

Face, Human Detection

Disclaimer: Many of the slides used here are obtained from online resources (including many open lecture materials) without specific acknowledgements. They are used here for the sole purpose of classroom learning. All copyrights belong to the original authors or publishers. You should not copy it, redistribute it, put it online, or use it for any purposes other than help you learning ENGN4528/6528.

Announcement

- Lab Sessions will happen at this and next week for Lab Assignment 2. Please attend the lab sessions and ask questions.
- Final exam timetable is to be confirmed. The initially suggested time by the exam office is on 9:00 am Canberra time, June 3rd, which might be too early for students overseas. We have requested changes.

Outline

- Review
- Viola-Jones Face detection- continue
- HOG features
- Bag of words for object detection

Review: Eigenfaces (PCA on face images)

1. Compute covariance matrix of face images
2. Compute the principal components (“eigenfaces”)
 - K eigenvectors with top K largest eigenvalues
3. Represent all face images in the dataset as linear combinations of eigenfaces
 - Perform nearest neighbor on these coefficients

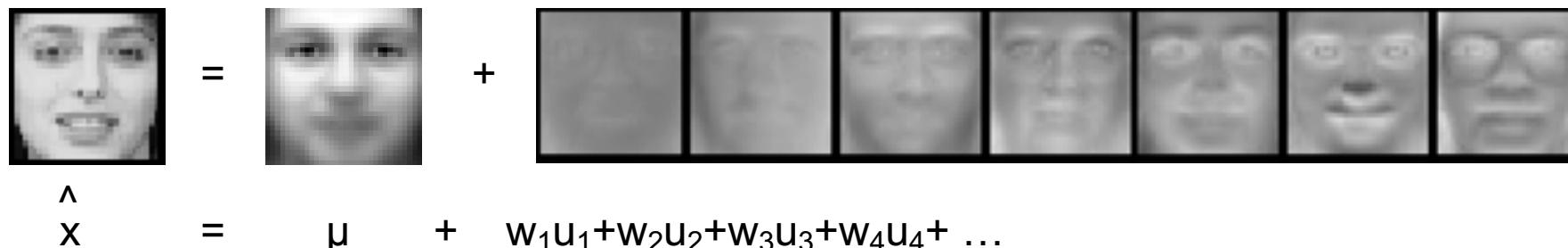
Review: Eigen Face, Representation and reconstruction

- Face \mathbf{x} in “face space” coordinates:

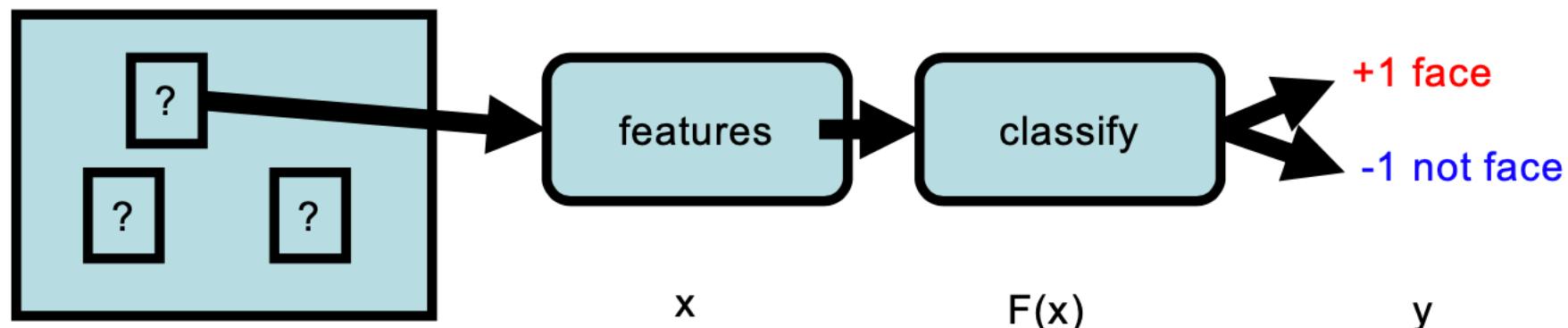


$$\begin{aligned}\mathbf{x} &\rightarrow [\mathbf{u}_1^T(\mathbf{x} - \mu), \dots, \mathbf{u}_k^T(\mathbf{x} - \mu)] \\ &= w_1, \dots, w_k\end{aligned}$$

