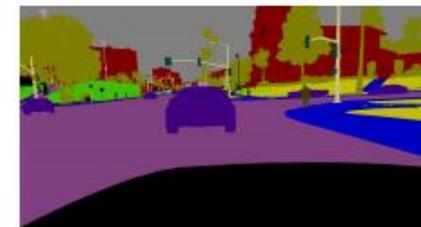
- 8 different drivers
- 3 different landscapes  
{Highway, Countryside, Downtown}
- 3 different weather's conditions:  
{Sunny, Cloudy, Rainy}
- 3 different light's conditions:  
{Morning, Evening, Night}

**74 videos of 5 minutes each!**

# Other datasets

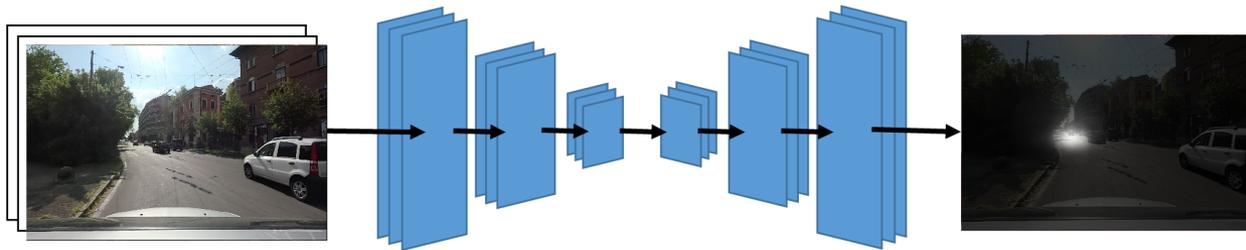
Dataset	Frames	Drivers	Scenarios	Annotations	Real-world	Public
Pugeault <i>et al.</i> [17]	158,668	n.d.	Countryside, Highway Downtown	9 classes in Environment Road, Junction, Attributes	Yes	No
Simon <i>et al.</i> [19]	40	30	Downtown	Gaze Maps	No	No
Underwood <i>et al.</i> [23]	120	77	Urban Motorway	n.d.	No	No
Fridman <i>et al.</i> [6]	1,860,761	50	Highway	6 Gaze Location Classes	Yes	No
Dr(eye)ve	555,000	8	Countryside, Highway Downtown	Gaze Maps	Yes	Yes

**Cityscapes** [2]: 25000 frames depicting street scenes. Each frame is annotated for both pixel-level and instance-level segmentation (20000 coarse, 5000 fine).

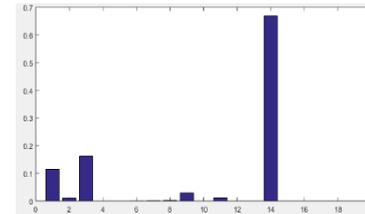
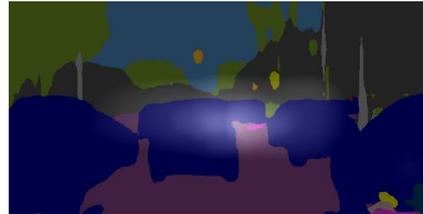


	frames	differentiation	real world	benchmark	addressed tasks
<b>Dr(eye)ve</b>	500000	landscapes, weather, daytime	yes	no	visual saliency
<b>Cityscapes</b>	25000	cities	yes	yes	semantic segmentation, instance level segmentation
<b>Kitti</b>	depends on the benchmark	landscapes	yes	yes	stereo, optical flow, visual odometry, object detection, tracking
<b>PfD</b>	25000	weather, daytime	no	yes	semantic segmentation





