

# ICCE2015, 15-18 March 2015 Kuwait City

*Invited speech @International Conference on Computer Science and Engineering:  
Big data science for the social goods*

## *Human Behavior understanding in video*

Prof. Ing. Rita Cucchiara

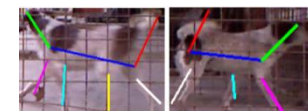
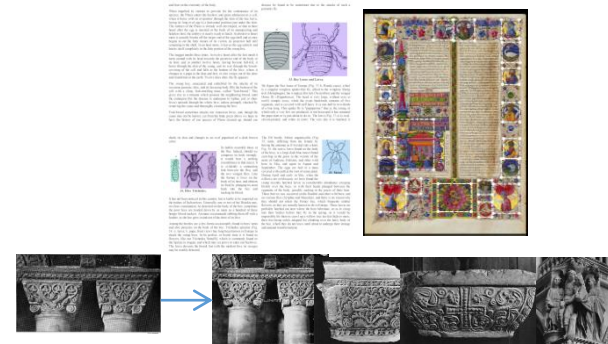
Imagelab, Dipartimento di Ingegneria «Enzo Ferrari», Modena, Italy

Director of the Research Center in ICT Softech-ICT



- Pattern recognition and Image processing
- Medical Imaging ( dermatology eu projects)
- Digitalized Document analysis ( Encycl.Treccani)
- Multimedia
- Multimedia big data annotation (RAI)
- 2D, 3D, wearable Computer vision
- Augmented experiences in culture and museums
- Experience with Wearable devices and IoT
- 2D and 3D augmented visits
- Computer vision for Behaviour analysis
- Children behaviour analysis
- Surveillance in crowd
- Animal behaviour

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



# Computer Vision and Human behaviour understanding

- **Computer Vision** is the scientific discipline studying how to perceive and understand the world through visual data by computers.

*Can computer vision provide effectively  
Human behaviour understanding?*

*Can computer vision is useful Big data?*





# The visual data cycle

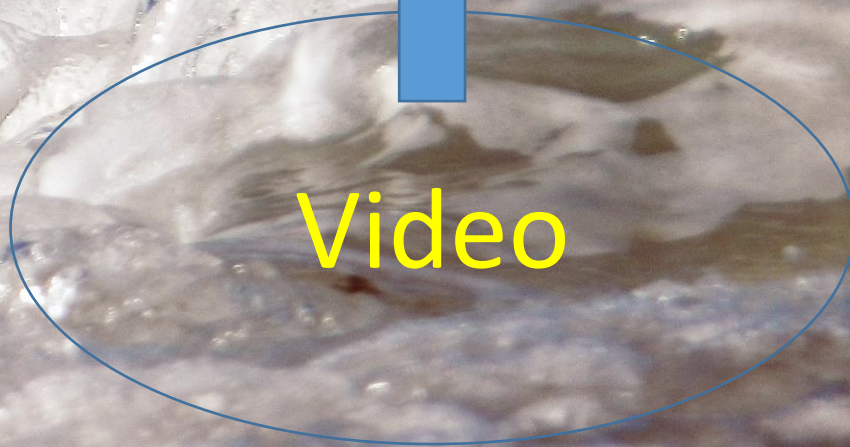
human behaviour  
data

Visual knowledge

Computer vision  
and pattern recognition

Video

Big data



- T. Huang, "Surveillance Video: The Biggest Big Data," *Computing Now*, vol. 7, no. 2, Feb. 2014 - See more at: <http://www.computer.org/web/computingnow/archive/february2014#sthash.b4UxnARn.dpuf>

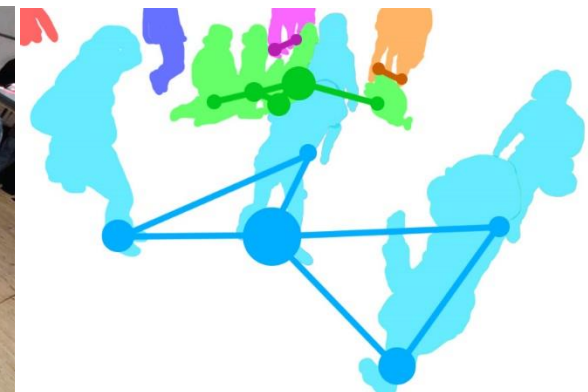


# Human behaviours

- What he is doing?
- What are they doing?
- Single and collective behaviours,
- Collaborative or not collaborative behaviours



Italian project  
Surveillance in  
cultural cities



Italian project “the educating city” cluster smart city 2014-2016

- Why:

- To support sociologists' and psychologists' work
- To provide support for a huge number of applications, services and systems
- For on-line and off-line knowledge extraction by visual data



A long story.....

- 1997-2000 MIT Alex Pentland: PFINDER projects and understanding interactions
- 2006- datasets for action analysis (Weizmann ICCV2005)
- 5 workshops on HBU (from 2010: IAPR, AMI, IROS, ACM MM, ECCV)  
<http://www.cmpe.boun.edu.tr/hbu/2014>
- Chalearn workshops 2011- 2015; CVPR 2015 challenge “Looking at People”
- Now many datasets

978

R. Poppe / Image and Vision Computing 28 (2010) 976–990

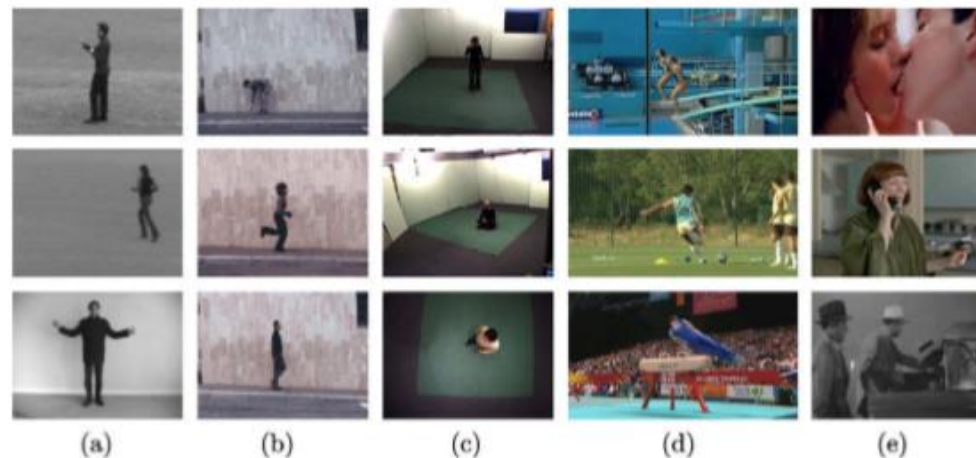
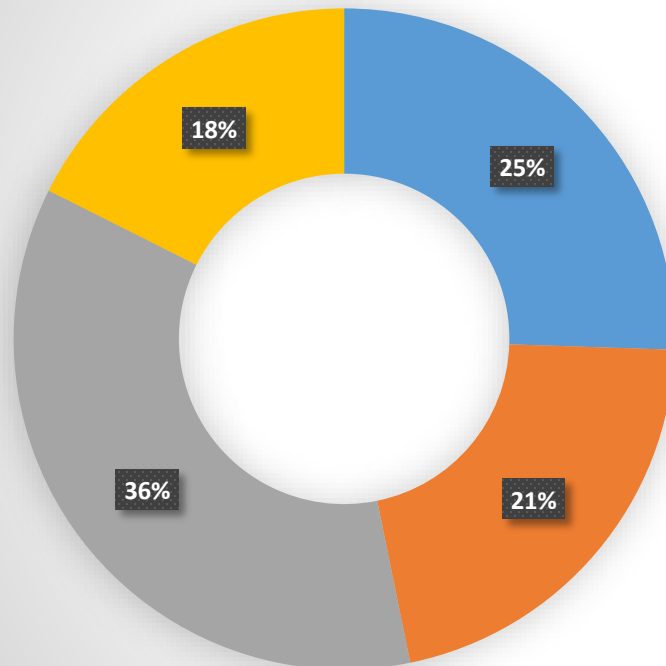


Fig. 1. Example frames of (a) KTH dataset, (b) Weizmann dataset, (c) Inria XMAS dataset, (d) UCF sports action dataset and (e) Hollywood human action dataset.

- From R. Poppe “A survey on vision based action recognition” Image and vision computing 2010

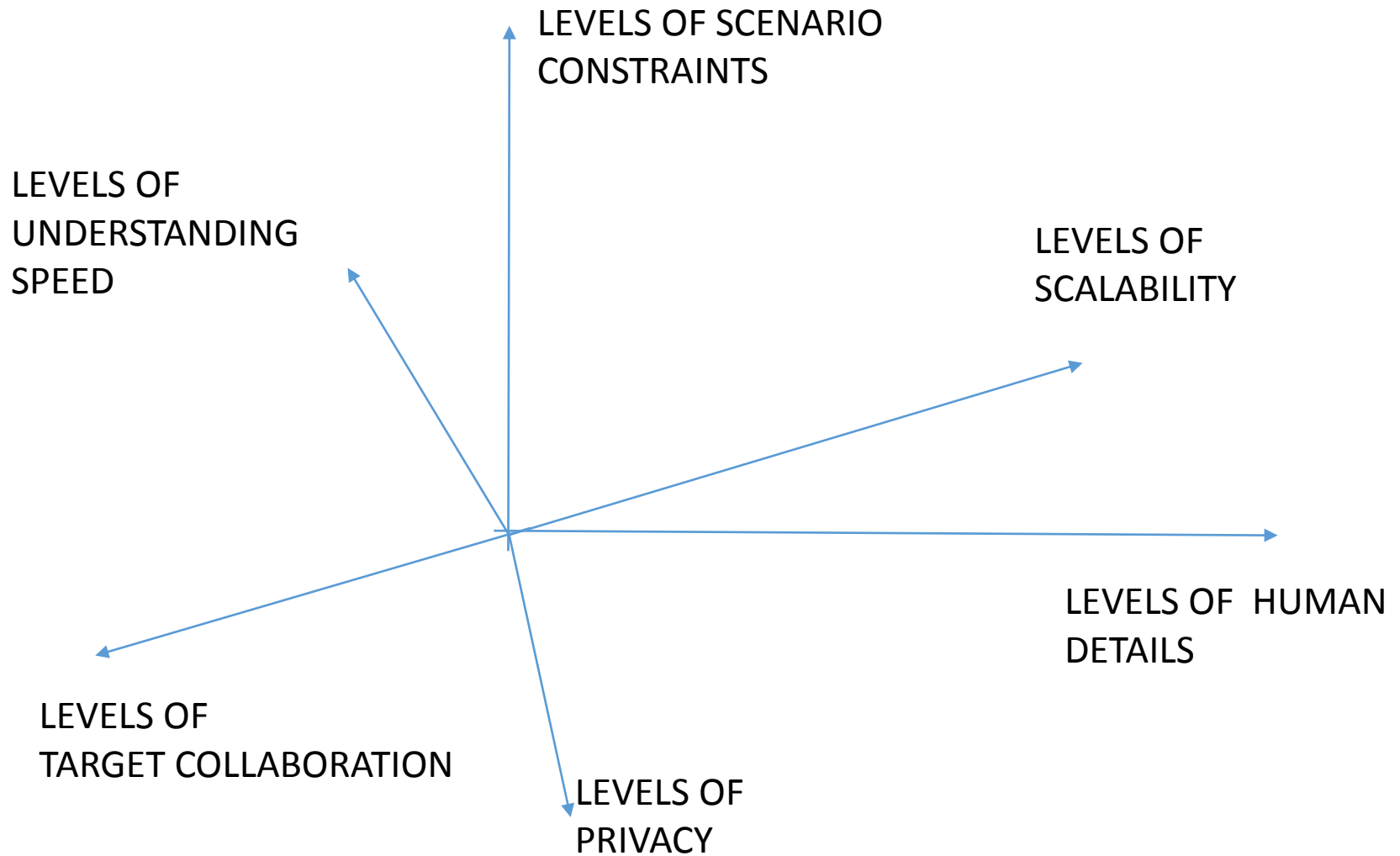
- HBU by vision:
- More than 1000 papers from 2010 to 2015

1000 scientific papers from google scholar 2015



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

# Different dimensions..



# A DIMENSION: the levels of details

- Levels of details



Human  
Full body



Human(s)  
In the environment



Humans  
In crowd





# Lev 1: Expressions and gestures

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



What are they doing?

Deaf Sign language



Picture from Benjamin Lewis UCLA

## Lev2: body actions



**Fig. 7.** Sample input frame of the Weizmann dataset



Roberto Vezzani, Davide Baltieri, and Rita Cucchiara, HMM Based Action Recognition with Projection Histogram Features ICPRW2010 supported by EU THIS Project

# Lev. 3 people in the environment

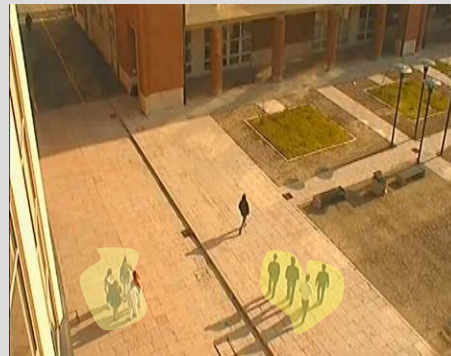
- What are they doing?



