# ELEC 2885:
# Image Processing and Computer Vision

Artificial Intelligence in computer vision.

christophe.devleeschouwer@uclouvain.be

# Outline

O **Introduction**
- **Terminology: AI / ML / DL;**
- AI4Vision:
  - From pixels to semantic;
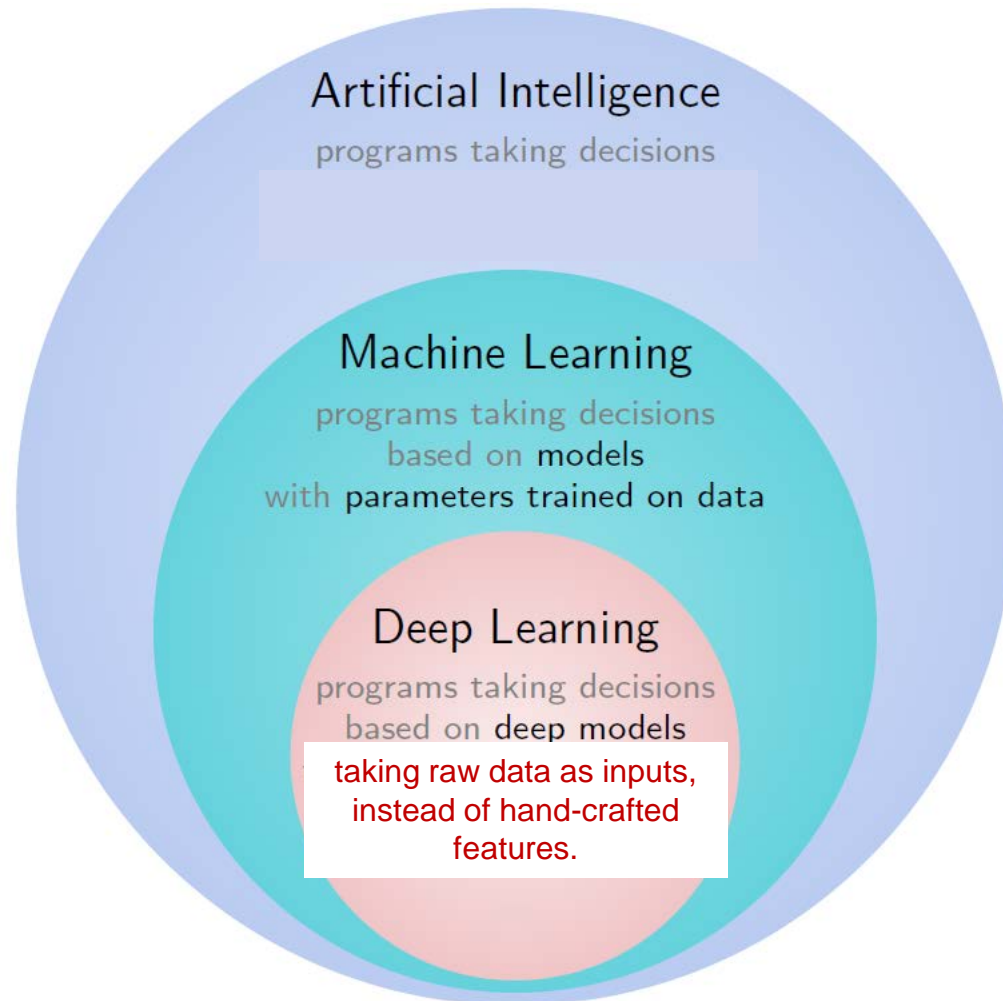  - Two paradigms: (AI or ML) vs. DL.

O Hand-crafted image features
- Pixel- or patch-based features;
- Interpreting the features:
  - AI: 3D Geometry;
  - AI: Graph-based topology;
  - ML: Random ferns, AdaBoost, Support Vector Machines.

O DL in vision : Convolutional Neural Networks (CNNs)
- Paradigm shift & revolution;
- CNNs basics;
- CNNs heads & backbones: network architecture examples;
- CNNs limitations.

# Terminology

**Artificial Intelligence**

programs taking decisions

**Machine Learning**

programs taking decisions
based on **models**
with **parameters trained on data**

**Deep Learning**

programs taking decisions
based on **deep models**

taking raw data as inputs,
instead of hand-crafted
features.

# Outline

O **Introduction**
- Terminology: AI / ML / DL;
- **AI4Vision:**
  - From pixels to semantic;
  - Two paradigms: (AI or ML) vs. DL.

O Hand-crafted image features
- Pixel- or patch-based features;
- Interpreting the features:
  - AI: 3D Geometry;
  - AI: Graph-based topology;
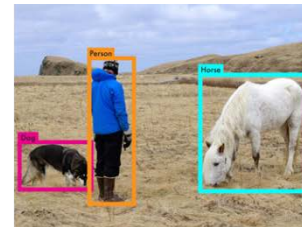  - ML: Random ferns, AdaBoost, Support Vector Machines.

O DL in vision : Convolutional Neural Networks (CNNs)
- Paradigm shift & revolution;
- CNNs basics;
- CNNs heads & backbones: network architecture examples;
- CNNs limitations.