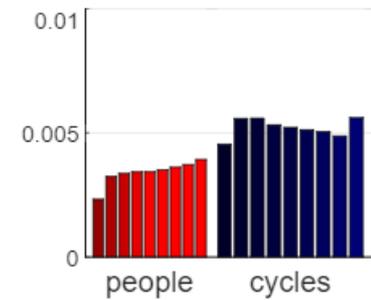
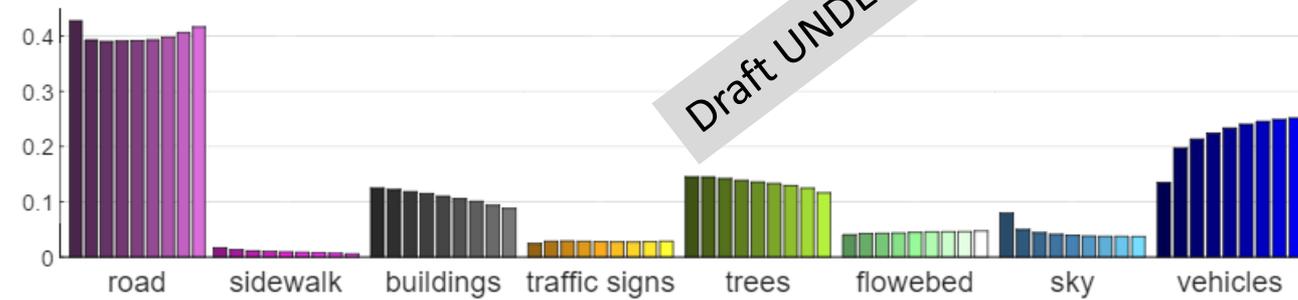
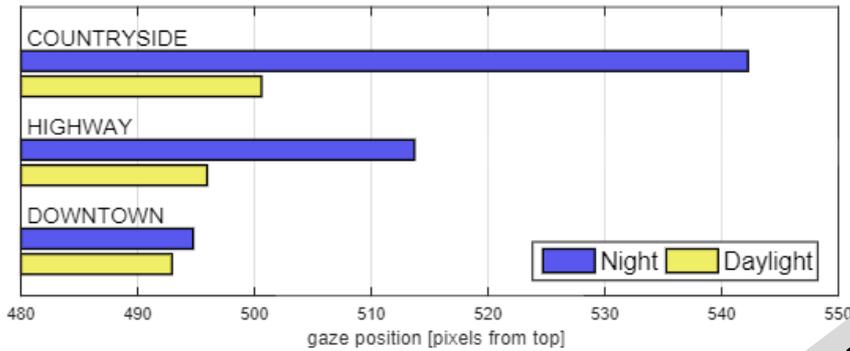
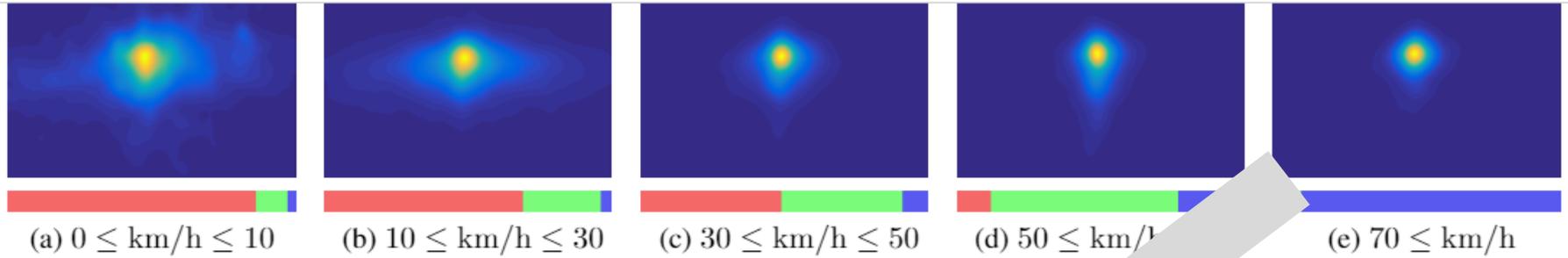
Good driving habits model: **where should we attend?**