Real-time surveillance

(big) Data analysis  
for digital forensics

ENVI-VISION



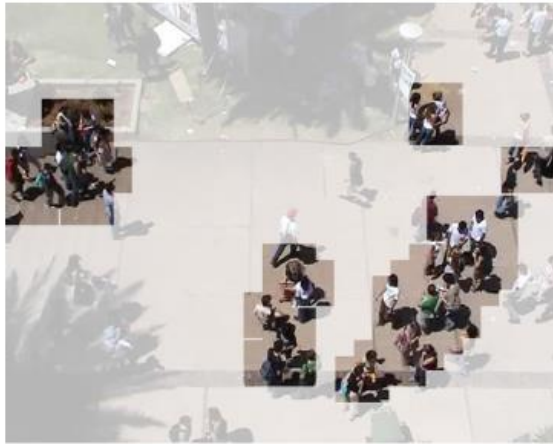
EGO-VISION



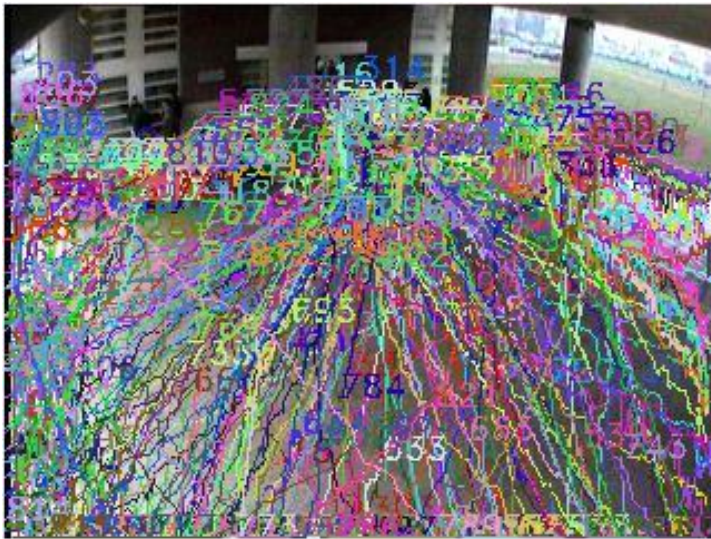


## Lev 4: social activities and behavior

- Detecting, tracking and understanding groups in crowd



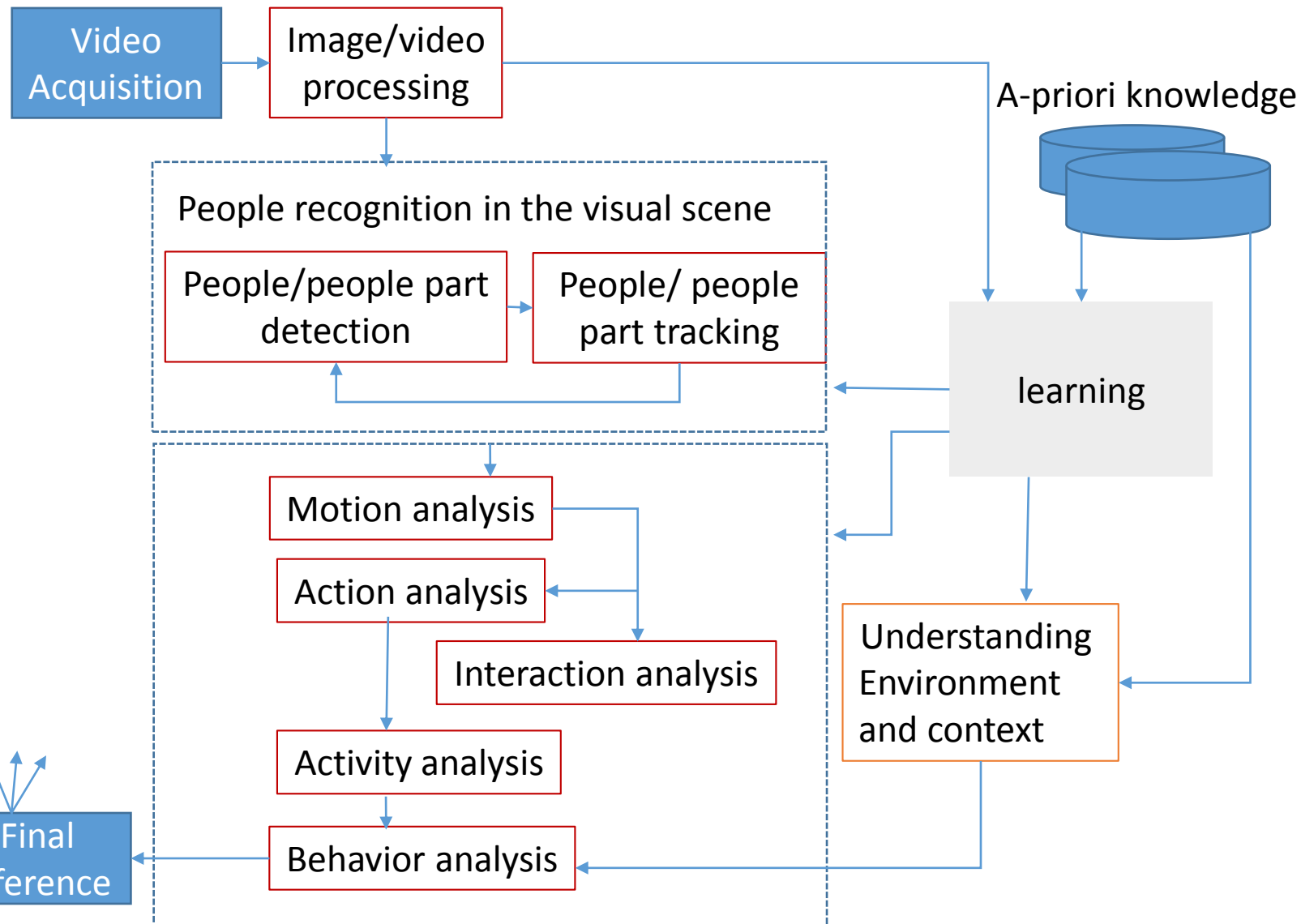
What are they doing?



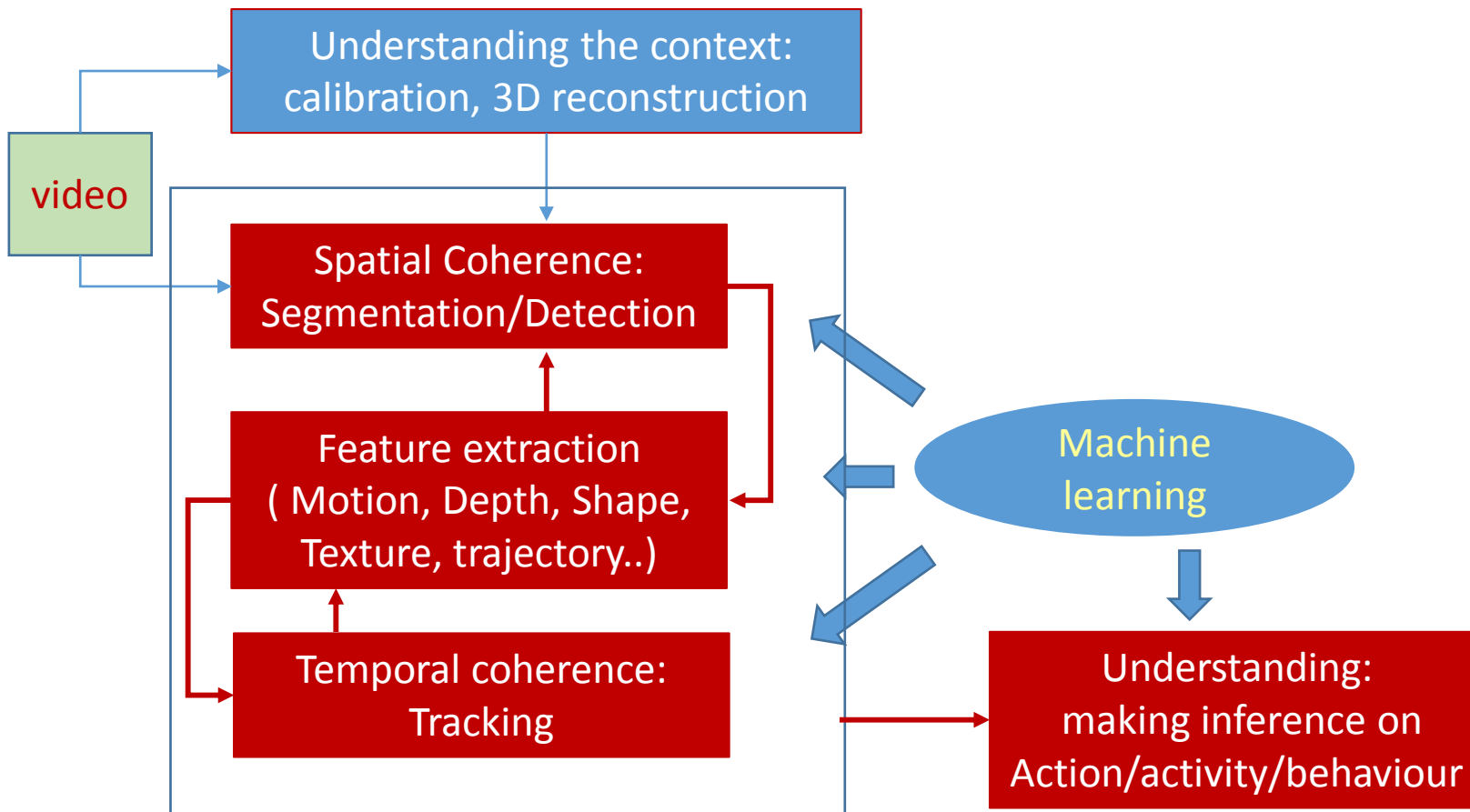
Alex Prager Crowd #8 (city Hall) 2013



# Common framework



# The big challenges



A big role of machine learning

2015 30.000\$ in prizes from Microsoft in Looking at People competition  
Special issue on PAMI 2015

# A second dimension: a constrained scenario

- The (1-...) constraints in the scenario

Health ,  
Sport training



Human  
Computer  
Interaction



Surveillance



Multimedia

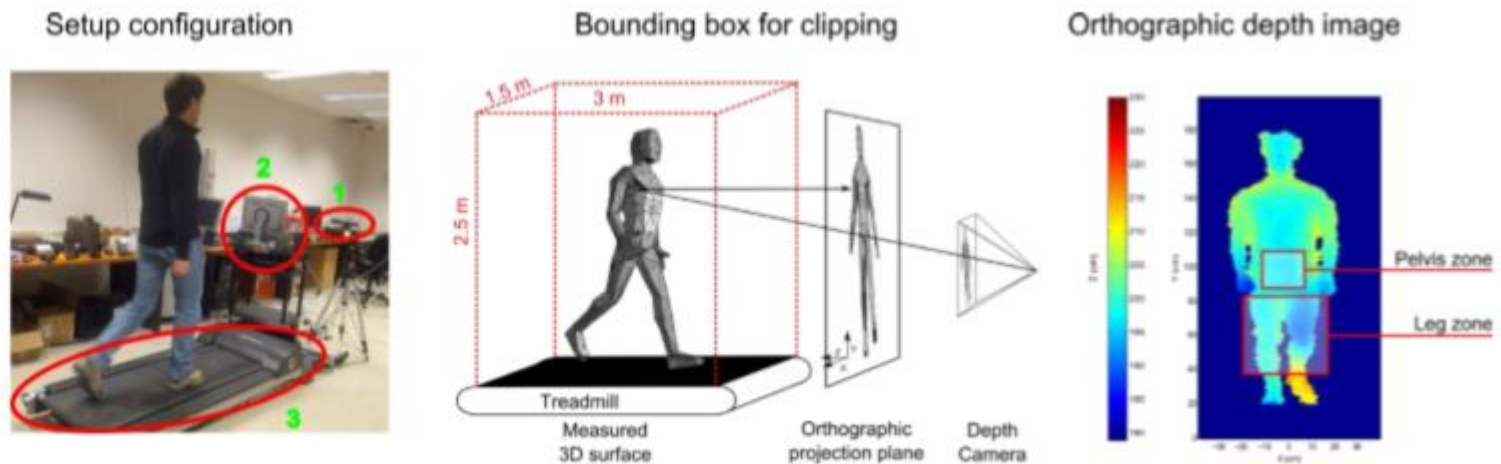


Humans  
In crowd



# 1 constrained applications: eg health

- Acquisition in a very constrained environment
- From marker-based to markerless 3D clouds
- Specific setup, Reconstruction and measures
- learning as a support of final diagnosis




Edouard Auvinet, Franck Multon and Jean Meunier; New Lower-Limb Gait Asymmetry Indices Based on a Depth Camera, Sensor 2015




# 2 multimedia video annotation

- Unconstrained scenario
- Detecting people behaviour in (educational) video
- From Italian national broadcast RAI
- Deep learning for concept detection


















UNIMORE  
UNIVERSITÀ DEGLI STUDI DI  
MODENA E REGGIO EMILIA


Scene segmentation demo



Detail: 100%

Hide annotations

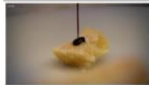
 blur · food · still life · nobody	 portrait · two · men · wheelchair	 energy · wire · sunset · evening
 two · adult · chef · men	 food · vegetable · meal · dish	 food · produce · fruit · nobody
 dairy product · fruit · food · sweet	 food · chef · commerce · store	 facial expression · restaurant · group · men
 european · female · one · portrait	 fashion · european · dress · museum	 fine art · painting · sculpture · european
 sculpture · fine art · statue · tuscany	 fine art · painting · god · religion	 sky · religion · architecture · building




Andrew Graham-Dixon and chef Giorgio Locatelli travel through Italy exploring the country's history, culture, food, art and landscape. Their journey begins in Bologna, the capital of Emilia-Romagna, one of the richest regions in Italy. They find out why the city is known as la Dotta (the Learned), la Grassa (the Fat) and la Rossa (the Red), while visiting its shops, art institutions and the oldest university in the world.

food

Drag your shots here (double click to remove)

  
dairy product · fruit · food · sweet

  
food · vegetable · meal · dish

## 2.HCI and collaborative applications

- **Human computer interaction**

- *( general assumptions... often too general)*

- **Less constrained environment** but with strong assumptions (e.g. the person knows what has to be done)
- **Real-time processing** (often with embedded solutions) but with an acceptance temporal window
- Collaborative environment: goal **high precision**; if recall is not enough, human-in-the-loop can handle it
- **General purpose** but a-priori defined features ( eg hand colors, point trajectories..)
- **Learning by very few examples** ( often by a single person only) but the Learning space is well defined (eg with few ambiguities)

→ all Kinect applications in HCI!