# AI4vision: a variety of tasks


**Classification**


**Detection**
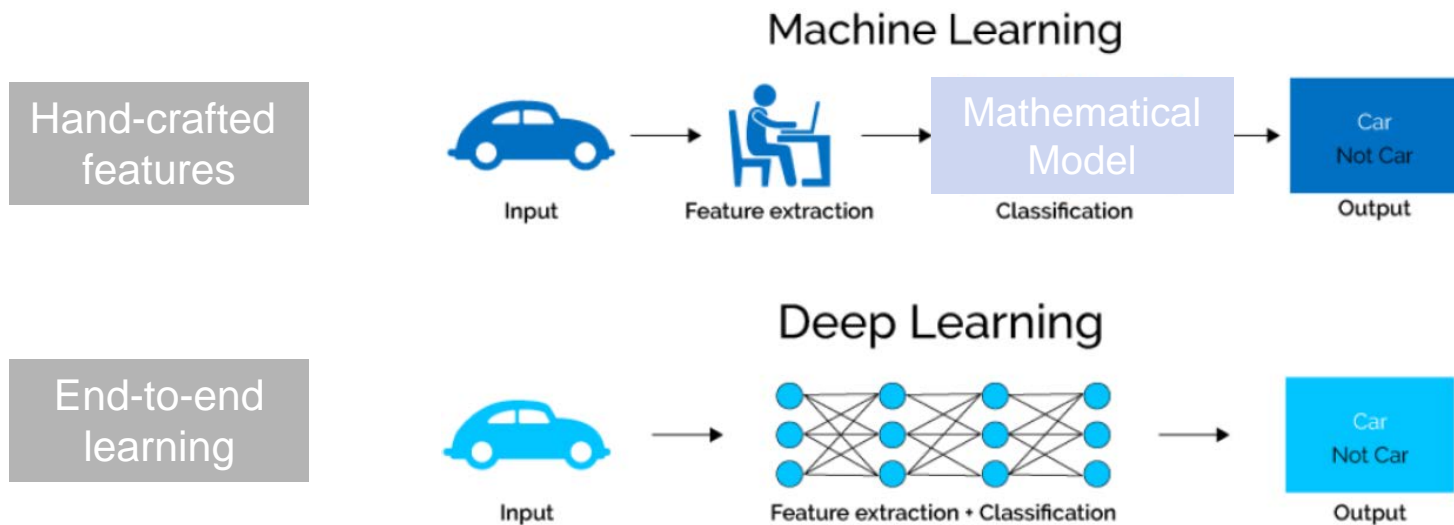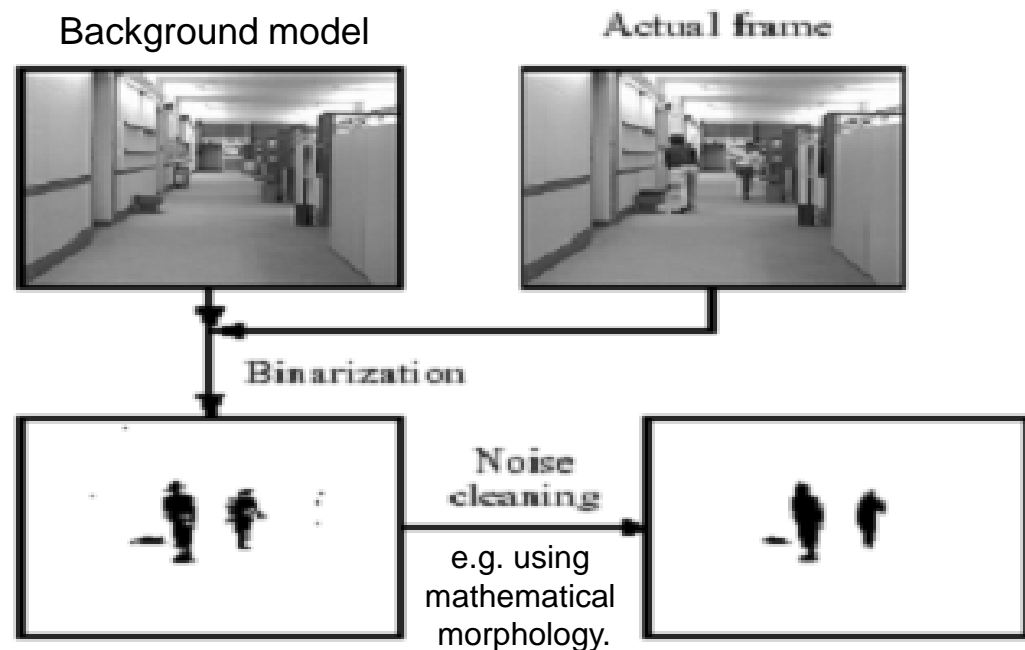

**Class or instance segmentation**


**Pose estimation**

# AI4vision: implementation

O **General principle: from pixels to semantic**
- collect pixel-wise features (foreground mask, HoG, local binary features, feature map CNN);
- build on topology / geometry / *training* to turn low-level local information into high-level semantic knowledge.

O **Two paragdims**

## Machine Learning

| Hand-crafted features | Input | Feature extraction | Mathematical Model | Car Not Car |
|---|---|---|---|---|
| | | | Classification | Output |

## Deep Learning

| End-to-end learning | Input | Feature extraction + Classification | Car Not Car |
|---|---|---|---|
| | | | Output |

# Outline

O  Introduction
- Terminology: AI / ML / DL;
- AI4Vision:
  - From pixels to semantic;
  - Two paradigms: (AI or ML) vs. DL.

O  **Hand-crafted image features**
- **Pixel- or patch-based features;**
- Interpreting the features:
  - AI:    3D Geometry;
  - AI:    Graph-based topology;
  - ML:    Random ferns, AdaBoost, Support Vector Machines.

O  DL in vision : Convolutional Neural Networks (CNNs)
- Paradigm shift & revolution;
- CNNs basics;
- CNNs heads & backbones: network architecture examples;
- CNNs limitations.

# Image pixel features

O  Color components

O  Gradient vector (see 'Segmentation' lecture)

O  Foreground label

# Foreground pixels in video

O  Key assumption: still cameras !!!
  - Realistic in many practical scenarios.

O  Principle:
  - Background pixel is estimated as the most frequent observation among past ones;
  - Foreground pixel detected if the observed pixel does not fit the background model.

O  Key advantage:
    low complexity.

Background model      Actual frame

Binarization

Noise cleaning

e.g. using mathematical morphology.

# Background model

O  Exponentially weighted moving average :

- Construct a background image B as average of few images, e.g. using an exponentially weighted moving average: $B(t) = (1-\alpha) B(t-1) + \alpha I(t)$ ;

- For each frame I, classify individual pixels as foreground if $|B-I| > \tau$.

O  Gaussian Mixture Model :                                    *Not covered this year.*

(Adaptive background mixture models for real-time tracking, Stauffer and Grimson CVPR98)

- Advantage: is able to model multiple appearances for a background pixel;
- Each pixel is modeled as the sum of  $K$ weighted Gaussians,
  with weight $w_{k,t}$ , mean $\mu_{k,t}$, and standard deviation $\sigma_{k,t}$, $k < K = 3{\sim}5$.

- The weights reflects the frequency of occurrence. Hence, large weights correspond to background modes ;

- Model is updated adaptively with new observations, and some learning rate $\alpha < 1$.

New Observation $X_t$ $\longrightarrow$ Matching and model updating $\longrightarrow$ background

$\searrow$ foreground

○ **Matching Criterion**

$$\left\| X_t - \mu_{k,t} \right\|^2 < \beta \cdot \sigma_{k,t}^2$$

- If no match found: *foreground*
- If match found:
  - *B first modes with highest $w_k$ are background,*

  $$B = argmin_b \left( \sum_{k=1}^{b} \omega_k > T \right)$$

  - *others are foreground.*

○ **Update Formula**

- Update active ($M_k$=1) and non-active ($M_k$=0) weights :

  $$w_k = (1 - \alpha) \cdot w_k + \alpha \cdot M_k$$

- Update Gaussian for active mode *m*:

Match found:

$$\mu_{m,t+1} = (1 - \alpha) \cdot \mu_{m,t} + \alpha \cdot X_t$$

$$\sigma_{m,t+1}^2 = (1 - \alpha) \cdot \sigma_{m,t}^2 + \alpha (X_t - \mu_{m,t})^T (X_t - \mu_{m,t})$$

No Match found:
   Replace Gaussian with smallest weight by new observation.

# Foreground mask weakness



## Noise !

• illumination changes;
• shadows;
• motion, e.g.due to wind.