Semantic segmentation: **what are we actually looking at?**

Look for us on <http://imagelab.ing.unimore.it/dreyeve>

# Some results: Attentive Behavior, measured



Draft UNDER SUBMISSION

\* Under submission

# Some results: Attentive Behavior, predicted

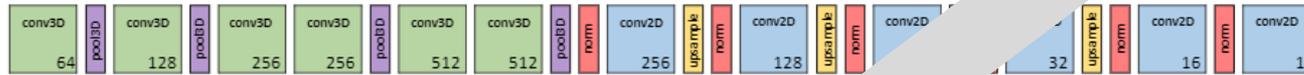
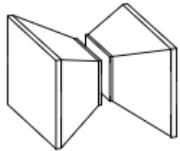
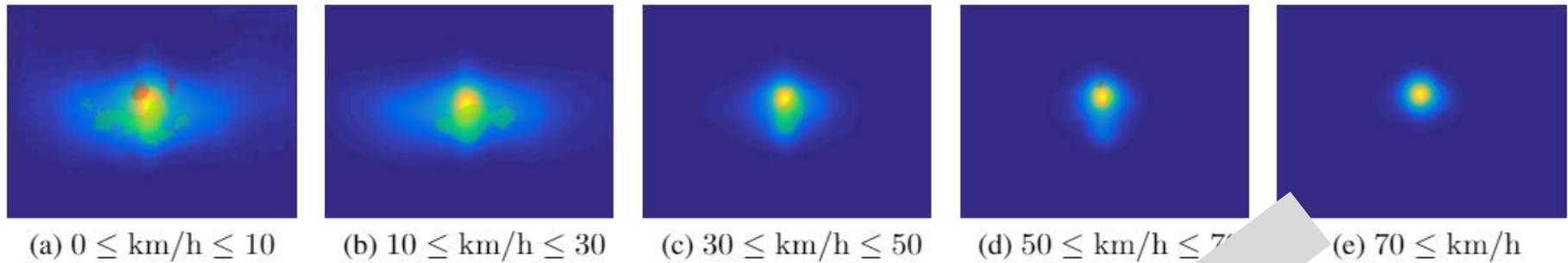
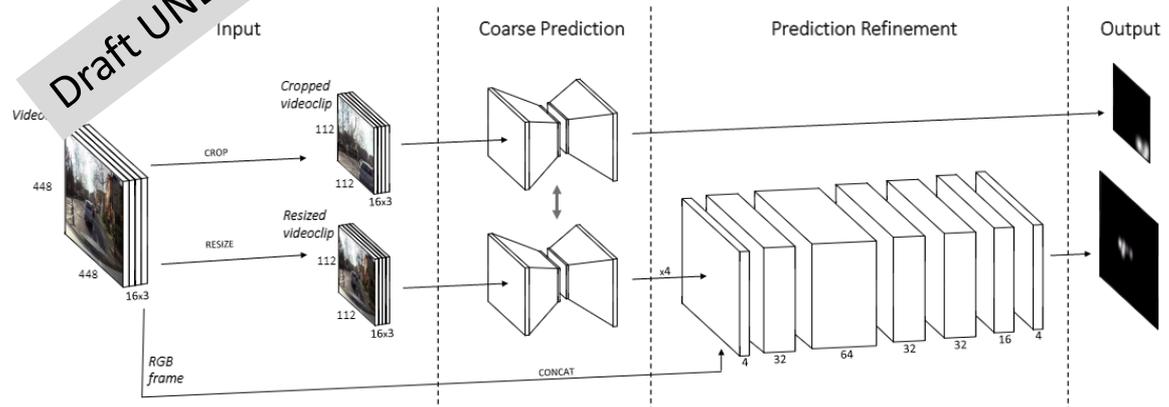


Fig. 9. Architecture of the coarse prediction module. The first part of the network performs feature encoding of the input videoclip. The input videoclip is a tensor of size  $3 \times 16 \times 112 \times 112$  that undergoes a sequence of conv3D and pool3D layers that gradually squeeze it to size  $512 \times 7 \times 7$ . All conv3D have kernel size (3,3,3) and ReLU activation units; all pool3D have pool size (2,2,2) except the first one that has pool size (1,2,2). In order to obtain a saliency map with the same spatial size of the input frame, the feature representation is decoded through a series of intertwined layers of batch normalization, conv2D and x2 upsampling on the spatial dimension. The conv2D have kernel size (3,3) and are followed by leaky ReLU activations with  $\alpha = .001$ . As a result, the output of the network is a tensor of size  $1 \times 16 \times 112 \times 112$ , i.e. the predicted grayscale saliency map.

Draft UNDER SUBMISSION



*Dr(eye)ve learned where the drivers see,  
and what the drivers pay attention on...*