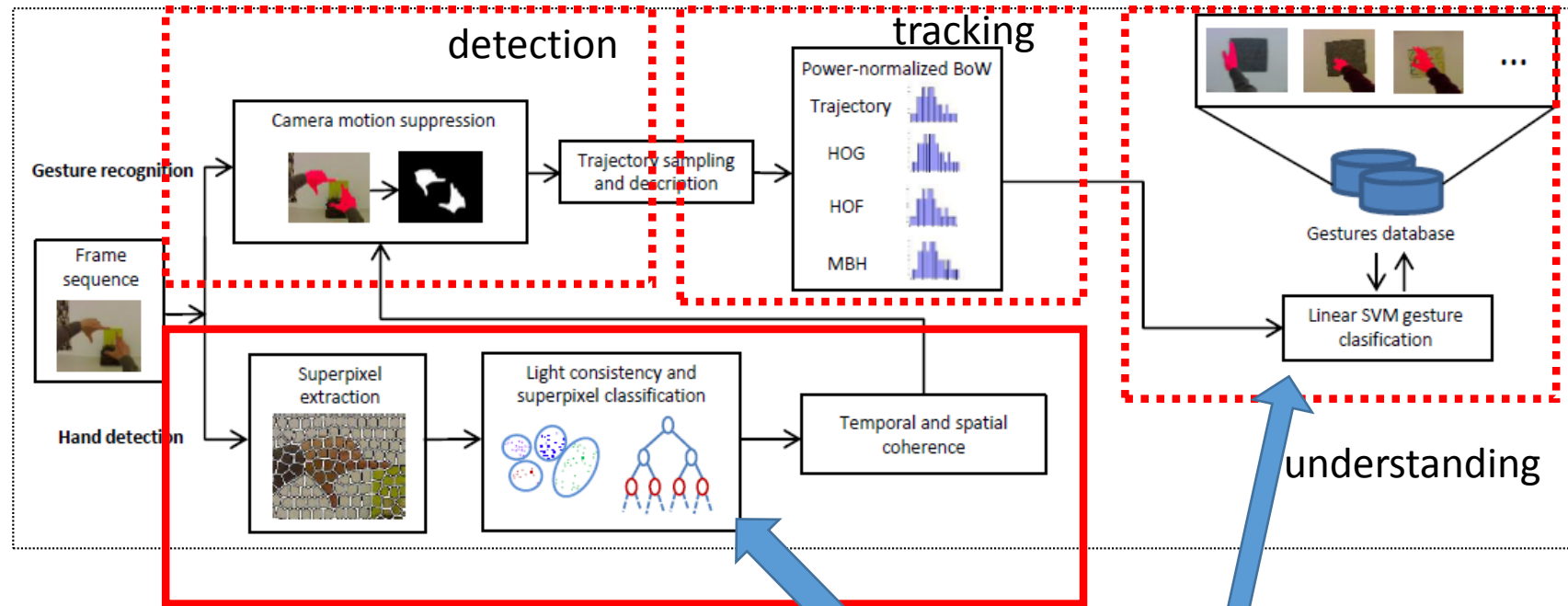
*the gesture recognition market is estimated to grow at a CAGR of 29.2 % year from 2013-2018 for Gaming and Entertainment, Healthcare, Automotive applications, Educations and serious gaming (By USA Markets and Markets 2014)*

## An example:

- Understanding human behaviour in interacting with artistic objects for augmented experiences
- HBI: Gesture analysis, associated with context



# Ego-Gesture recognition



- 1) (Ego-)Hand detection
- 2) (Ego-)Camera motion suppression
- 3) Feature extraction
- 4) Classification and gesture recognition

Machine learning



### 3. Surveillance

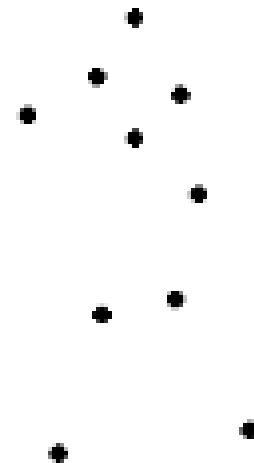
#### People Surveillance

one of the biggest topic in computer vision

- On-line, fast, not collaborative at all.
- Partially constrained ( eg. calibration)
- Segmentation/ detection
- Simple features (**motion**, appearance)
- Tracking: the big challenge
- Learning and high reasoning: with many examples



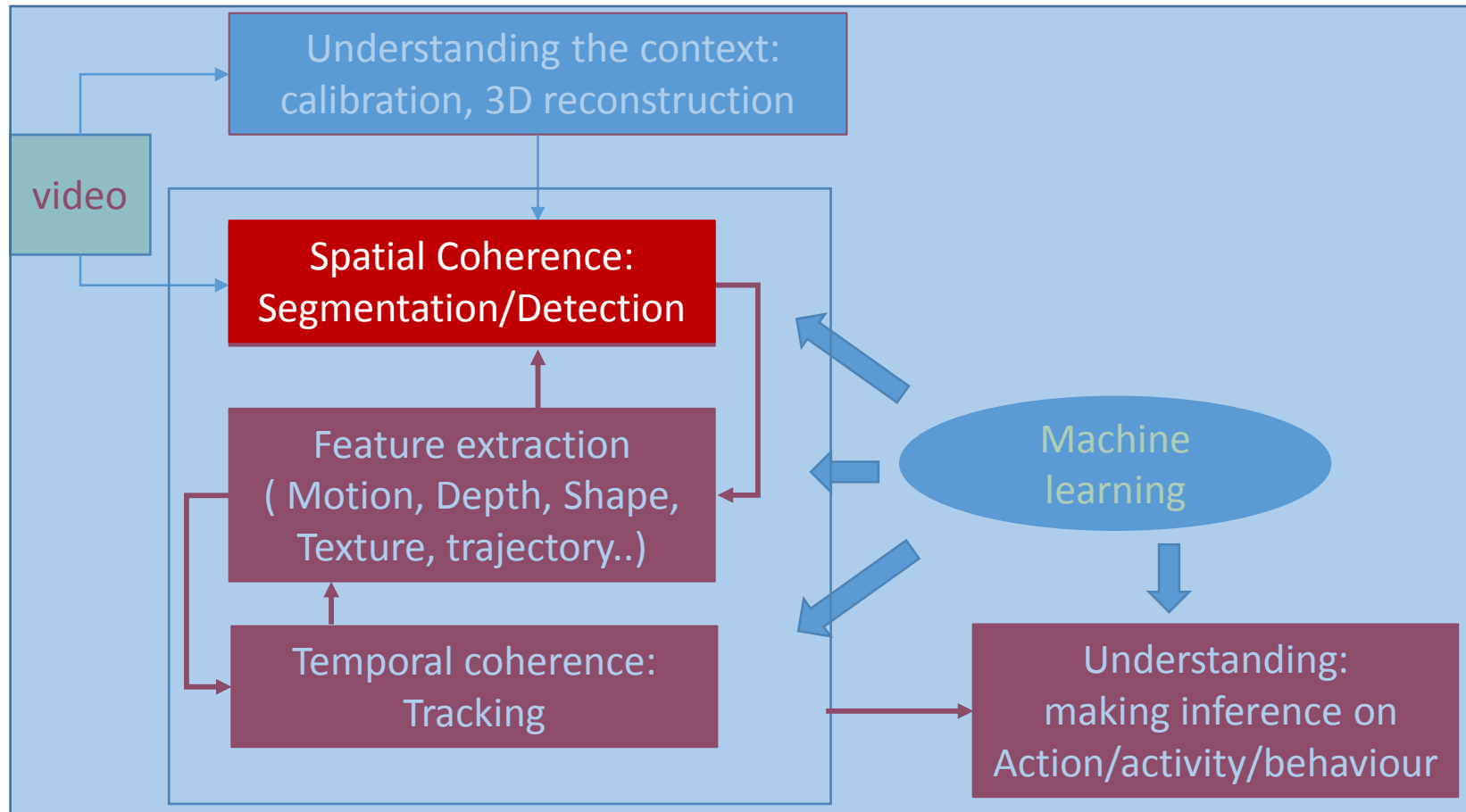
*Humans recognize  
motion and recognize  
by motion*



For single people and crowd:

## Why HBU for surveillance?

- Monitoring dangerous/forbidden zone
- Access control
- Single people iterative activity recognition
- Extracting common behaviour in crowd
- Understanding anomalous behaviours ( eg unfrequent trajectories)
- Recognizing specific behaviour (e.g. suspicious behaviour, for terrorism, social engagement for children..)