- Reconstruction:

$$\begin{aligned}\hat{\mathbf{x}} &= \mu + \mathbf{w}_1\mathbf{u}_1 + \mathbf{w}_2\mathbf{u}_2 + \mathbf{w}_3\mathbf{u}_3 + \mathbf{w}_4\mathbf{u}_4 + \dots \\ &= \mu + w_1\mathbf{u}_1 + w_2\mathbf{u}_2 + w_3\mathbf{u}_3 + w_4\mathbf{u}_4 + \dots\end{aligned}$$


Face detection -- basic scheme



- We slide a window over the image
- Extract features for each window
- Classify each window into face/non-face

Viola-Jones face detection

Paul A. Viola and Michael J. Jones

Intl. J. Computer Vision

57(2), 137–154, 2004

Some slides (*slides adapted from Bill Freeman, MIT 6.869, April 2005*)

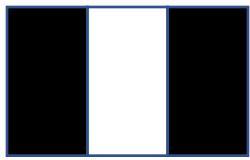
Characteristics of the algorithm

- Feature set (...is huge about 16M features)
 - Efficient feature selection using AdaBoost
 - New image representation: Integral Image
 - Cascaded Classifier for rapid detection
- Fastest known face detector for gray scale images (circa 2001, Marr Prize Winner).

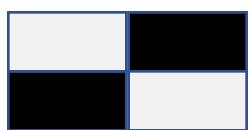
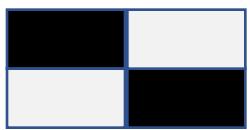
Viola-Jones face detection

- 3 key ideas
 - Use of Haar features
 - Integral Image
 - Cascade classifier
- For more information, read the original paper.
- Also, lots of videos on the web.
 - <https://www.youtube.com/watch?v=uEJ71VIUmMQ&vl=en-GB> (Google: Viola-Jones Computerphile)