*it is learning an intelligent visual behavior!*

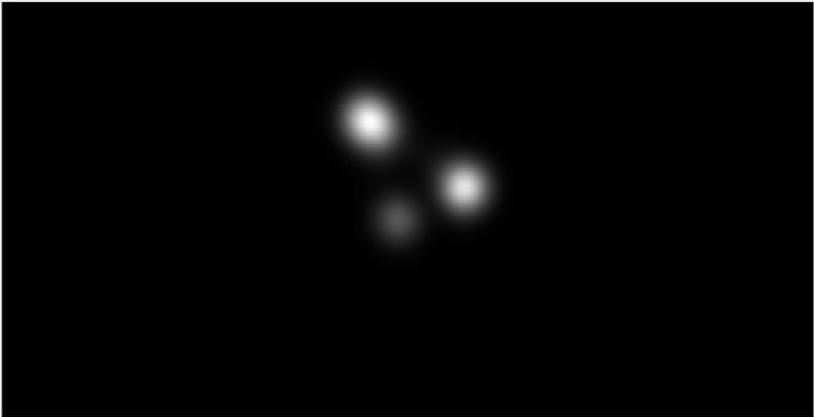
roof mounted camera



semantic segmentation



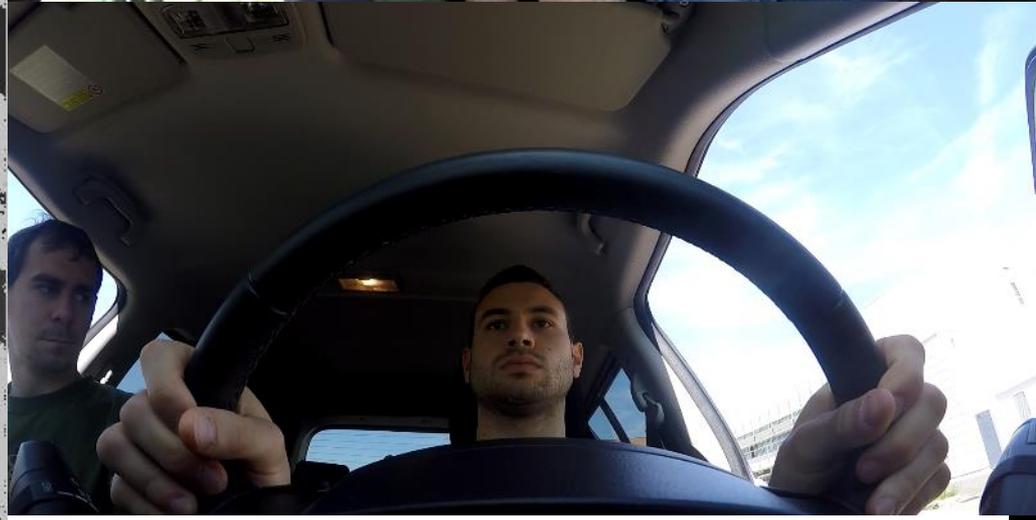
attention model prediction



overlay



# HBU IN THE CAR by depth images



Look outside

- Head Pose estimation
- Large literature on methods for Head Pose estimation\*\*
- Approaches:
  - Feature-based,
    - Nose, eyes
    - Landmark
  - Appearance-based
    - Pixels classifiers
    - CNNs\*
  - 3D model registration \*\*\*
  - Optimization based



\*\*E. Murphy-Chutorian and M. M. Trivedi. Head pose estimation in computer vision: A survey. IEEE Trans. Pattern Anal. Mach. Intell., 31(4):607–626, Apr. 2009.

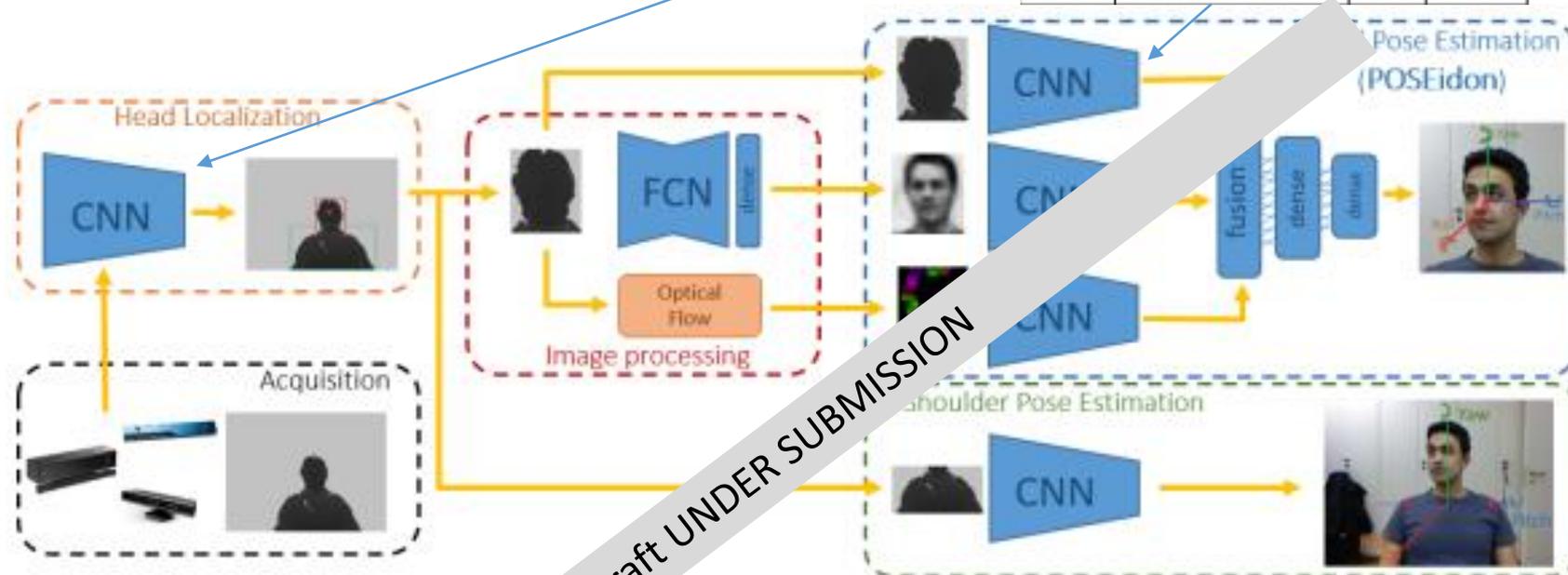
\*S. S. Mukherjee and N. M. Robertson. Deep head pose: Gaze-direction estimation in multimodal video. IEEE Transactions on Multimedia, 17(11):2094–2107, 2015.