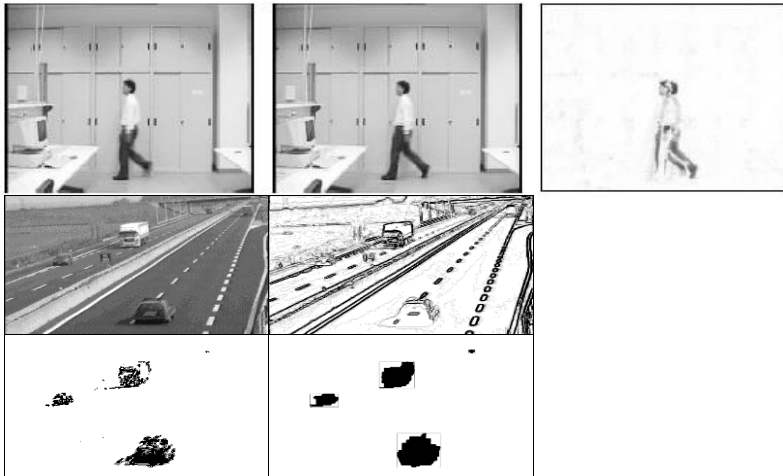


# Segmentation/detection in surveillance

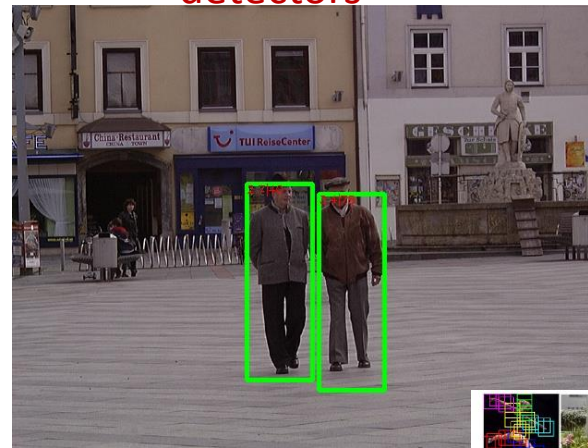
Segmentation by  
motion

Detection with  
learning and  
classifiers

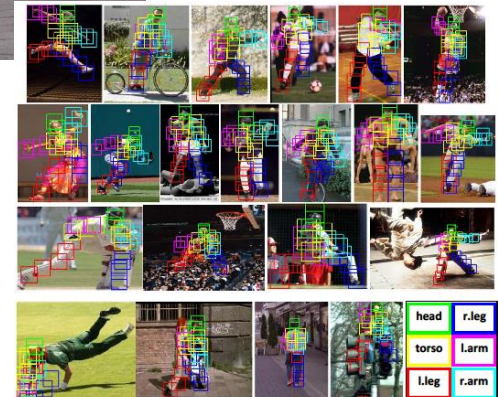
## Differential Motion



## HOG, part based.. People detectors



## Background suppression

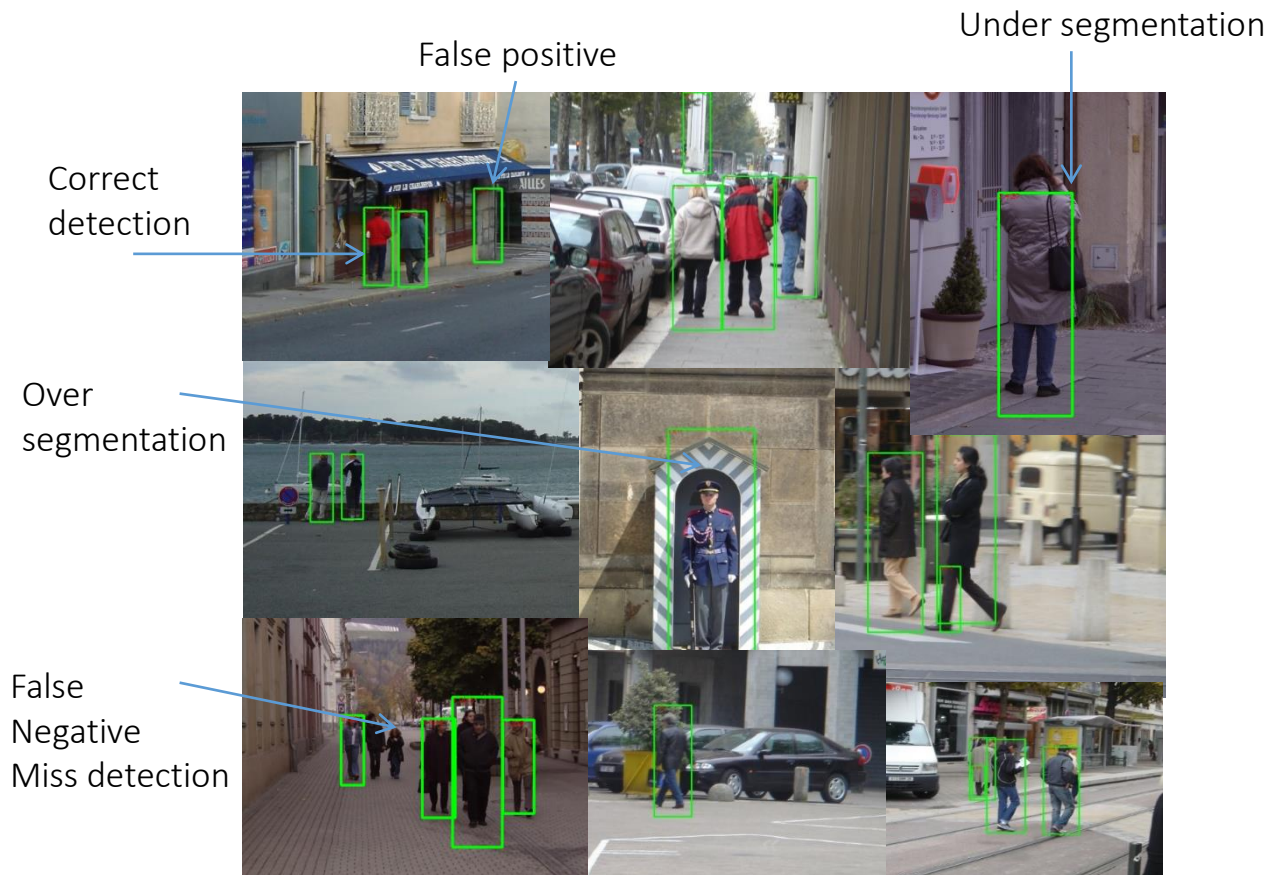




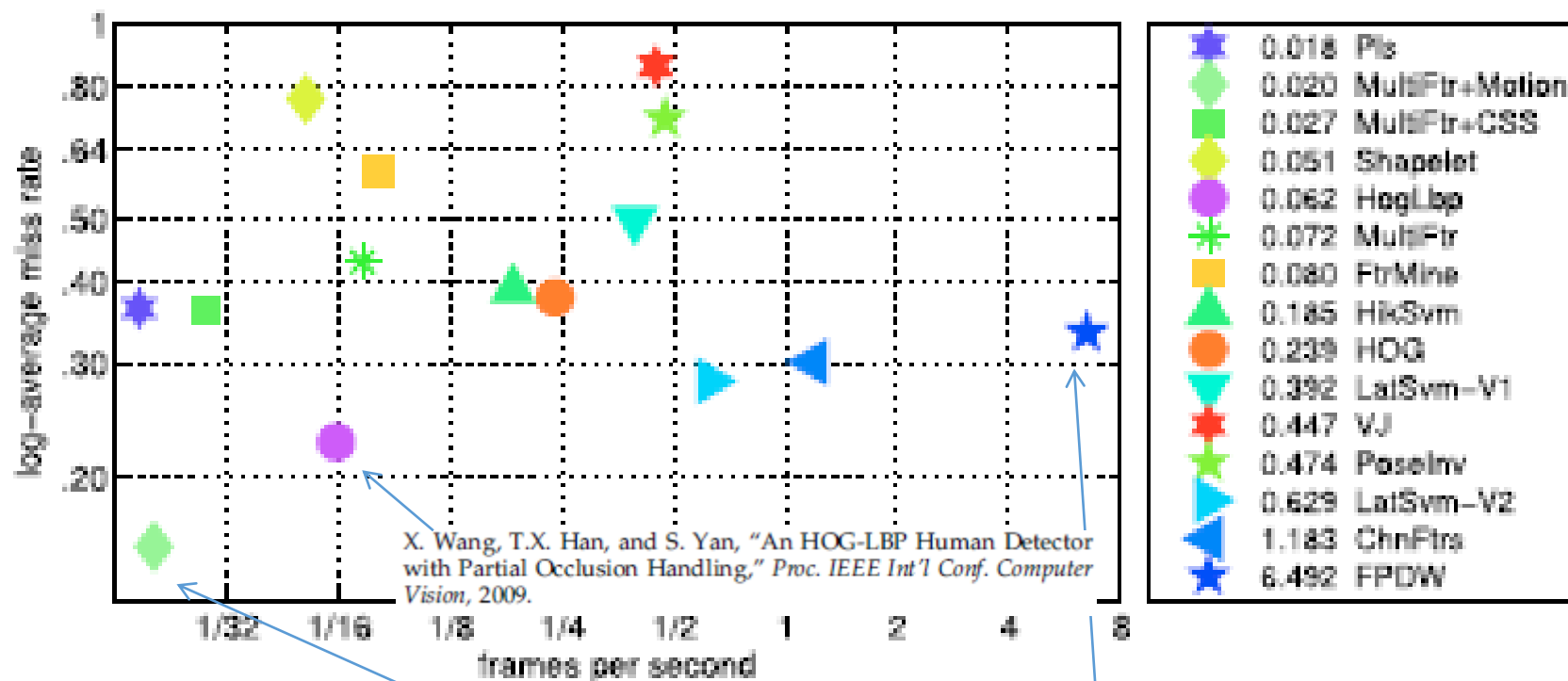
# People detection

- 😊 very general for pedestrian
- Without any constraints: from mobile, moving cameras, wearable...)

- 😞 someErrors....



# Speed and accuracy



(a) accuracy versus runtime for pedestrians over 100 pixels



(a) Caltech [3]

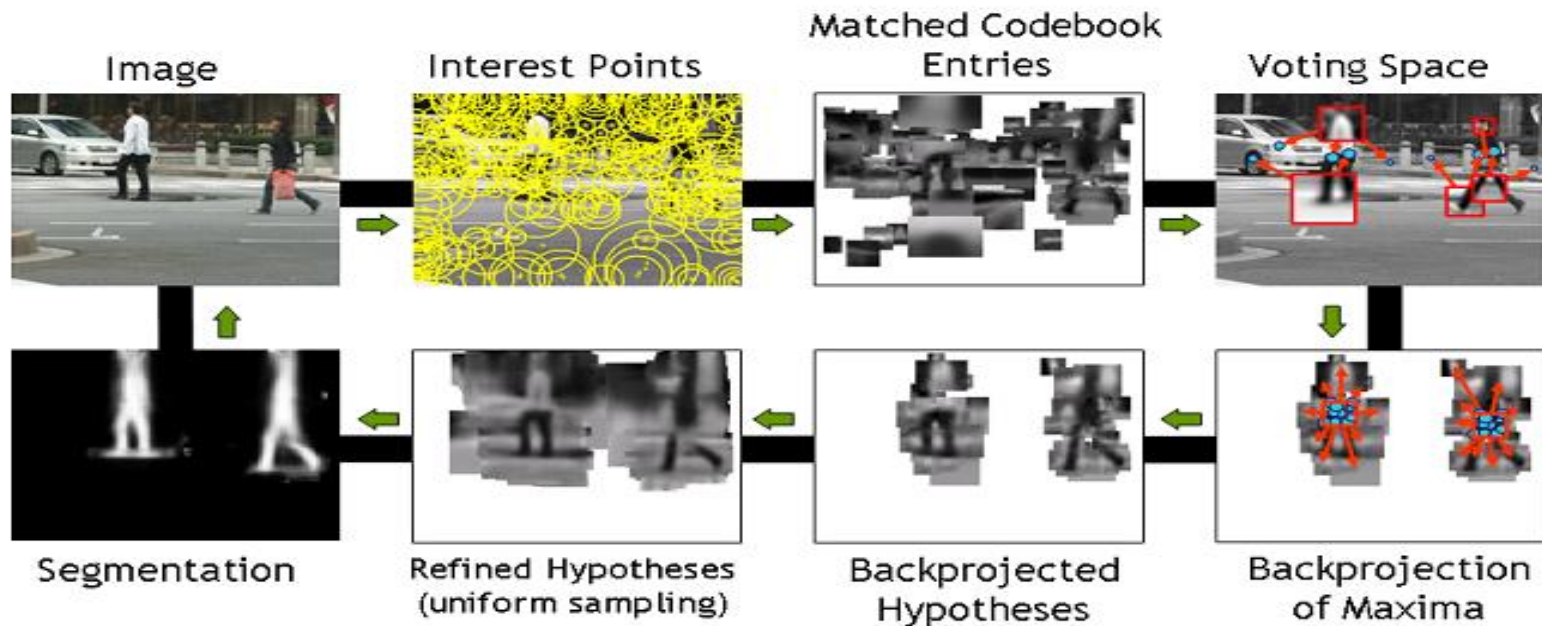
(b) Caltech-Japan [3]

S. Walk, N. Majer, K. Schindler, and B. Schiele, "New Features and Insights for Pedestrian Detection," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2010.

# Detecting people & target

1) The initial attempts: Generate hypotheses (local detector)

- SIFT, HOG etc



# Detecting people & target

- 2) Most of the approaches: better to see everywhere: **Sliding windows**



Thanks to Derek Hoiem Illinois univ.



# What the detector sees

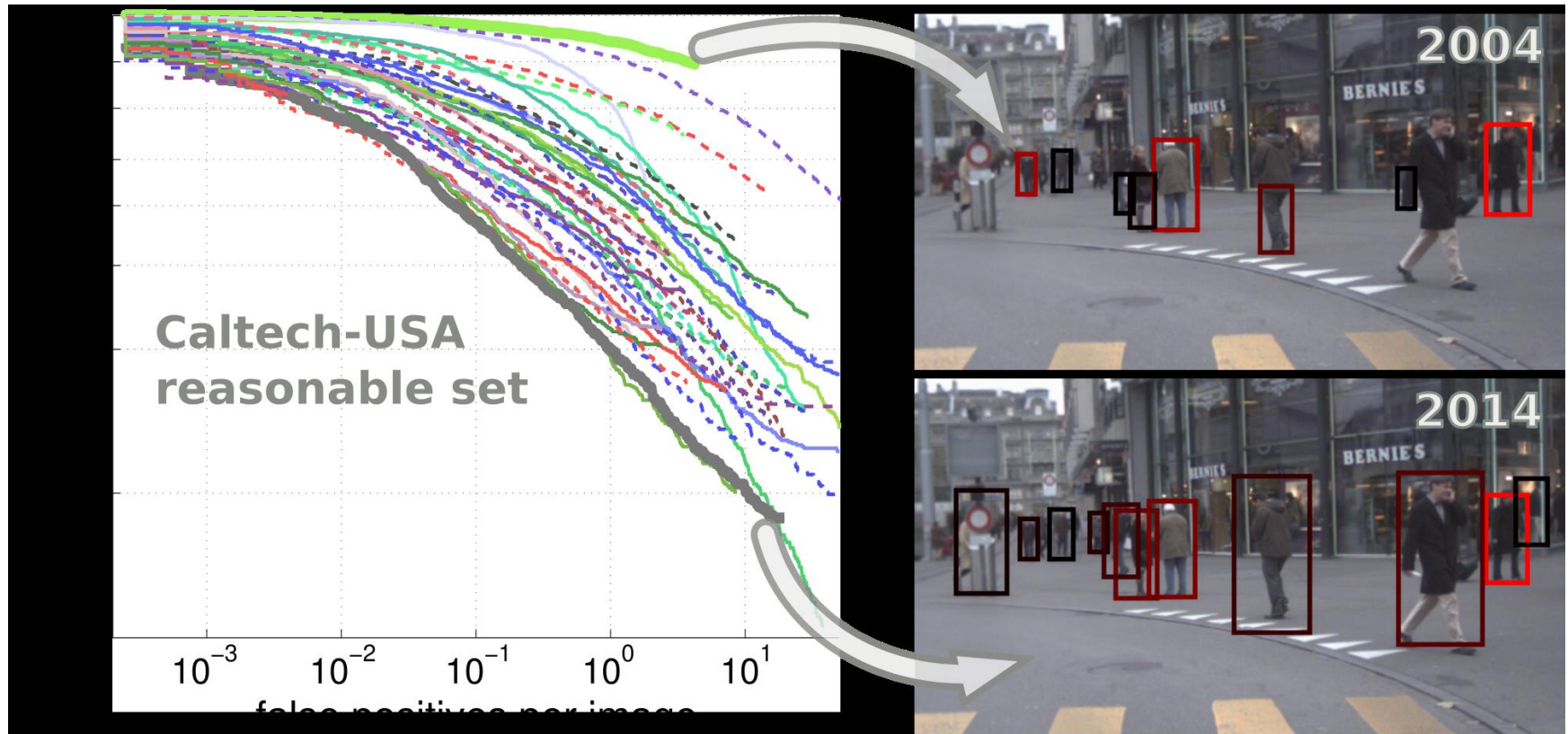


- And repeated at each possible scale... and then, learning

Thanks to Derek Hoiem Illinois univ.