# Image patch features

## Why image patches ?

**The patch might correspond to a _pixel_ neighborhood, and the purpose is to decide whether the pixel lies inside, outsise or between two cells.**





**SLIDING WINDOW ALGORITHM**

**JUDYBATS**

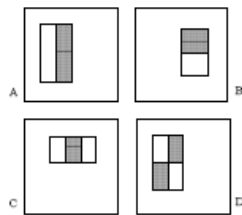**The patch might be a _window of semantic interest_ scanned across the image at multiple scales and locations to detect objects (eg. faces or pedestrians).**

# Image patch _binary tests_, considered as weak classifiers

O Pixel comparison (as for BRIEF):



$$b_i = \begin{cases} 1 & if\ I(\bullet) < I(\bullet) \\ 0 & else \end{cases}$$

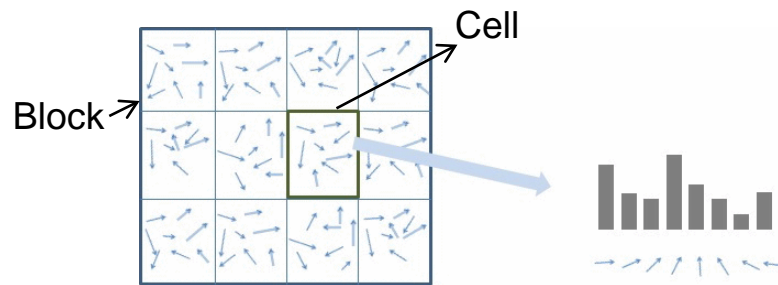O Integral image features (as for SURF):



Binary feature $h_i(x)$ defined by selecting a filter and a threshold in a window around pixel x.

Efficient computation through integral images.
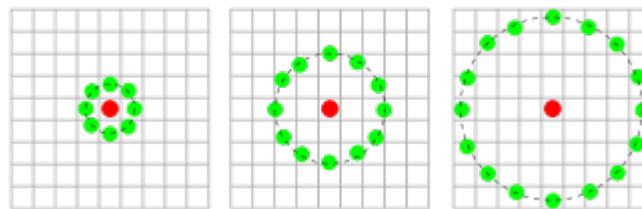
# Image patch *feature vectors*, to feed ML models

O Histogram of Oriented Gradients



Cell

Block

! normalization on blocks of cells.

Histograms of Oriented Gradients for Human Detection, Dalal and Triggs, CVPR2005.

O Local binary patterns:



Histogram of eight-bit vector (256 bins),
each vector being associated to
one pixel in the histogram cell/bock.

Dynamic texture recognition using local binary patterns with an application to facial
expressions, Zhao, Guoying, and Matti Pietikainen.  IEEE TPAMI, 2007.

# Outline

O **Introduction**
- Terminology: AI / ML / DL;
- AI4Vision:
  - From pixels to semantic;
  - Two paradigms: (AI or ML) vs. DL.
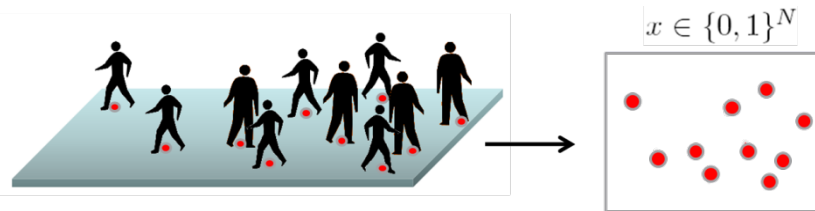
O **Hand-crafted image features**
- Pixel- or patch-based features;
- **Interpreting the features:**
  - **AI:        3D Geometry;**
  - AI:        Graph-based topology;
  - ML:        Random ferns, AdaBoost, Support Vector Machines.

O **DL in vision : Convolutional Neural Networks (CNNs)**
- Paradigm shift & revolution;
- CNNs basics;
- CNNs heads & backbones: network architecture examples;
- CNNs limitations.

# Detecting objects from foreground mask

- Assumption: ***Prior knowledge about the 3D shape of the object + camera calibration information***.

- Principle = Detect from object silhouette models :

    - Ground occupancy map: Measure the likelihood that an object is located in a ground plane position X.

    - The contribution of one view to the ground occupancy map in position X is obtained by integrating the foreground mask over the projected silhouette of the (expected) 3D object in X.

- This interpretation allows for
    - An efficient implementation, in which accumulation is based on integral images [1].
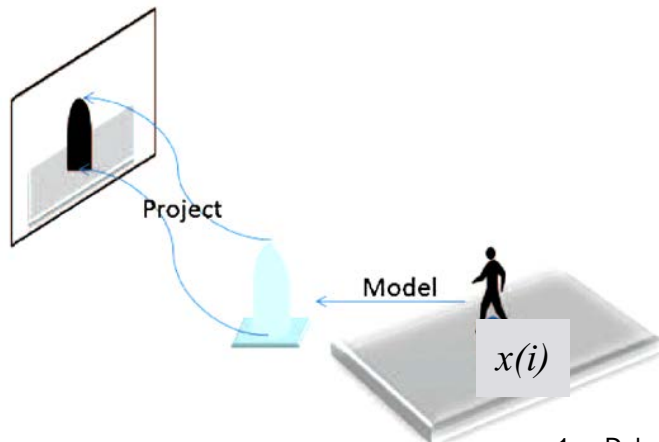    - An inverse problem formulation to reconstruct the sparse ground occupancy map [2].

$x \in \{0,1\}^N$

$$\underset{x \in \{0,1\}^N}{\arg\min} \; \|x\|_0 \;\; \text{s.t.} \;\; \|y - Q(Dx)\|_2^2 < \varepsilon$$

$$\underset{x \in \{0,1\}^N}{\arg\min} \; \|y - Q(Dx)\|_2^2 \;\; \text{s.t.} \;\; \|x\|_0 < \varepsilon_p$$