\*\*\*C. Papazov, T. K. Marks, and M. Jones. Real-time 3d headpose and facial landmark estimation from depth images using triangular surface patch features. CVPR2015

- A new approach:\*
- only row depth images
- POSEidon a trident CNN

Order	Layer Type	# f/n	Kernel
1-2	2D Convolution (x2)	30	5x5
3	2D Convolution	30	4x4
4	2D Convolution	30	3x3
5	2D Convolution	120	3x3
6-7	2D Convolution (x2)	256	3x3
8-9	Dense (x2)	256	-
10	Dense	2	-

Order	Layer Type	# f/n	Kernel
1-2	2D Convolution (x2)	32	5x5
3	2D Convolution	32	4x4
4	2D Convolution	32	3x3
5	2D Convolution	128	3x3
6	Dense	128	-
7	Dense	84	-
8	Dense	3	-



Draft UNDER SUBMISSION

- Preliminary comparative results

Method	Year	Data	Location	Pitch	Roll	Yaw	Mean	Acc.
Fanelli [16]*	2011	Depth	14.0	$8.5 \pm 9.9$	$7.9 \pm 8.3$	$8.9 \pm 13.0$	$8.5 \pm 10.4$	0.790%
Yang [51]	2012	RGB + Depth	$3.97 \pm 2.18$	$9.1 \pm 7.4$	$7.4 \pm 4.9$	$8.9 \pm 11.0$	$8.5 \pm 6.9$	-
Padeleris [33]	2012	Depth	13.8	6.6	6.7	6.7	8.1	76.0
Rekik [38]	2013	RGB + Depth	5.1	4.3	5.2	5.1	4.9	-
Baltrusaitis [2]	2012	RGB + Depth	7.6	5.1	11.2	6.3	7.6	-
Ahn [1]*	2014	RGB	-	$3.4 \pm 2.9$	$2.6 \pm 2.1$	$2.8 \pm 2.4$	$2.9 \pm 2.6$	-
Martin [28]*	2014	Depth	5.8	2.5	2.7	3.6	2.9	-
Saeed [39]	2015	RGB + Depth	-	$5.0 \pm 5.0$	$3.5 \pm 4.6$	$3.9 \pm 4.2$	$4.4 \pm 4.9$	-
Papazov [35]	2015	Depth	8.4	$2.5 \pm 2.8$	$3.8 \pm 16.0$	$3.0 \pm 9.6$	$4.0 \pm 11.0$	-
Drouard [13]	2015	RGB	-	$2.8 \pm 2.8$	$4.7 \pm 4.6$	$4.9 \pm 4.1$	$5.2 \pm 4.5$	-
Meyer [30]	2015	RGB	-	2.4	2.1	2.1	2.2	0.946
Liu [25]	2016	RGB	-	$6.0 \pm 5.8$	$5.7 \pm 7.3$	$6.1 \pm 5.2$	$5.9 \pm 6.1$	-
<b>POSEidon</b>	2016	Depth	5.9	<b><math>1.6 \pm 1.7</math></b>	<b><math>1.8 \pm 1.8</math></b>	<b><math>1.7 \pm 1.5</math></b>	<b><math>1.7 \pm 1.7</math></b>	<b>0.950</b>

Table 4. Results on *Biwi Dataset*. In this case, no head localization is performed for head pose estimation task. Location is expressed in millimeters. The last column report the accuracy, established as the number of angle prediction below a certain threshold ( $10^\circ$ ).

# Understanding the human pose by depth only (with DL)

The main image displays a person's head and shoulder pose estimation results. The **HEAD POSE** section shows:

ROLL	: +01.64	█
PITCH	: +02.25	█
YAW	: +00.87	█

The **SHOULDER POSE** section shows:

ROLL	: +00.83	█
PITCH	: +03.66	█
YAW	: -29.34	█

On the right, there are three face-related visualizations:

- FACE DEPTH**: A grayscale depth map of the face.
- FACE GRAY**: A grayscale image of the face.
- FACE OF**: A color image of the face.