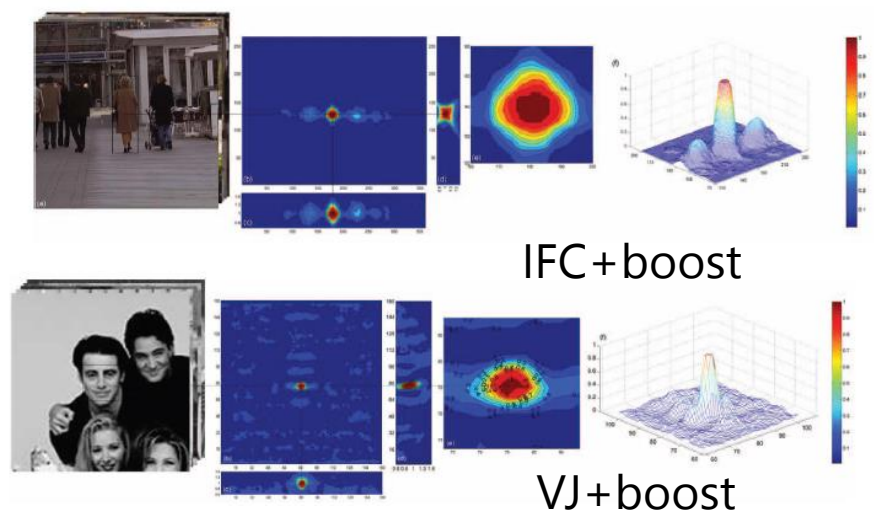
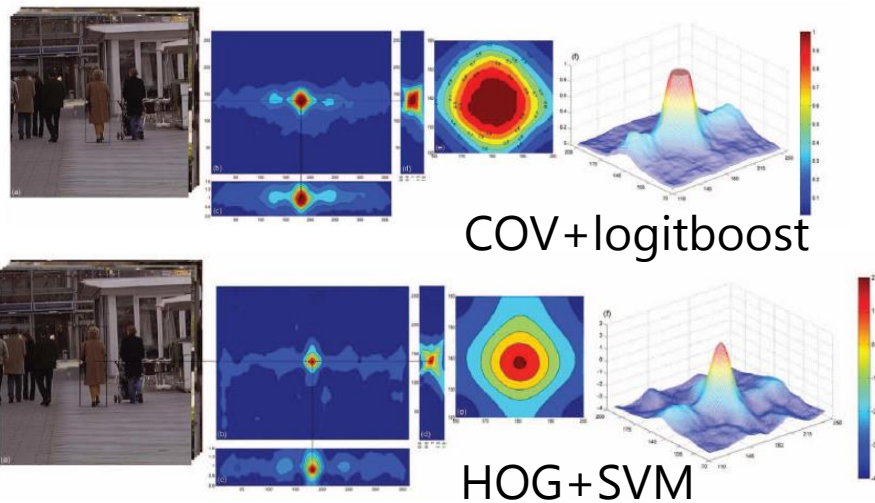
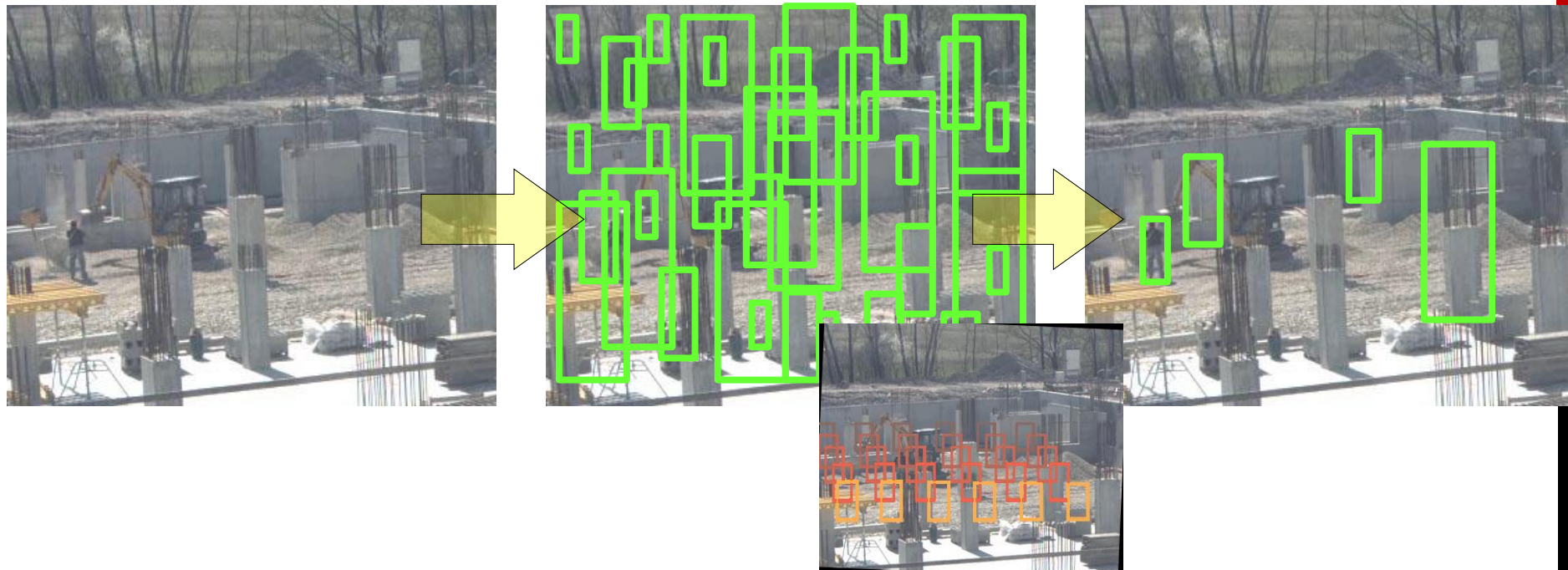
# A conclusive study

- Benenson.. Schiele. Ten years of pedestrian detection, what have we learned ? ECCV2014





# Speed and accuracy



A probabilistic bayesian paradigm for object detection:  
“*estimate obj. detection as a pdf*”

Set  $q_0(\mathbf{X}) = U(\mathbf{X})$

for  $i=1..m$  do

Draw  $N_i$  samples from  $q_{i-1}(\mathbf{X})$

Assign a Gaussian kernel to each sample

Compute the measurement on each sample  $s_i^{(j)}$

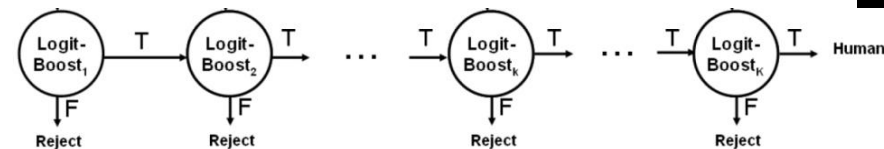
Obtain the measurement density function at step

$$p_i(\mathbf{Z}|\mathbf{X}) = \sum \pi_i^{(j)} \cdot \mathcal{N}(s_i^{(j)}, \Sigma_i^{(j)})$$

Compute the new proposal distribution:

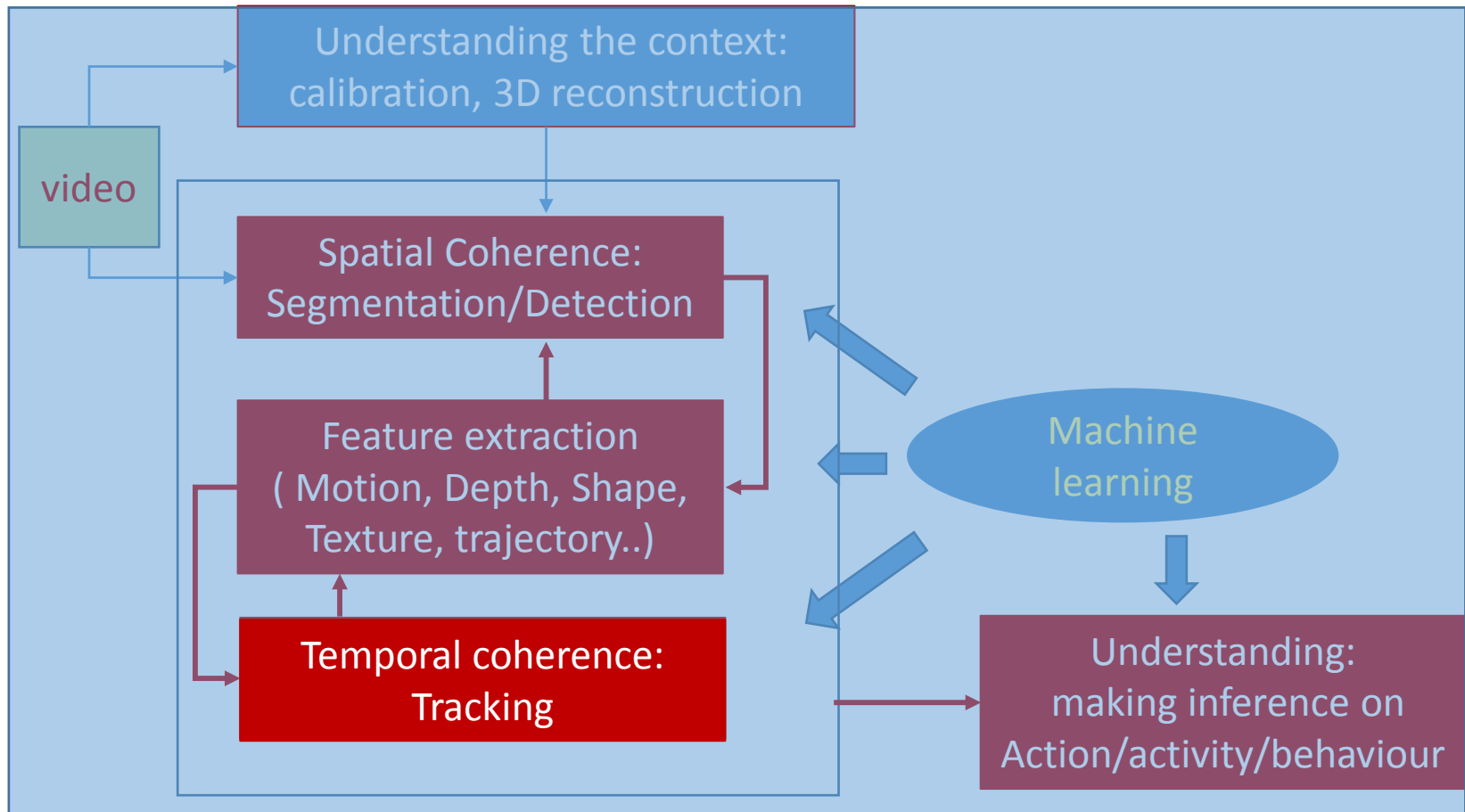
$$q_i(\mathbf{X}) = (1 - \alpha_i) q_{i-1}(\mathbf{X}) + \alpha_i \frac{p_i(\mathbf{Z}|\mathbf{X})}{\int p_i(\mathbf{Z}|\mathbf{X}) d\mathbf{X}}$$

end for



**G.Gualdi, A.Prati, R.Cucchiara** Multi-Stage Particle Windows for Fast and Accurate Object Detection  
IEEE Transactions on PAMI Aug. 2012





# Detection and tracking

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

In semi-constrained world  
Tracking is possible



- Tracking is the hardest problem.
- Finding the visual invariance among frames
  - Appearance?
  - Motion?
  - Space continuity?

- **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)

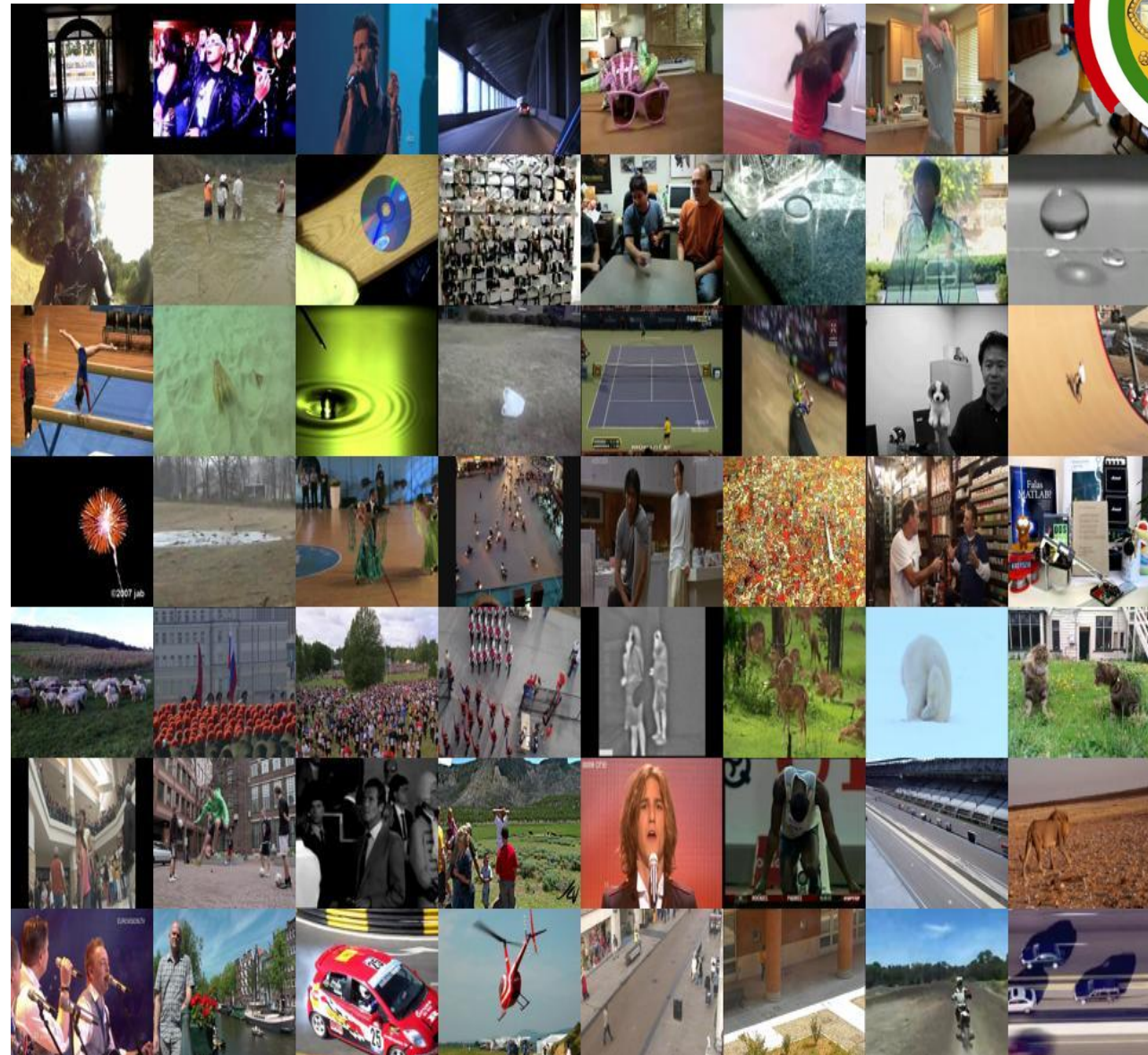
315 video ALOV++  
<http://www.alov300.org>  
<http://imagelab.ing.unimo.it/dsm>