Haar features

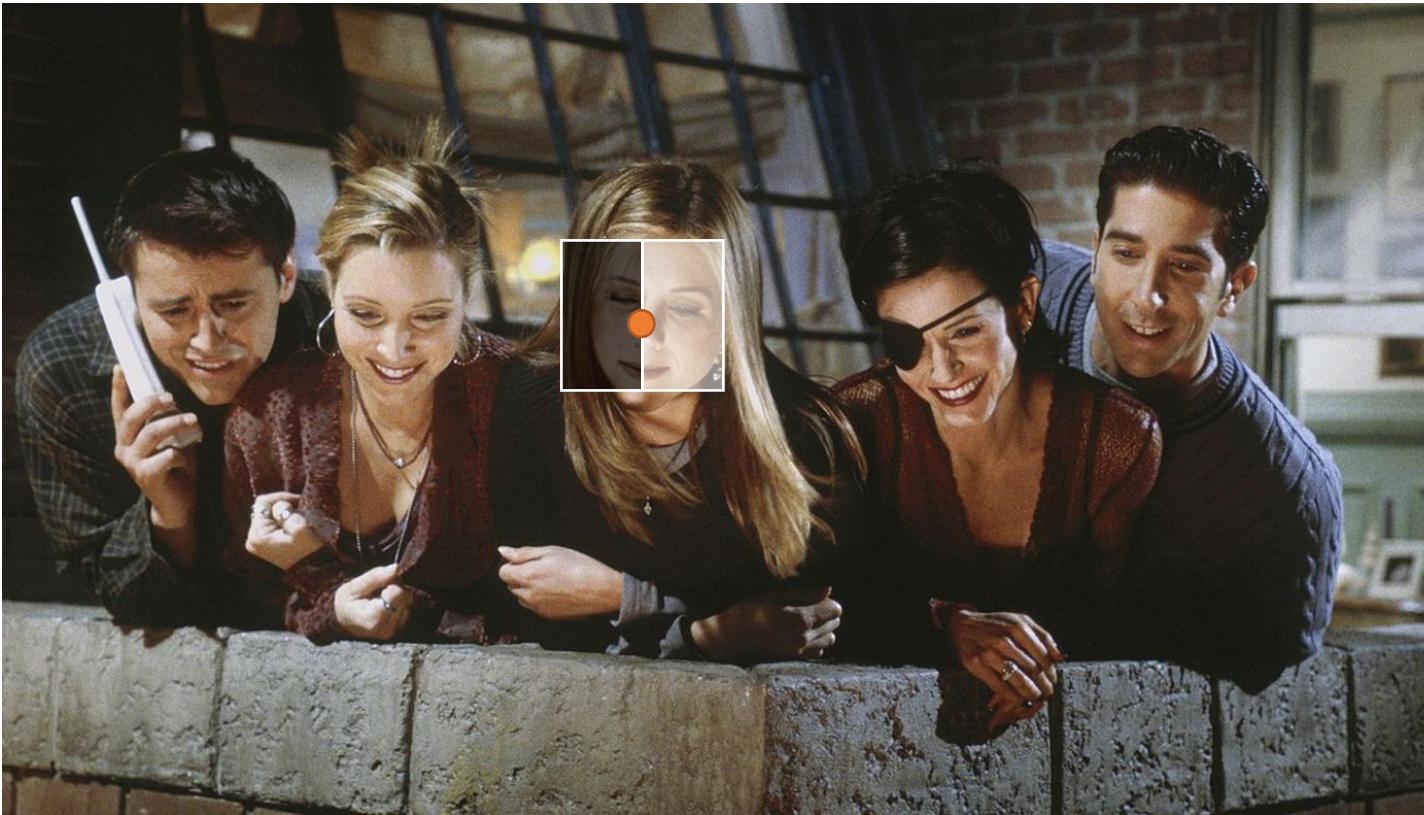


Black = add
White = subtract.



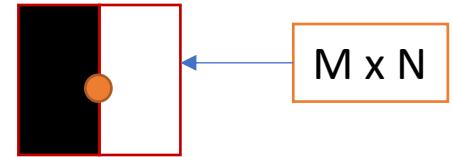
Lots of kernels at
different scales and
shapes.

Haar Features



$$V_A[i, j] = \sum(\text{pixel intensities in black area}) - \sum(\text{pixel intensities in white area})$$

Haar Features



Computation Cost
= $(NxM-1)$ additions per pixel, per filter, per scale.

How could we compute the features Efficiently?

$$V_A[i, j] = \sum_{\text{black area}} (\text{pixel intensities}) - \sum_{\text{white area}} (\text{pixel intensities})$$

Integral Image

Integral Image

162	405	243	243	486
486	162	324	324	405
162	162	243	243	324
486	324	243	324	243
405	405	162	162	243

Image I

162	567	810	1053	1539
648	1215	1782	2349	3240
810	1539	2349	3159	4374
1296	2349	3402	4536	5994
1701	3159	4374	5670	7371

Integral Image

- A table that holds the sum of all pixel values to the left and top of a given pixel, inclusive.

Integral Image

162	405	243	243	486
486	162	324	324	405
162	162	243	243	324
486	324	243	324	243
405	405	162	162	243

Image I

162	567	810	1053	1539
648	1215	1782	2349	3240
810	1539	2349	3159	4374
1296	2349	3402	4536	5994
1701	3159	4374	5670	7371

Integral Image

- A table that holds the sum of all pixel values To the left and top of a given pixel, inclusive.