The **Head Pose Estimation** pipeline diagram is as follows:

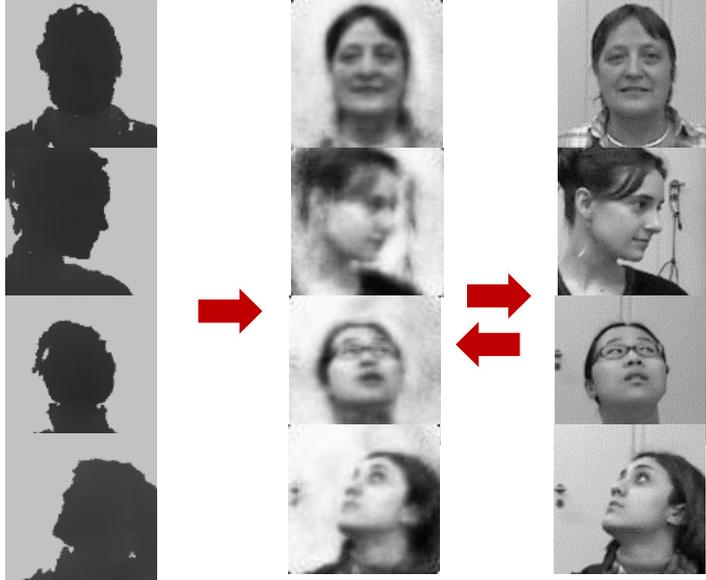
```
graph LR
    Input[Input] --> HL[Head Localization]
    subgraph HL_Box [Head Localization]
        HL_CNN[CNN]
    end
    HL --> HL_Box
    HL_Box --> HL_CNN
    HL_Box --> HPE[Head Pose Estimation]
    subgraph HPE_Box [Head Pose Estimation]
        HPE_CNN1[CNN]
        HPE_FC[FCN]
        HPE_CNN2[CNN]
        HPE_CNN3[CNN]
        HPE_Concat[concat]
        HPE_Dense1[Dense]
        HPE_Dense2[Dense]
    end
    HPE_CNN1 --> HPE_Concat
    HPE_FC --> HPE_Concat
    HPE_CNN2 --> HPE_Concat
    HPE_CNN3 --> HPE_Concat
    HPE_Concat --> HPE_Dense1
    HPE_Dense1 --> HPE_Dense2
    HPE_Dense2 --> Output[Output]
```

# Depth-to-face an impressive side effect

DEPTH

RGB from DEPTH

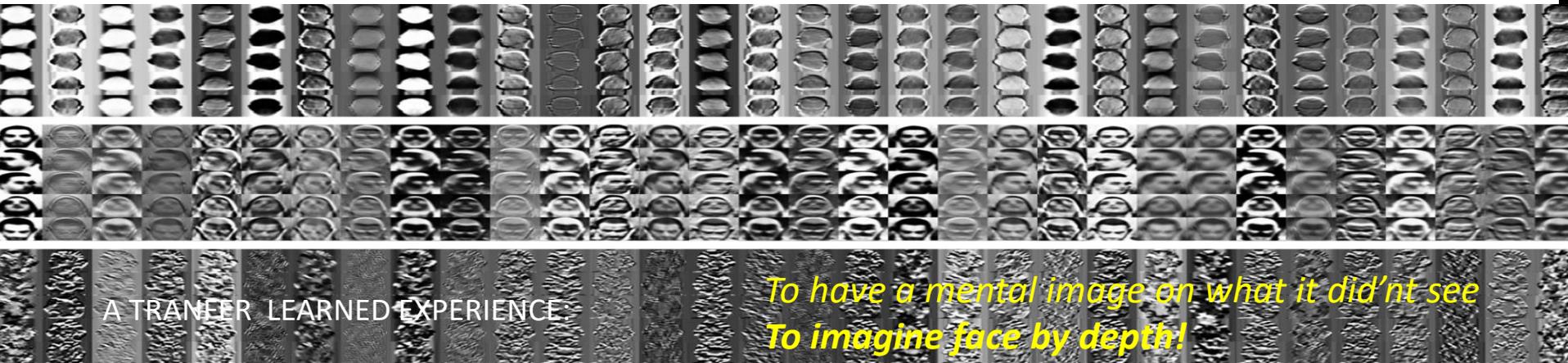
Original unseenRGB



Learned by the POSEidon Net @Imagelab



*POSEidon learned something more...*



A TRANSFER LEARNED EXPERIENCE

*To have a mental image on what it didn't see  
To imagine face by depth!*

- A new CNN architecture ( thanks to Guido Borghi and Marco Venturelli, Rob Vezzani)
- fuses the key aspects of autoencoder and fully connected deep networks
- The loss function works on centred images with a multivariate Gaussian on prior mask (parameters 3.5, 2,5)
- And Adadelta optimizer