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**  
[TPAMI 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

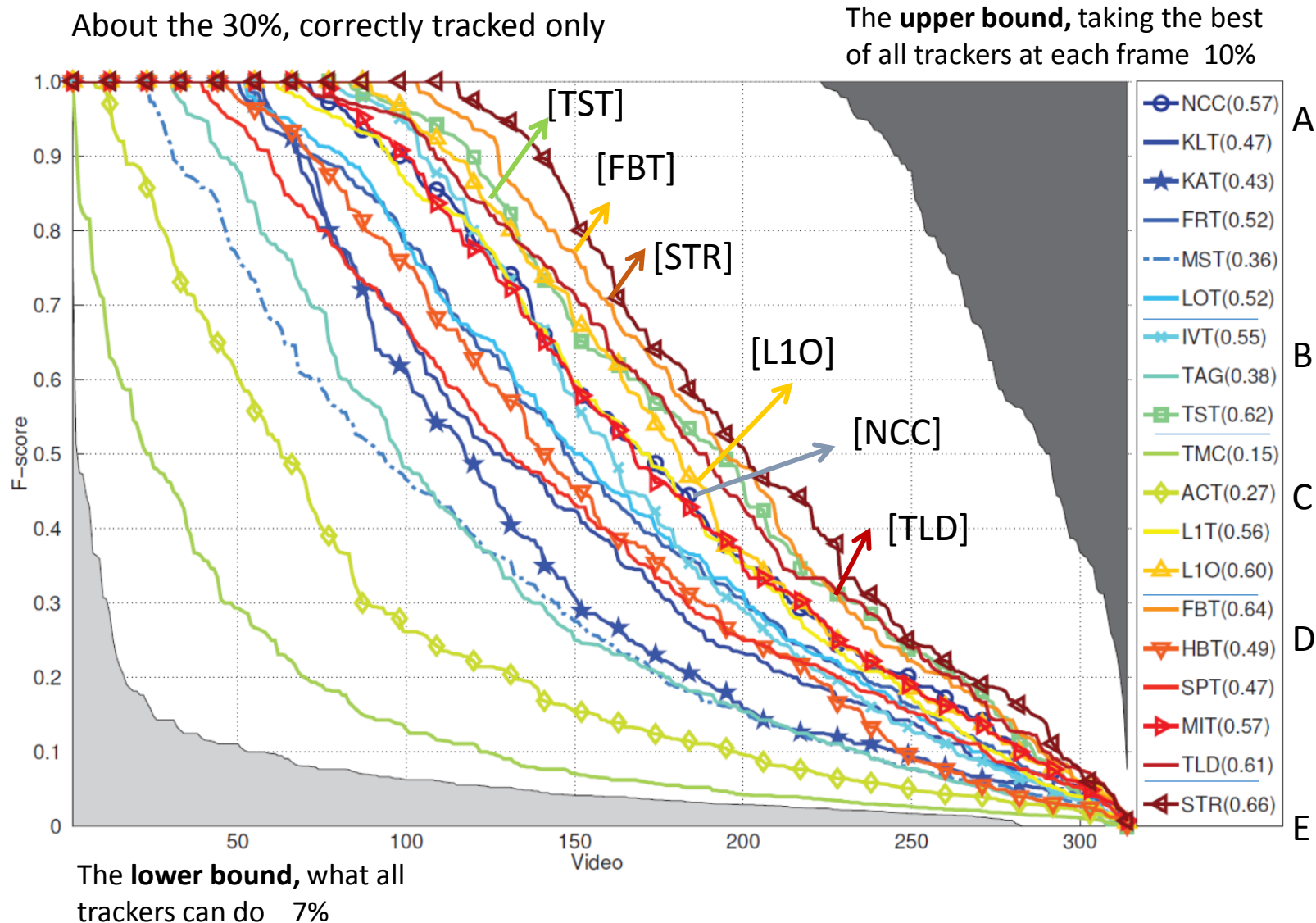
12-MOVINGCAMERA

13-ZOOMINGCAMERA

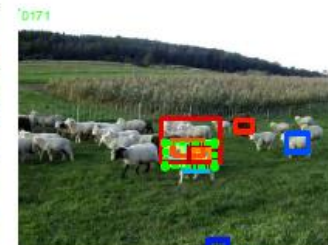
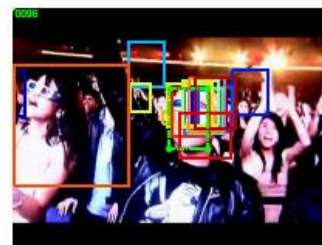
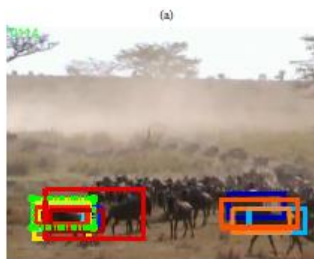
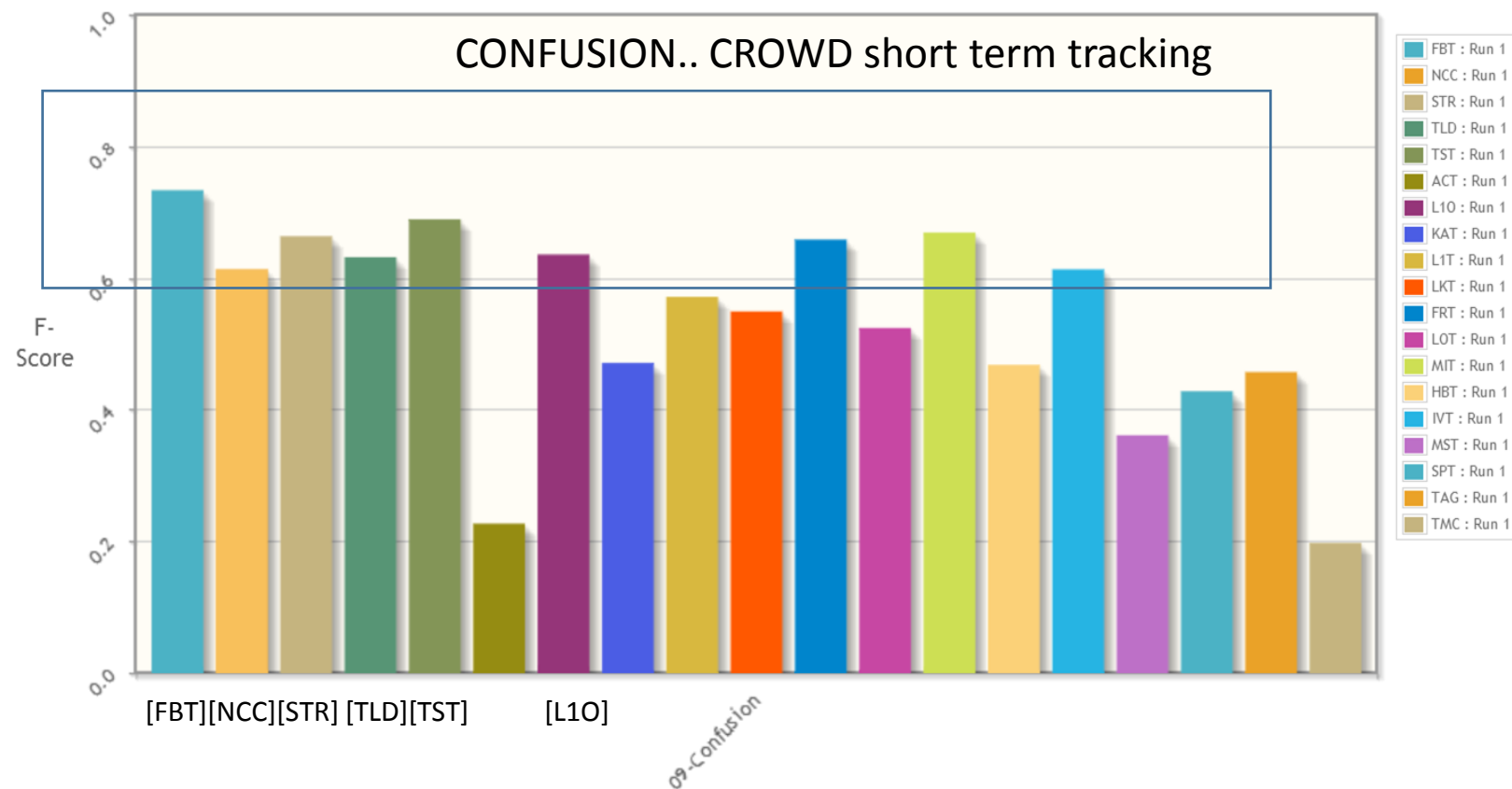
14-LONGDURATION



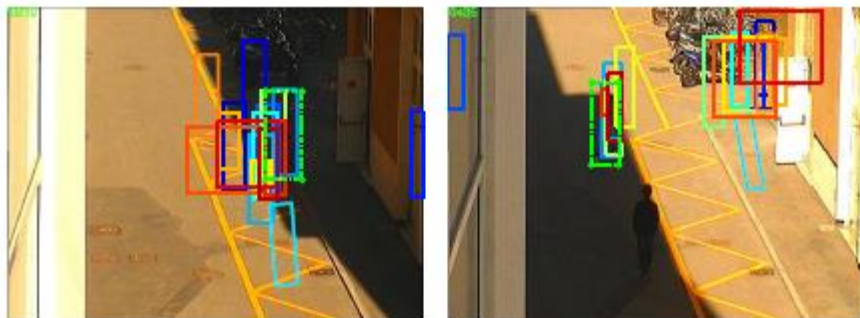
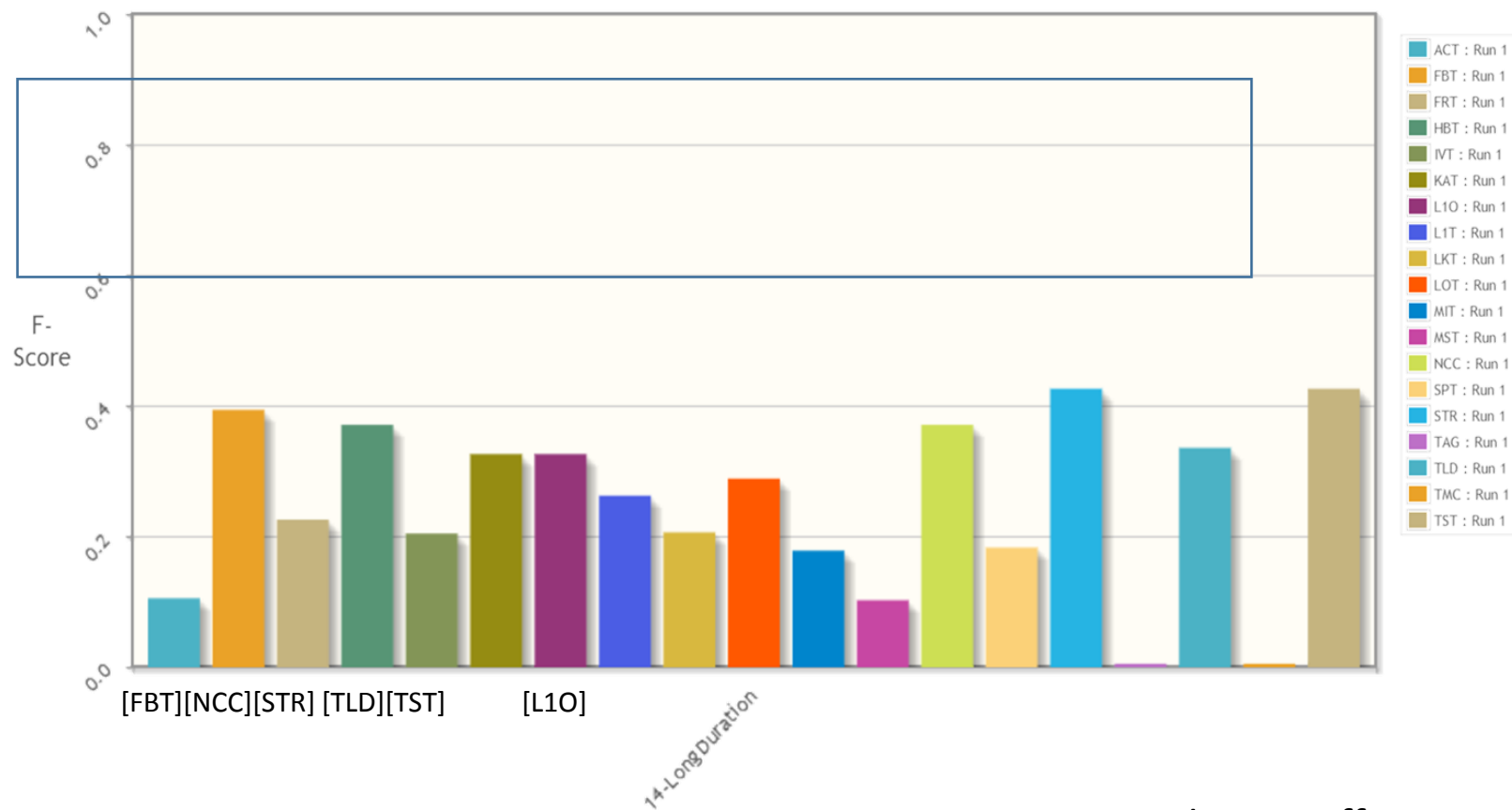
# A comprehensive view Survival curve