Integral Image

162	405	243	243	486
486	162	324	324	405
162	162	243	243	324
486	324	243	324	243
405	405	162	162	243

Image I

162	567	810	1053	1539
648	1215	1782	2349	3240
810	1539	2349	3159	4374
1296	2349	3402	4536	5994
1701	3159	4374	5670	7371

Integral Image

- A table that holds the sum of all pixel values To the left and top of a given pixel, inclusive.

Integral Image

- Fast computation of arbitrary rectangles using integral image.

162	405	243	243	486
486	162	324	324	405
162	162	243	243	324
486	324	243	324	243
405	405	162	162	243

Image I

162	567	810	1053	1539
648	1215	1782	2349	3240
810	1539	2349	3159	4374
1296	2349	3402	4536	5994
1701	3159	4374	5670	7371

Integral Image

Integral Image

- Fast computation of arbitrary rectangles using integral image.

162	405	243	243	486
486	162	324	324	405
162	162	243	243	324
486	324	243	324	243
405	405	162	162	243

Image I

162	567	810	1053	1539
648	1215	1782	2349	3240
810	1539	2349	3159	4374
1296	2349	3402	4536	5994
1701	3159	4374	5670	7371

Integral Image II

P

$$\text{Sum} = \text{II}_P \dots \\ = 7371$$

Integral Image

- Fast computation of arbitrary rectangles using integral image.

162	405	243	243	486
486	162	324	324	405
162	162	243	243	324
486	324	243	324	243
405	405	162	162	243

Image I

162	567	810	1053	1539
648	1215	1782	2349	3240
810	1539	2349	3159	4374
1296	2349	3402	4536	5994
1701	3159	4374	5670	7371

Integral Image II

$$\text{Sum} = II_P - II_Q \dots \\ = 7371 - 3240$$

Q

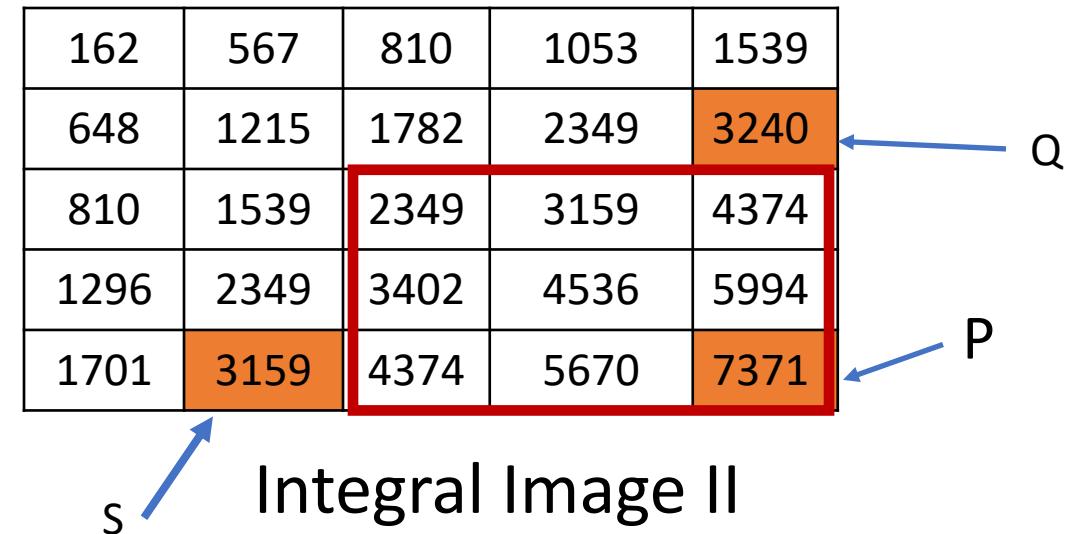
P

Integral Image

- Fast computation of arbitrary rectangles using integral image.