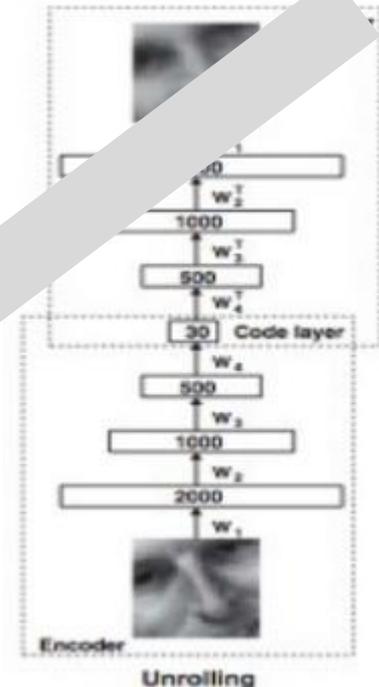
$$L = \frac{1}{N} \sum \sum \left[ \frac{1}{ch} \sum (y_{ik} - \bar{y}_{ik}) \right] w_{ij}$$

$$w = (2\pi)^{-\frac{k}{2}} |\Sigma|^{-\frac{1}{2}} e^{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)}$$

$$\mu, \Sigma = \left[ \begin{pmatrix} \frac{R}{2} \\ \frac{C}{2} \end{pmatrix}, \begin{pmatrix} \left(\frac{R}{\alpha}\right)^2 & \frac{R-C}{\alpha\beta} \\ \frac{R-C}{\alpha\beta} & \left(\frac{C}{\beta}\right)^2 \end{pmatrix} \right]$$

Order	Layer Type	#f/n	Kernel
1-2	2D Convolution (x2)	30	5x5
3	2D Convolution	60	4x4
4	2D Convolution	60	3x3
5-6	2D Convolution (x2)	120	3x3
7	2D Convolution	256	1x1
8	2D Convolution	256	1x1
9-10	2D Convolution (x2)	120	3x3
11	2D Convolution	60	3x3
12	2D Convolution	60	4x4
13-14	2D Convolution (x2)	30	5x5
15	Dense	$r \times c \times ch$	-