# challenge: trackers comparison

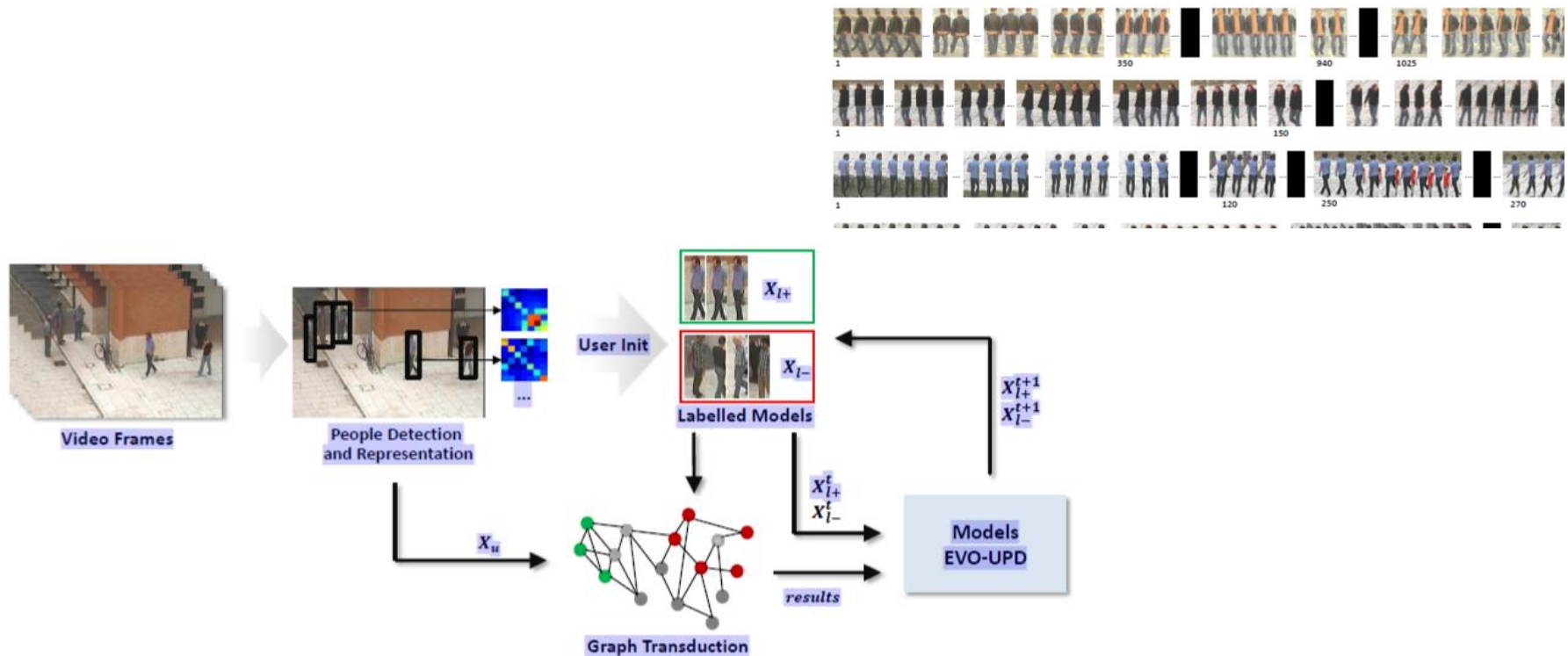


# challenge: trackers comparison



We need more effort

1. Tracking by background suppression: appearance based tracking
2. Tracking by detection









## Tracking (few) people

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

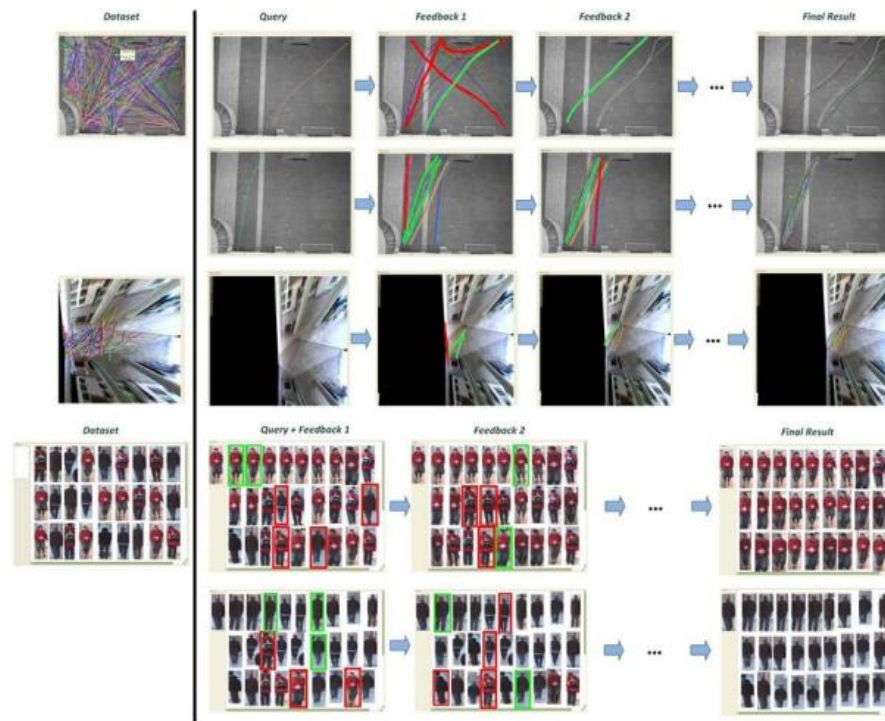


S. Calderara, R. Cucchiara, A. Prati, "Bayesian-competitive Consistent Labeling for People Surveillance" in IEEE Trans. On PAMI , 2008

- Searching from anomalous behaviour
- query systems for forensics applications

16

Dalia Coppi et al.



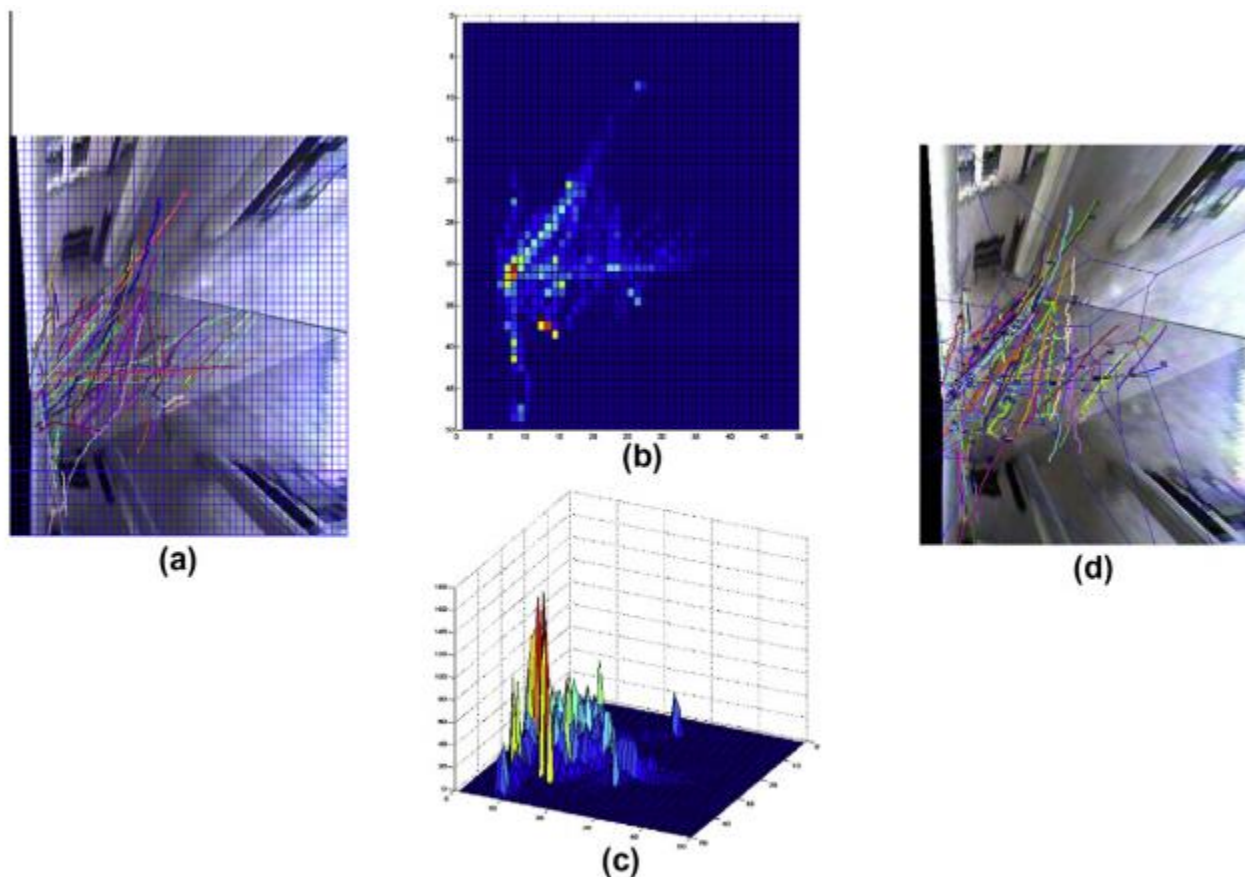
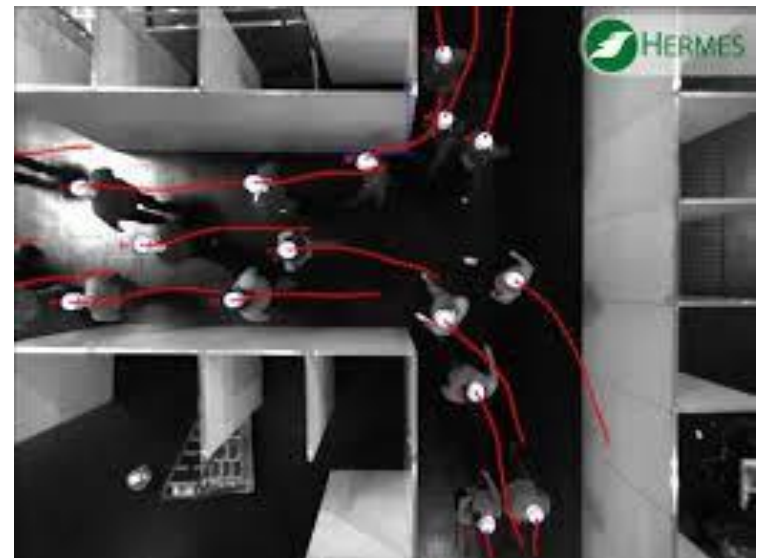
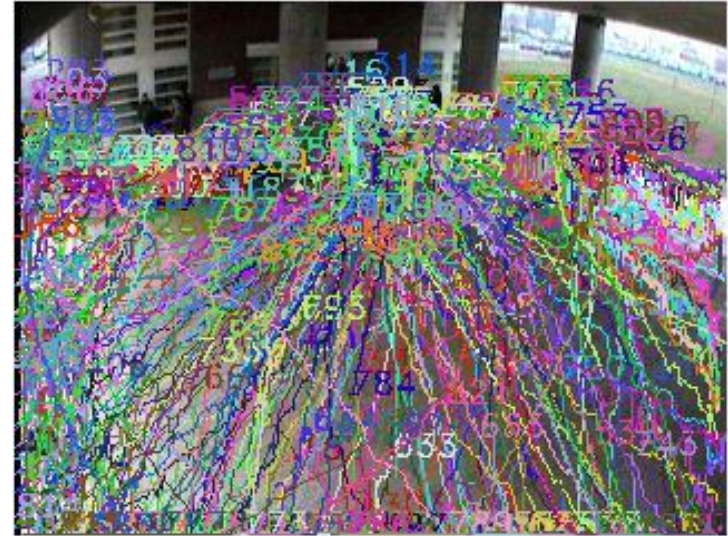