162	405	243	243	486
486	162	324	324	405
162	162	243	243	324
486	324	243	324	243
405	405	162	162	243

Image I



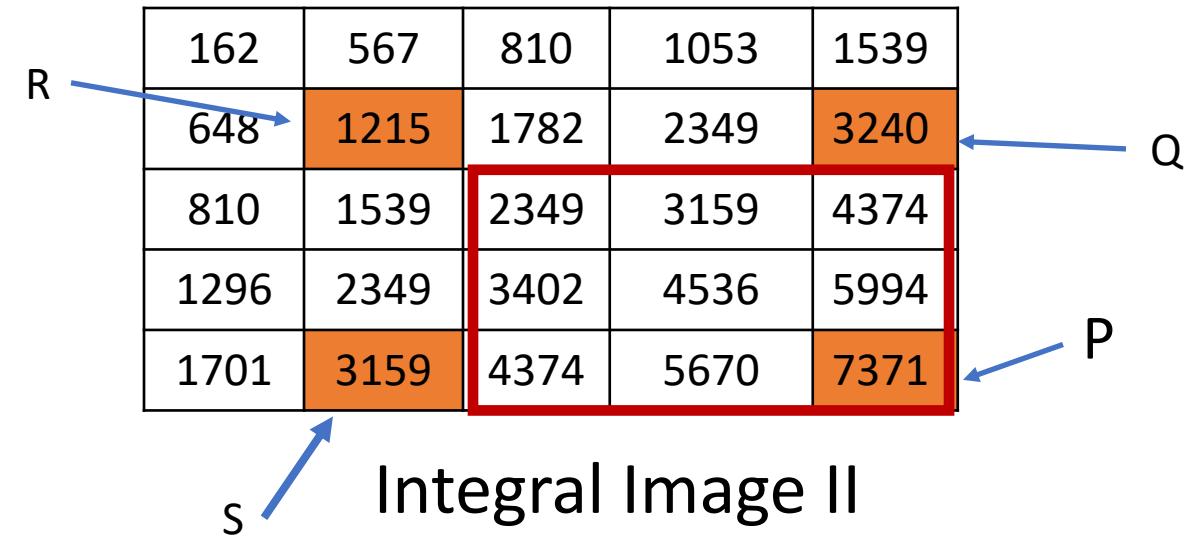
$$\begin{aligned} \text{Sum} &= II_P - II_Q - II_S \dots \\ &= 7371 - 3240 - 3159 \dots \end{aligned}$$

Integral Image

- Fast computation of arbitrary rectangles using integral image.

162	405	243	243	486
486	162	324	324	405
162	162	243	243	324
486	324	243	324	243
405	405	162	162	243

Image I



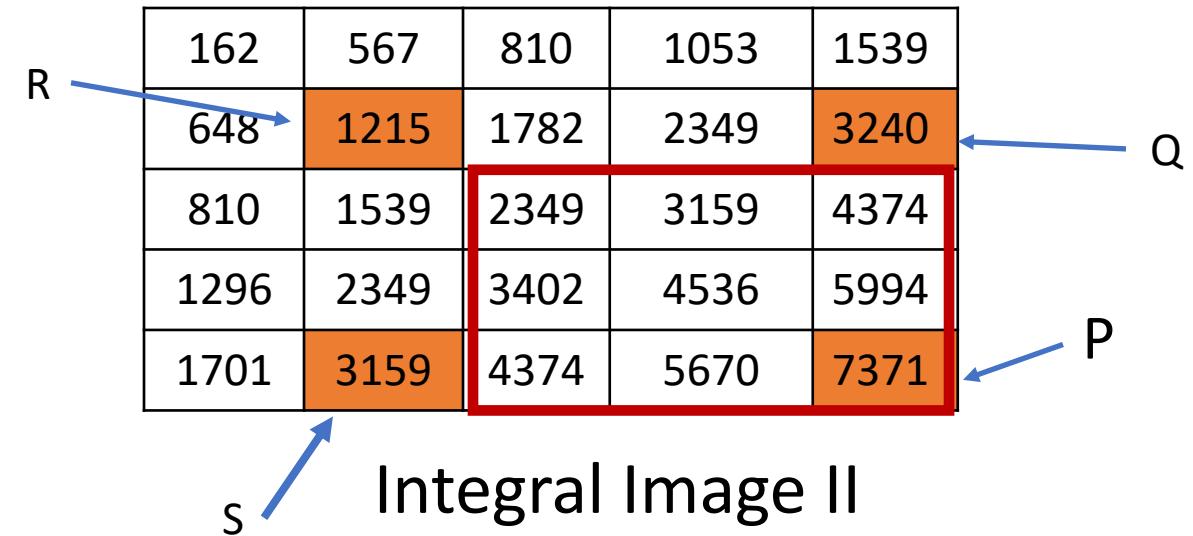
$$\begin{aligned}\text{Sum} &= II_P - II_Q - II_S + II_R \\ &= 7371 - 3240 - 3159 + 1215 \\ &= 2187\end{aligned}$$