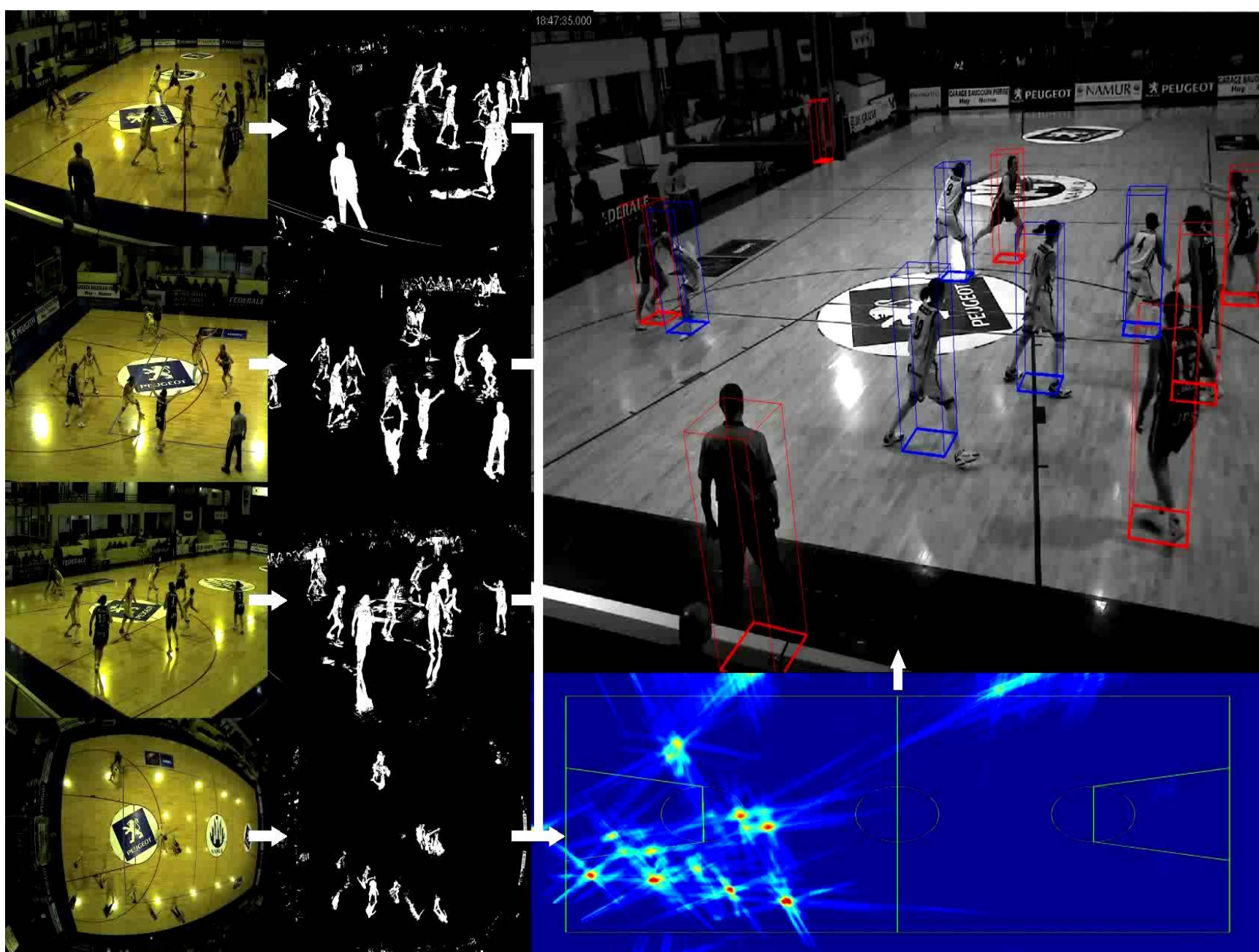$y$ = foreground mask;
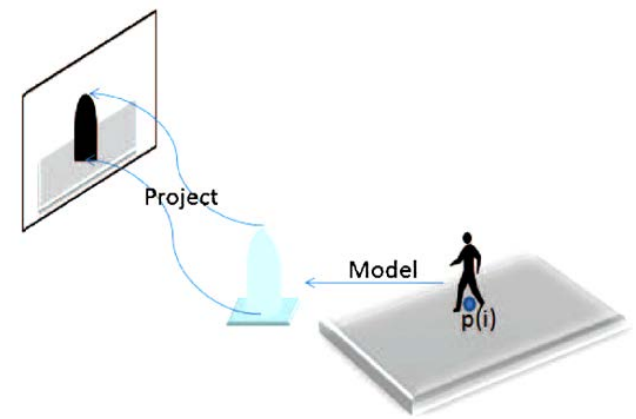$D$ = silhouette projection operator;
$Q$ = quantization (because silhouettes are not additive)

Project

Model

$x(i)$

1. Delannay and al., Detection and recognition of sports (wo)men from multiple views, ICDSC 2009.
2. Alahi and al., Sparsity driven people localization with a heterogeneous Network of Cameras, JMIV, 2009.
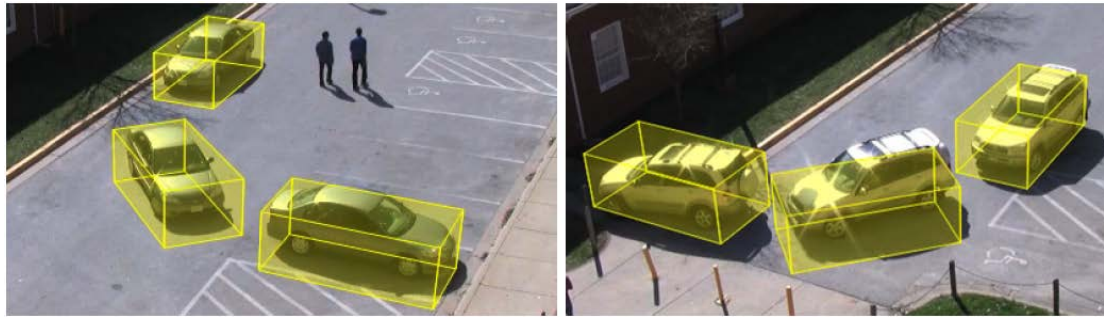3. I. Matthews, P. Carr, Y. Sheikh, Monocular Object Detection Using 3D Geometric Primitives, ECCV 2012.

**VIDEO**

## Mature & efficient but…

… requires **multiple models** for non-isotropic 3D objects;



**… unreliable in single-view** case.

Ground occupancy mask is ambiguous.

# Outline

O **Introduction**
- Terminology: AI / ML / DL;
- AI4Vision:
  - From pixels to semantic;
  - Two paradigms: (AI or ML) vs. DL.

O **Hand-crafted image features**
- Pixel- or patch-based features;
- **Interpreting the features:**
  - AI:        3D Geometry;
  - **AI:        Graph-based topology;**
  - ML:        Random ferns, AdaBoost, Support Vector Machines.

O **DL in vision : Convolutional Neural Networks (CNNs)**
- Paradigm shift & revolution;
- CNNs basics;
- CNNs heads & backbones: network architecture examples;
- CNNs limitations.

O  Mathematical morphology and Watershed (see 'Segmentation' lecture)

e.g. to clean a foreground mask          e.g. to segment a region of connected pixels with similar color.

O  Energy minimization via graph-cuts [Boykov, ECCV 2006]

Image segmentation can be regarded as computing a min-cut (or a max flow) in a graph.
The nodes of our graphs represent image pixels. S and T represent 'object' and 'background' labels (extendable to N labels).
Each pixel is connected to each terminal node (with t-links). Neighbouring pixels are interconnected by edges (n-links).