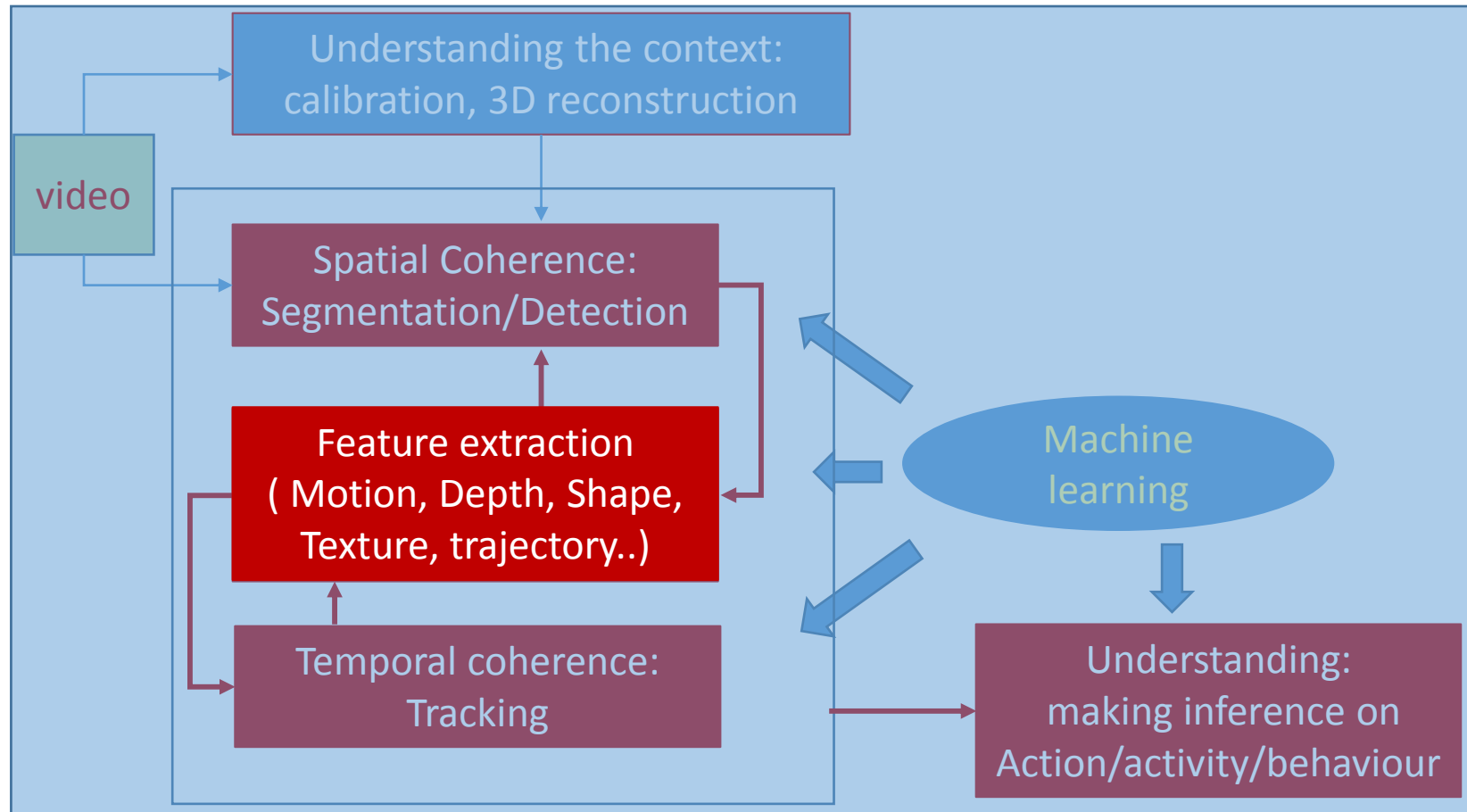
Fig. 2. Irregular partitioning of the image area through Voronoi diagrams: (a) Reports the first regular division of the image ( $50 \times 50 = 2500$  cells in this example); (b) shows the top view of the 2D histogram, while (c) shows a side view; and (d) shows the resulting Voronoi diagram with 50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

# Trajectory analysis

- Detection , tracking , trajectory acquisition
- With few people in the scene
- with controlled experiments  
( helmets) PEtrack

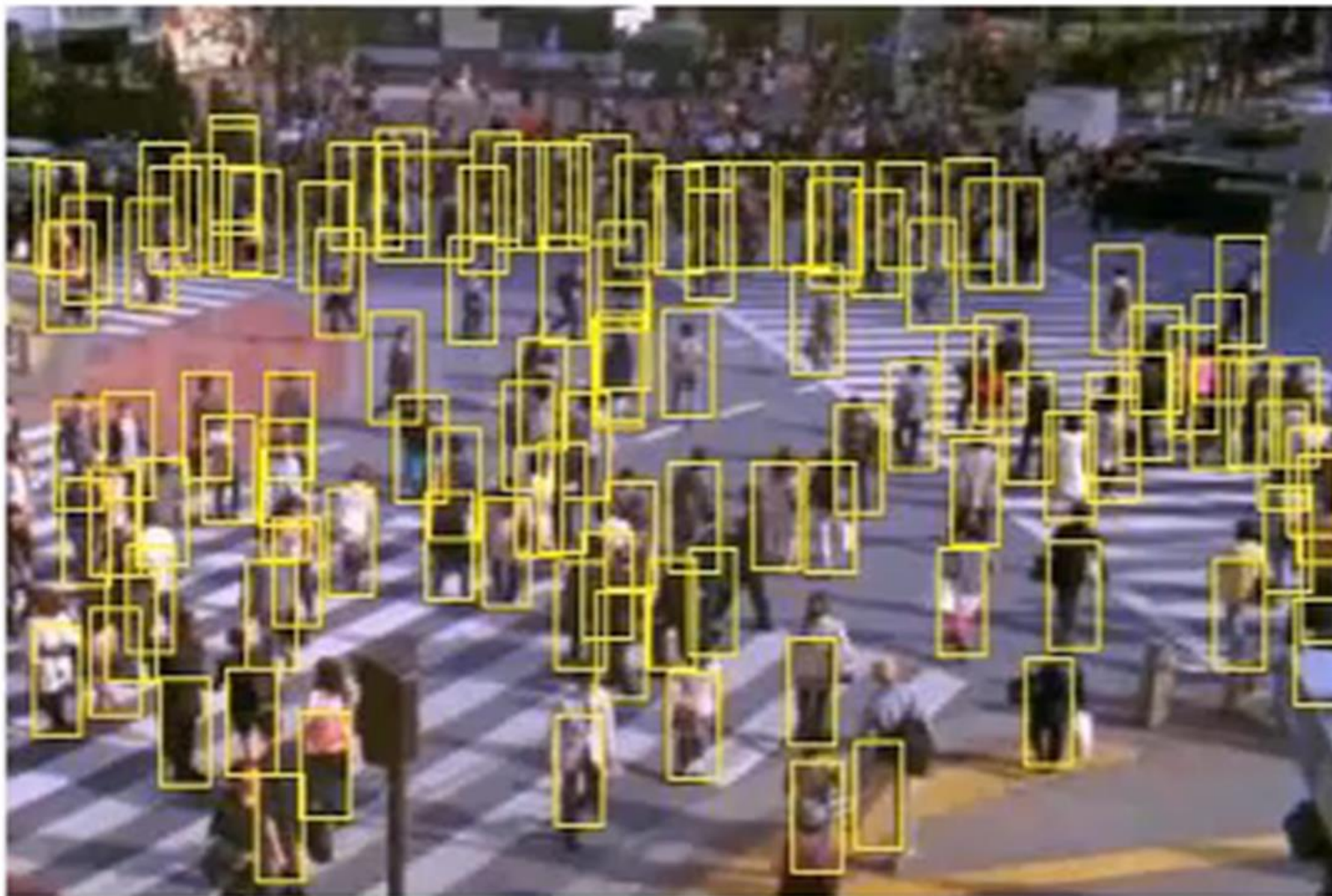


But in real crowd the problem  
is far to be solved..





# Detection and tracklets



- 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?**



# Detecting social groups in crowds

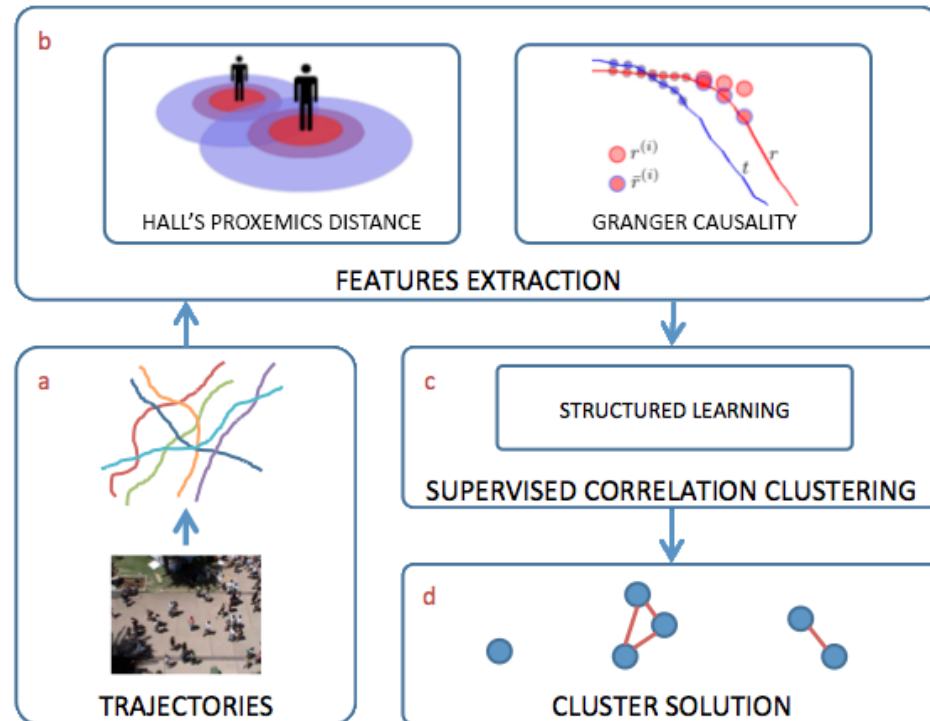
- Group detection: **learn to partition** into groups the pedestrians being part of a crowd observing pairwise relations and transivities.\*
- Integrating two cues:

## 1. HALL'S PROXEMICS

- Hall's proxemics theory<sup>1</sup> defines **reaction bubbles** around every individual and
- the **interaction between pairs** of individuals can be classified according to a quantization of their **mutual distance**

## 2. GRANGER CAUSALITY

- Intuition: two pedestrian belonging to the same group will probably influence each other position and direction!<sup>2</sup>
- The Granger causality test is a statistical hypothesis test for determining whether one time series is useful in forecasting another



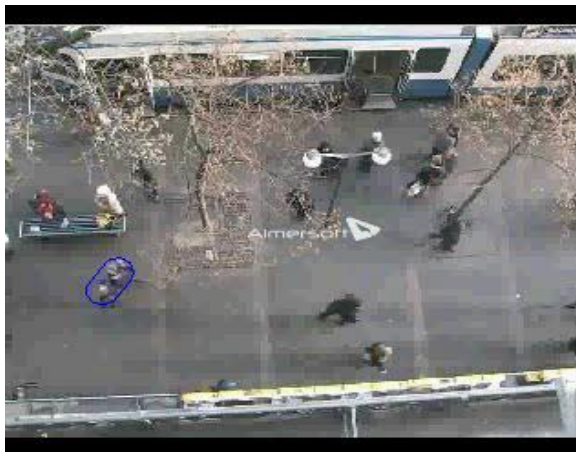
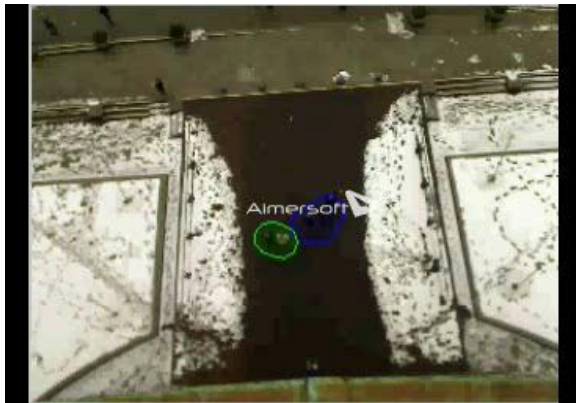
\* Structured learning for detection of social groups in crowd Solera, Calderara, Cucchiara, AVSS 2013



Features: Proxemics and Granger causality

Structure function: pair-wise correlation clustering

Group detection: Structured SVM [\[groups\]](#)





- Human behavior understanding in video: **a big challenge for big data**
- The problem is not in reasoning, or in data mining  
but in *extracting the correct spatial and temporal visual knowledge*
- Enormous improvements in very few years
- Some challenges “have been solved”
- Most of them have been not solved yet:
  - Human interaction with the environment in an unconstrained and not collaborative way
  - Understanding intention before behavior
  - Working on crowd and in cluttered scenario
  - Working in big footage from the web without text annotation
  - Working on streaming and scalable data....
  - .....

*The way ahead is still long...*



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

Interdipartimental Research Center in ICT  
Tecnopolo di Modena  
Emila Romagna High Technology Network



## PEOPLE



Rita Cucchiara



Costantino\_Grana



Roberto Vezzani



Simone Calderara



Augusto Pieracci



Giuseppe Serra



Paolo Santinelli



Martino Lombardi



Michele Fornaciari



Dalia Coppi



Marco Manfredi



Francesco Solera



Simone Pistocchi



Fabio Battilani



Patrizia Varini



Stefano Alletto,