Integral Image

- Fast computation of arbitrary rectangles using integral image.

162	405	243	243	486
486	162	324	324	405
162	162	243	243	324
486	324	243	324	243
405	405	162	162	243

Image I

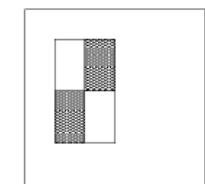
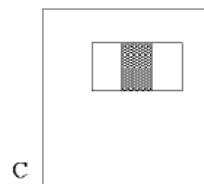
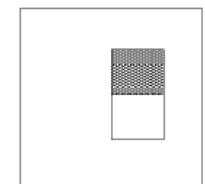
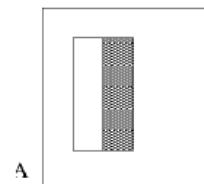


$$\begin{aligned}\text{Sum} &= \text{II}_P - \text{II}_Q - \text{II}_S + \text{II}_R \\ &= 7371 - 3240 - 3159 + 1215 \\ &= 2187\end{aligned}$$

Computational cost: only 3 additions which is Invariant to the size of the rectangle.

Image features

- “Rectangle filters”
- Differences between sums of pixels in adjacent rectangles



$$h_t(x) = \begin{cases} +1 & \text{if } f_t(x) > \theta_t \\ -1 & \text{otherwise} \end{cases}$$



Threshold determined by training

Cascade classifier

Slides adapted from Prof. Richard Hartley

Constructing the classifier

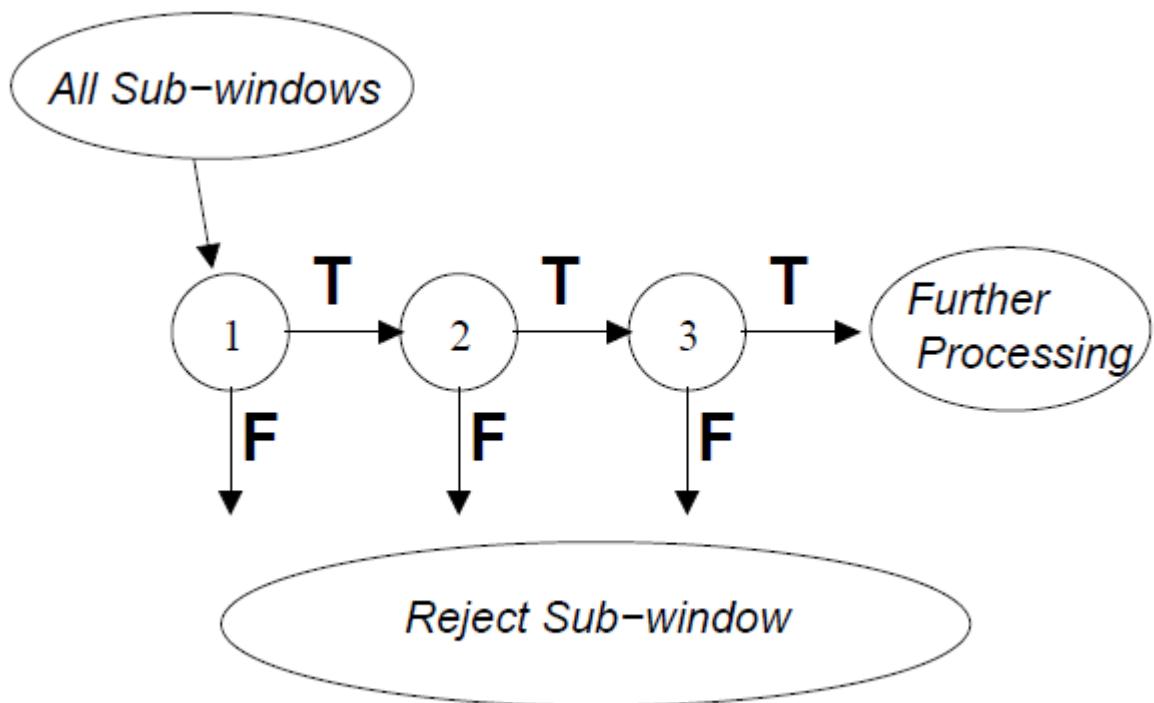
- Perceptron yields a sufficiently powerful weak classifier