A s-t cut is a subset of edges such that the terminals S ad T become completely separated in the induced graph.
A minimum s-t cut minimizes the sum of the cost associated to severed edges.
This is equivalent to finding an object/background pixel assignment $A = (A_1, \ldots, A_p, \ldots, A_{|\mathcal{P}|})$ at minimizes :

$$E(A) = \lambda \cdot R(A) + B(A)$$

with
$$R(A) = \sum_{p \in \mathcal{P}} R_p(A_p) \quad \text{(regional term)}$$

$$B(A) = \sum_{\{p,q\} \in \mathcal{N}} B_{p,q} \cdot \delta_{A_p \neq A_q} \quad \text{(boundary term)}$$
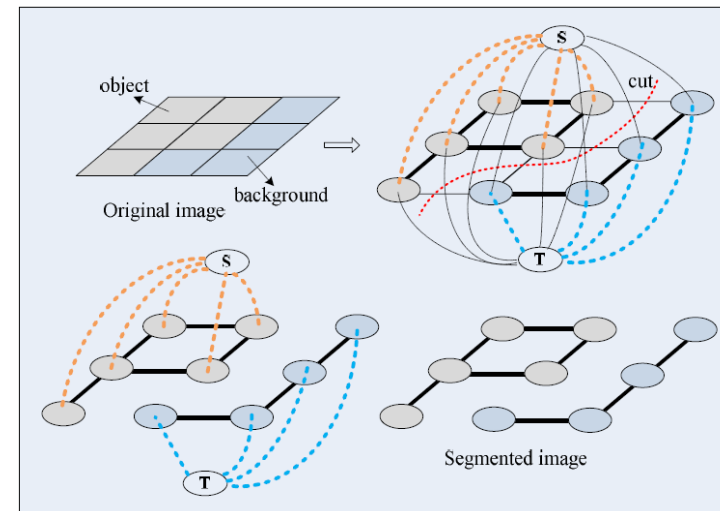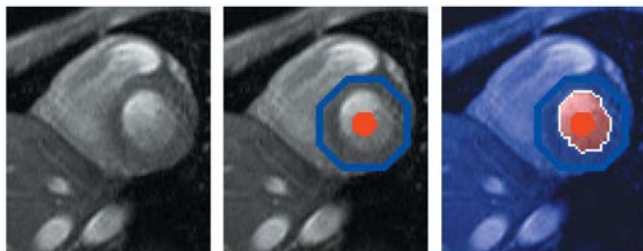
$R_p(.)$ measures how the pixel $p$ matches the object/background histogram.
$B_{p,q}$ coresponds to a penalty in case of label discontinuity between $p$ and $q$.

For example :  $R_p(\text{"obj"}) = -\ln \Pr(I_p|\text{"obj"})$   and   $B_{p,q} \propto \exp\left(-\frac{(I_p - I_q)^2}{2\sigma^2}\right)$

! NP-hard problem → approximation via iterative expansion & swap moves
See Y. Boykov, et al., Fast Approximate Energy Minimization via Graph Cuts. IEEE TPAMI, 2001

object

Original image

background

S

cut

T

S

T

Segmented image

(a) Original image      (b) Initialization      (c) Segmentation

# Outline

O  Introduction
- Terminology: AI / ML / DL;
- AI4Vision:
  - From pixels to semantic;
  - Two paradigms: (AI or ML) vs. DL.

O  Hand-crafted image features
- Pixel- or patch-based features;
- **Interpreting the features:**
  - AI:         3D Geometry;
  - AI:         Graph-based topology;
  - **ML:         Random ferns, AdaBoost, Support Vector Machines.**

O  DL in vision : Convolutional Neural Networks (CNNs)
- Paradigm shift & revolution;
- CNNs basics;
- CNNs heads & backbones: network architecture examples;
- CNNs limitations.

# Machine Learning in Vision: outline

O  Combining weak binary classifiers (= binary tests):

- Random ferns (Note: RFs share pros/cons with random trees)
- AdaBoost

O  Discriminating feature vectors :

- SVM ;
- Linear Discriminant Analysis. *Not covered this year.*

## Combining weak classifiers with **random ferns**.

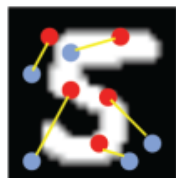**Goal:** associate a class label to a patch.

**Bayesian approach**, with uniform prior:

$$p(c|o) = \frac{p(o|c) \cdot p(c)}{p(o)} \qquad \Rightarrow \qquad p(c|o) \sim p(o|c)$$

Estimated from training samples

$$\hat{c} = \underset{c \in \{0,1\}}{\arg\max}\, P(C = c|b_1, ..., b_N) \;\; = \underset{c \in \{0,1\}}{\arg\max}\, P(b_1, ..., b_N|C = c)$$



binary code

**Problem:** observation space dimensionality.

**Semi-Naive Bayesian approach**: group the *N* binary tests into *M* sets, named ferns, of size *N/M,* and assume independence between ferns:
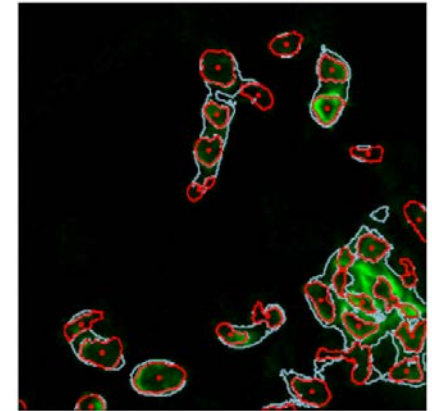
$$P(b_1, ..., b_N|C = c) = \prod_{k=1}^{M} P(F_k|C = c)$$