Table 2. Architecture for depth-to-gray reconstruction



## 3D Pandora dataset @Imagelab



# 2D Pandora dataset @Imagelab



# Learned by the Poseydon Net @Imagelab



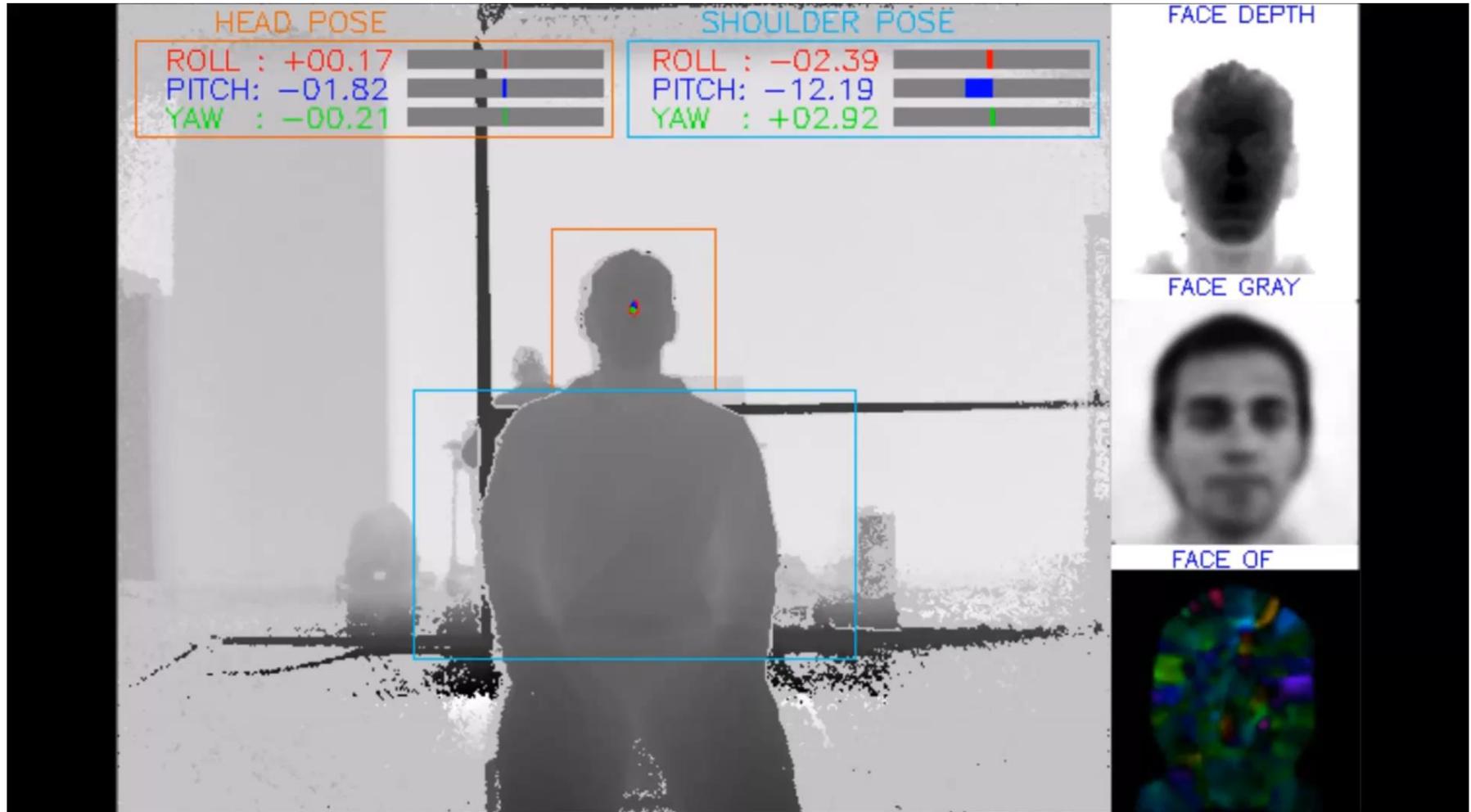


- Results on PANDORA Dataset

Method	Fusion	Head			Accuracy
		Pitch	Roll	Yaw	
depth (no crop)	-	$8.0 \pm 7.0$	$6.2 \pm 5.3$	$11.7 \pm 11.2$	0.553
depth	-	$6.5 \pm 6.6$	$5.4 \pm 5.1$	$10.5 \pm 11.8$	0.646
gray	-	$6.8 \pm 7.0$	$5.7 \pm 5.7$	$10.5 \pm 14.6$	0.647
rgb	-	$7.1 \pm 6.6$	$5.6 \pm 5.7$	$9.0 \pm 10.9$	0.639
optical flow	-	$7.7 \pm 7.5$	$5.7 \pm 5.7$	$10.0 \pm 12.5$	0.609
depth + gray	concat	$5.6 \pm 5.0$	$4.7 \pm 5.0$	$9.8 \pm 13.4$	0.698
depth + OF	concat	$6.0 \pm 6.1$	$4.5 \pm 4.8$	$9.2 \pm 11.5$	0.690
POSEidon	conv	$5.7 \pm 5.6$	$4.9 \pm 5.1$	$9.0 \pm 11.9$	0.715
POSEidon	concat	$6.3 \pm 6.1$	$5.0 \pm 5.0$	$10.6 \pm 14.2$	0.657
POSEidon	mul+concat	<b><math>5.6 \pm 5.6</math></b>	<b><math>4.9 \pm 5.2</math></b>	<b><math>9.1 \pm 11.9</math></b>	0.712

Draft UNDER SUBMISSION

# A demo



- Human behavior understanding (by vision)
  - A lot of stuff in computer vision and pattern recognition
  - **Geometry, graphs and statistical data analysis is unavoidable**
  - Features are mostly CNN-based
  - **Detection is not enough**
  - Spatio temporal coherence is needed (often tracking)
  - **New forms of input data are useful (sensors 3D...)**
  - A lot of learning ... the importance of datasets
  - **A strong knowledge of the context with experts (drivers, automotive industries, security persons, psychologists)**
  - **We are becoming truly multidisciplinary and our systems truly intelligent.**





<http://imagelab.ing.unimo.it>



Inter-dipartimental Research Center in ICT  
Tecnopolo di Modena  
Emilia Romagna High Technology Network

**Thanks to:**

**Rita Cucchiara**, Costantino Grana, **Roberto Vezzani**, **Simone Calderara**, [Giuseppe Serra],  
**Stefano Aletto**, Fabrizio Balducci, **Guido Borghi**, **Andrea Palazzi**,  
Federico Bolelli,  
[Marco Manfredi], Francesco Paci, **Francesco Solera**, Patrizia Varini,  
Lorenzo Baraldi,  
Andrea Corbelli, Marcella Cornia, Augusto Pieracci, Paolo Santinelli,  
Silvia Calio and **Marco Venturelli**