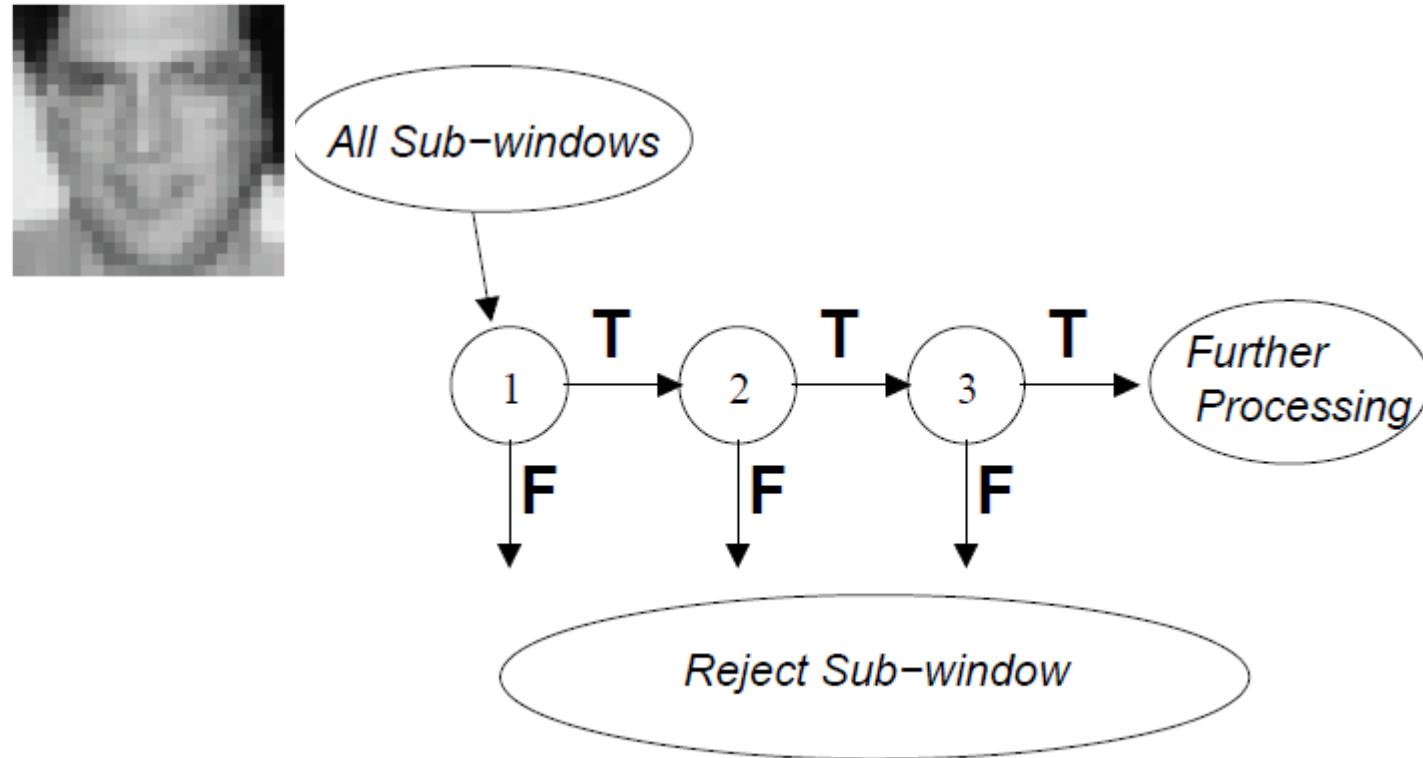
$$C(x) = \theta\left(\sum_i \alpha_i h_i(x) + b\right)$$