**Example : Binary test usage for cell segmentation.** https://arxiv.org/abs/1602.05439

Assign a label (interior/boundary/exterior) to each pixel based on binary tests in a small squared neighborhood.
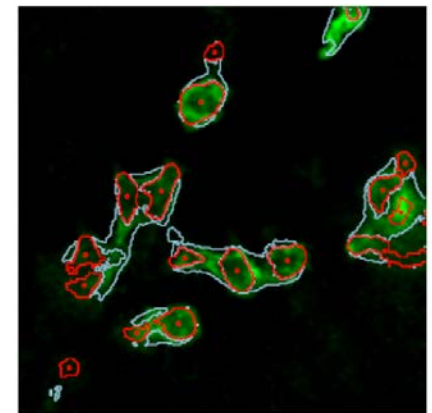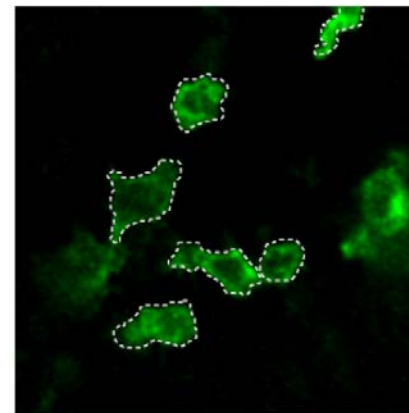
Estimated probability (log scale, zero mean) of being in the background class for 3 types of pixels
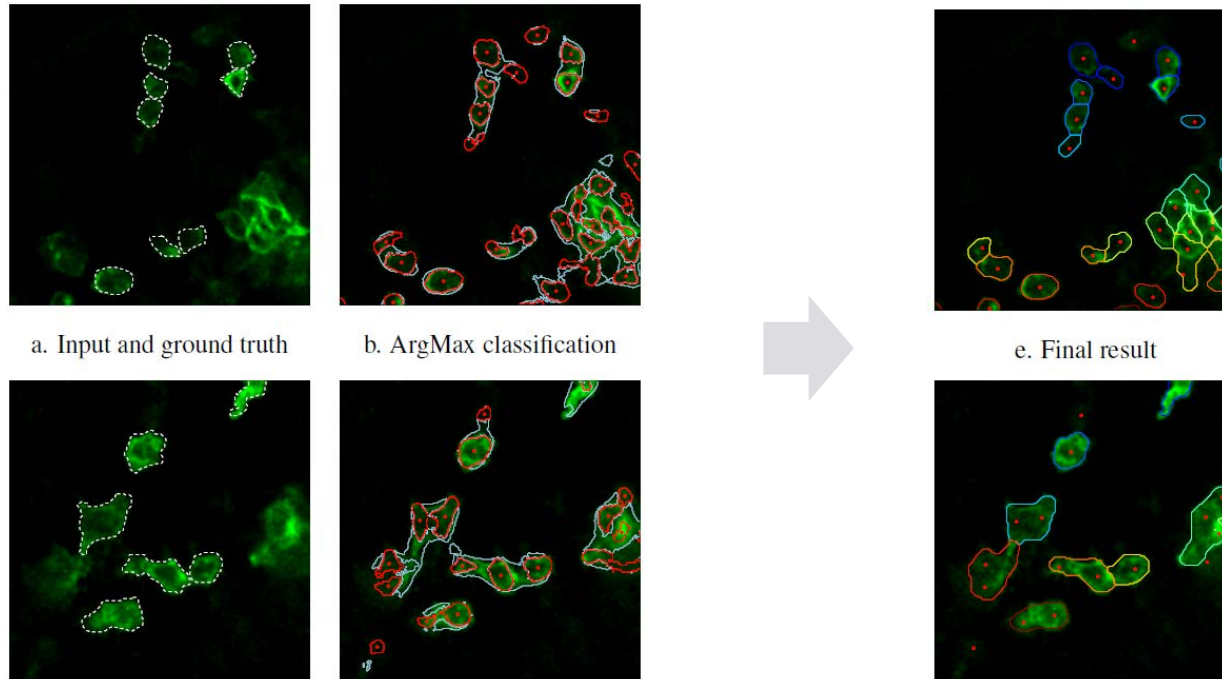


a. Input and ground truth

b. ArgMax classification

Note: Graph cuts can be used to improve argmax segmentation. [Browet, IEEE ICIP 2016]



a. Input and ground truth    b. ArgMax classification        e. Final result

Seeds = centers of argmax interior regions.

One label per seed, and graph-cut to define label assignment $f_p$ to pixel $p$, so as to minimize an energy function:

$D_p(f_p)$ is small when seed $f_p$ is close to $p$ ('inside' path).

Penalizes transitions occuring in non-boundary pixels.

Penalizes the use of many labels.

$$E\left(f\right) = \sum_{p\in\mathcal{P}} D_p\left(f_p\right) + \sum_{(p,q)\in\mathcal{N}} W\left(p,q\right)\left(1 - \delta(f_p, f_q)\right) + \sum_{l\in\mathcal{L}} h_l\left(f\right)$$

Goal: Define $H^T(x) = \sum_{t=1}^{T} \alpha_t \cdot h_t(x)$ to maximize $\sum_i y_i \cdot \text{sign}(H^T(x_i))$

when summing over $N$ training samples $(x_i, y_i)$, with $y_i = 1$ or $-1$.

Trick: Find $H^T(x)$ to minimize $E_T = \sum_i e^{-y_i \cdot H^T(x_i)}$

Exponential in $E_T$ implies that the energy change caused by an additional term in the sum of weak classifiers is larger for wrongly classified samples than for correctly classified ones. At each iteration, this favors the selection of a weak classifier $h_m$ that classifies correctly the samples that are wrongly classified by $H^{m-1}$.

O  Proceed iteratively, knowing $H^{m-1}(x)$,

select $h_m(x)$ and $\alpha_m$ such that $H^m(x) = H^{m-1}(x) + \alpha_m \cdot h_m(x)$ minimizes $E_m$.

O
$$E_m = \sum_i e^{-y_i \cdot H^{m-1}(x_i)} \cdot e^{-y_i \cdot \alpha_m \cdot h_m(x_i)}$$

To minimize since $e^{\alpha_m} > e^{-\alpha_m}$

$$= \sum_{h_m(x_i)=y_i} w^m(i) \cdot e^{-\alpha_m} + \sum_{h_m(x_i) \neq y_i} w^m(i) \, e^{\alpha_m} \quad \text{with} \quad w^m(i) \equiv e^{-y_i \cdot H^{m-1}(x_i)}$$

O  $h_m = \underset{h \in \mathrm{H}}{\text{argmin}} \sum_i w^m(i)[h(x_i) \neq y_i]$ and $\dfrac{\partial E_m}{\partial \alpha_m} = 0 \rightarrow \alpha_m = \dfrac{1}{2} \cdot \ln\left( \dfrac{\sum_{h_m(x_i)=y_i} w^m(i)}{\sum_{h_m(x_i) \neq y_i} w^m(i)} \right) > 0$
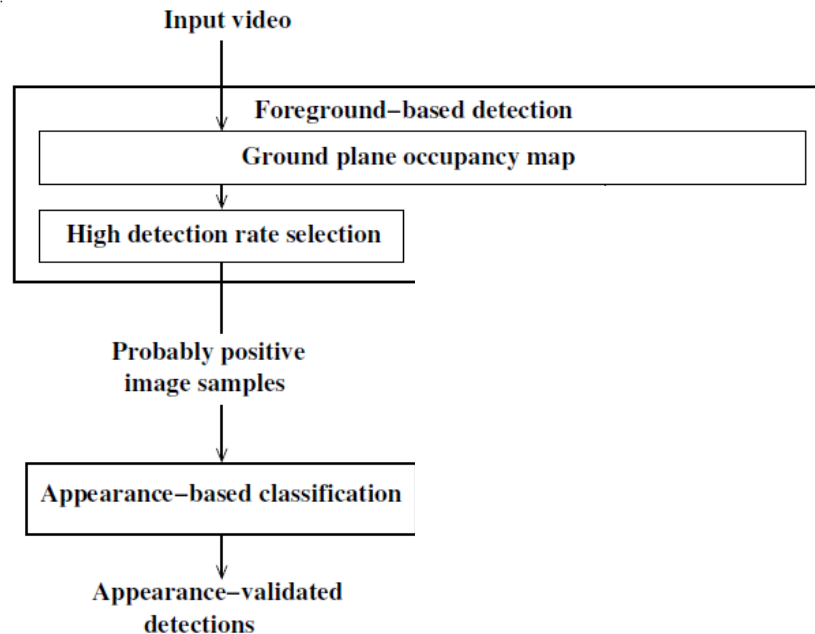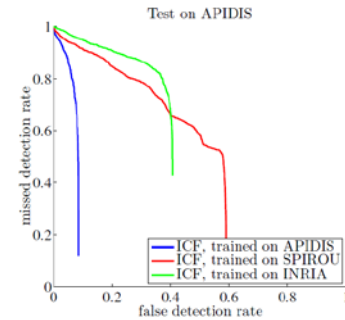
Random subset if H is huge.

Minimizing the sum requires to select a $h_m$ that classifies correctly samples that are wrongly classified by $H^{m-1}$, because they have a large $w^m(i)$.

# Ensemble of classifiers: Adaboost or Random Ferns ?


Test on APIDIS

O **Key concern: label corruption.**

O Example:

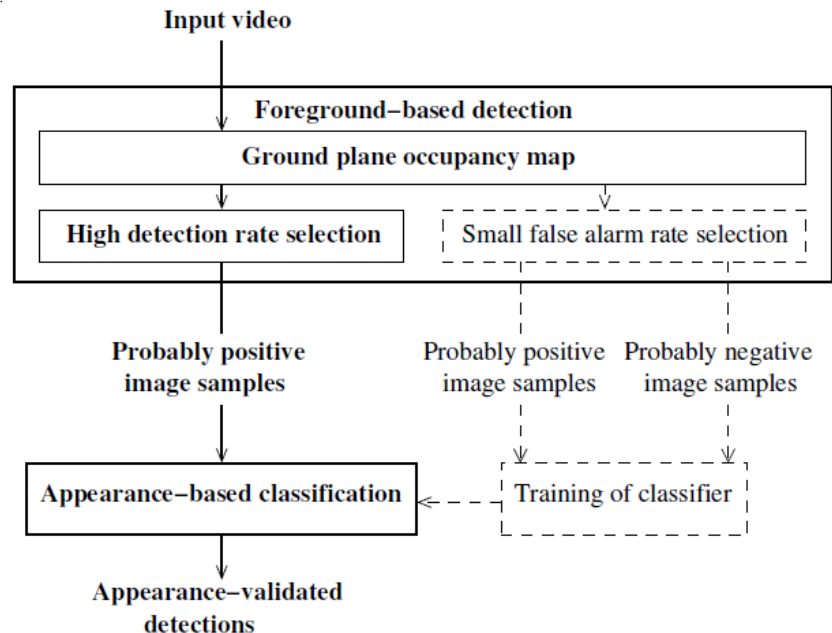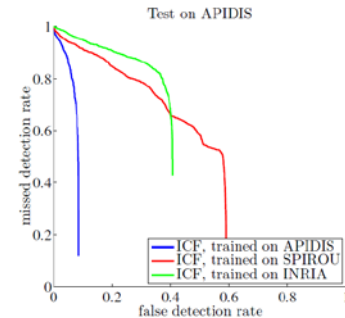Foreground detector + *SCENE-SPECIFIC* appearance-based classifier.



**Input video**

**Foreground−based detection**

**Ground plane occupancy map**

**High detection rate selection**

**Probably positive image samples**

**Appearance−based classification**

**Appearance−validated detections**

Scene-specific classifier for effective and efficient team sport players
detection from a single calibrated camera, CVIU2016.

# Ensemble of classifiers: Adaboost or Random Ferns ?
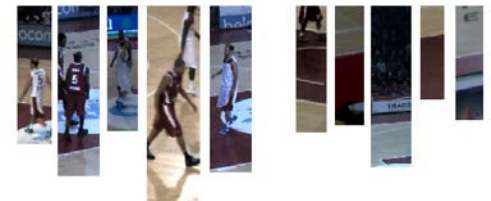

Test on APIDIS

O **Key concern: label corruption.**

O Real-life example:

Foreground detector + *SCENE-SPECIFIC* appearance-based classifier.



To adapt the classifier to the game/scene at hand, the training is based on samples labelled *with high confidence* by the foreground detector.

Scene-specific classifier for effective and efficient team sport players detection from a single calibrated camera, CVIU2016.
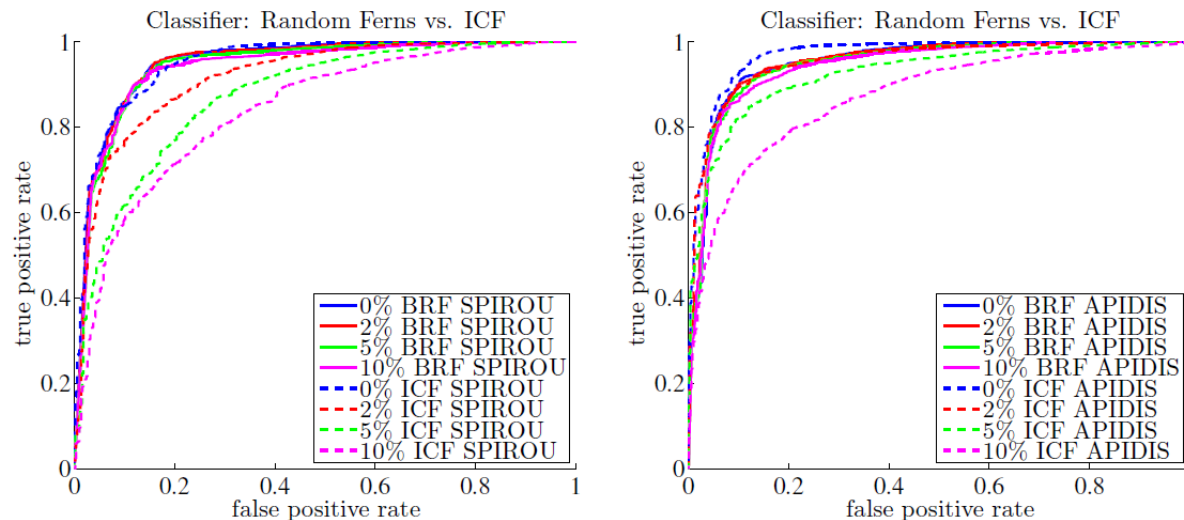
ICF = Integral (color) Channel Features, combined with AdaBoost [BMVC2009; CVPR2012; TPAMI2014].

BRF = Block Random Ferns:

Bounding box investigated by the classifier is split in a grid of square blocks;
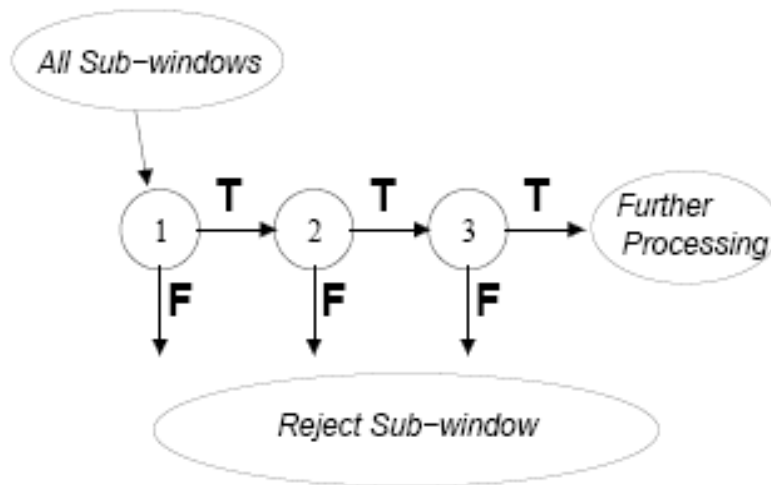
All the binary tests (= pixels value comparison) in a given ferns are related to the same block.

[CVIU2016] : Scene-specific classifier for effective and efficient team sport players detection from a single calibrated camera.

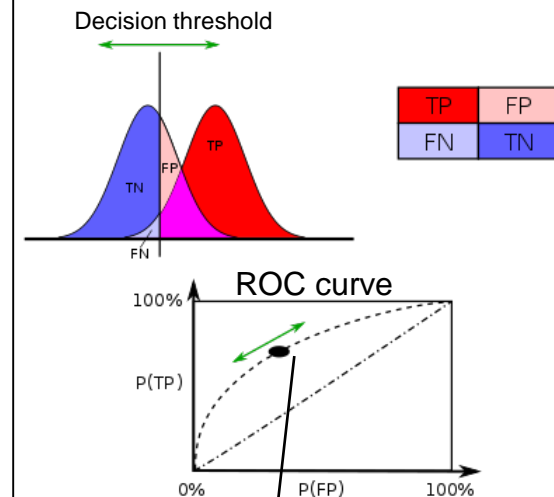# Ensemble of classifiers & **computational complexity**.

**Cascade of classifiers** allows for very fast testing:

How to select operating points for the binary classifiers in the cascade?



Schematic depiction of a the detection cascade. A series of classifiers are applied to every sub-window. The initial classifier eliminates a large number of negative examples with very little processing. Subsequent layers eliminate additional negatives but require additional computation. After several stages of processing the number of sub-windows have been reduced radically. Further processing can take any form such as additional stages of the cascade (as in our detection system) or an alternative detection system.

The operating point moves along the ROC curve when changing the classifier decision threshold.

Note: ROC = receiver operating characteristic.

Note: imbalance of training samples. 5000 vs. 350.000.000!

*Figure from 'Learning object waveforms: Face detection' Viola & Jones, 2001, and https://en.wikipedia.org/wiki/Receiver_operating_characteristic*

❑ Histogrammes d'orientations de gradients (HoG)
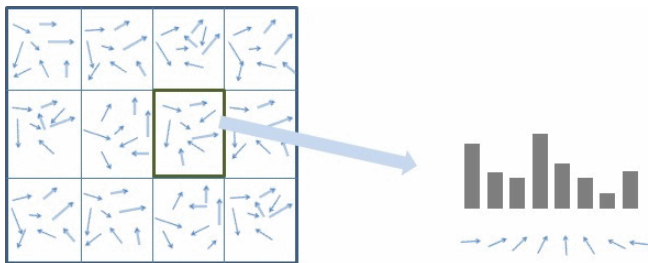


Fast and Accurate Digit Classification, Maji and Malik, 2009.



Image window
=
object of a given class when gradient orientations in window blocks correspond to the class model.

# Support Vector Machine
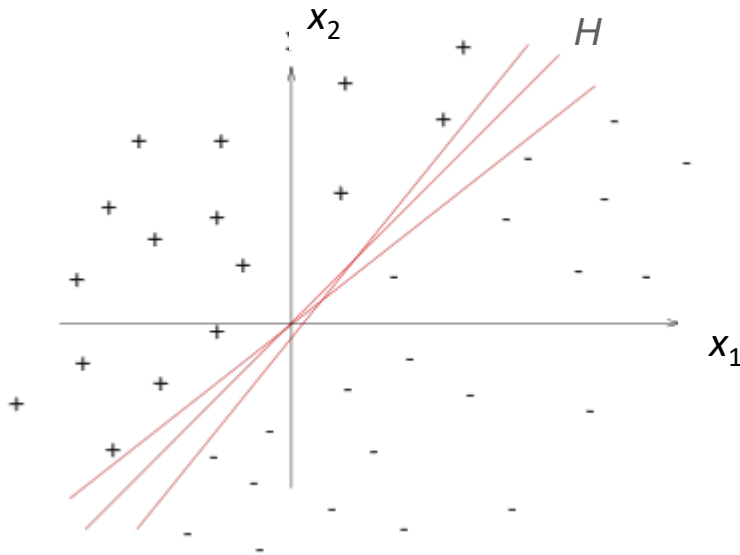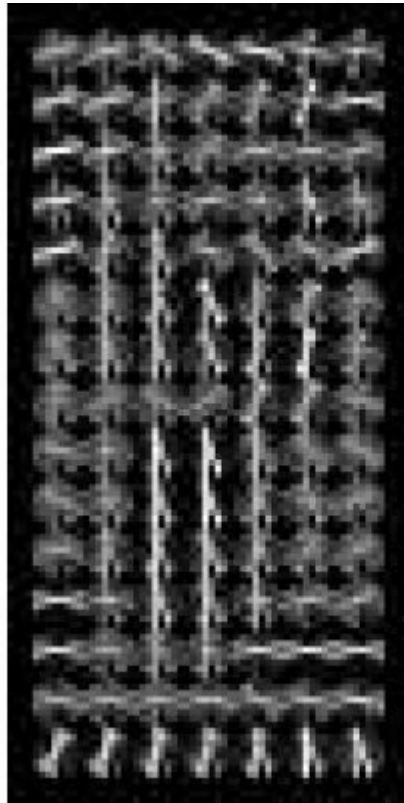
O  Linear SVM:

$$decision = sign(\mathbf{w}^T \cdot \mathbf{x} + b)$$



The SVM optimization problem
is formulated so as to find the
hyperplane that best separate
the training data
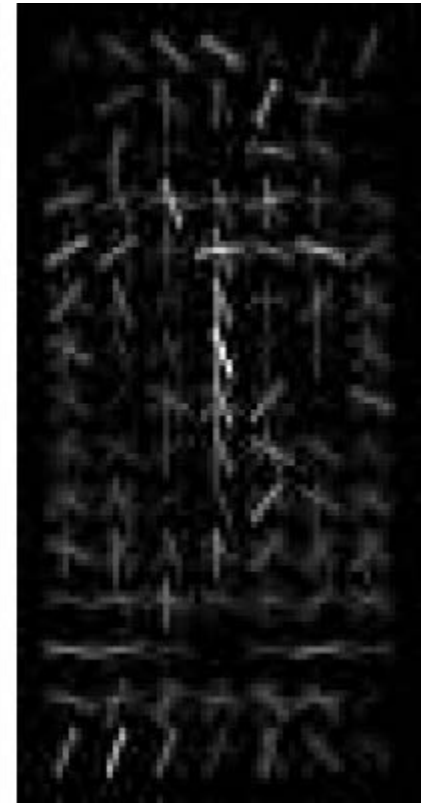into their two classes.

# O HOG + Linear SVM :



HoG

HoG components corresponding to positive SVM weights

HoG components corresponding to negative SVM weights