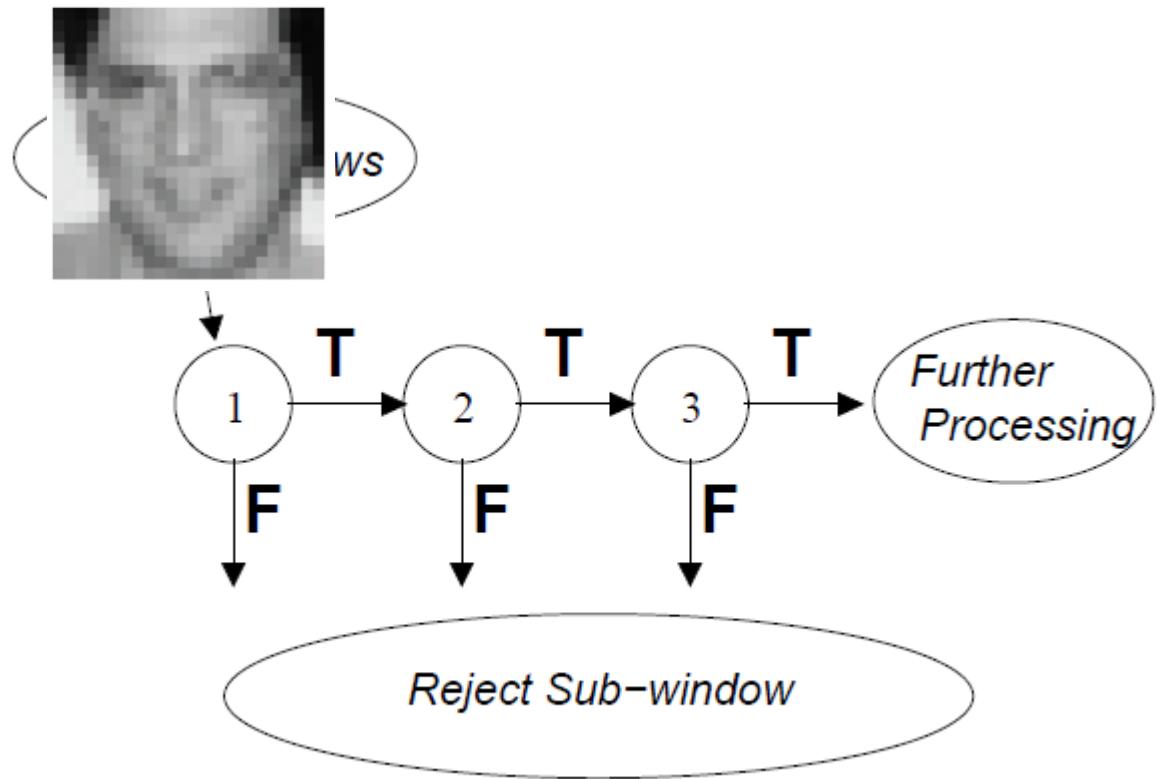
- Use AdaBoost to efficiently choose best features among the feature library.
- add a new $h_i(x)$ at each round
- each $h_i(x_k)$ is a “decision stump”

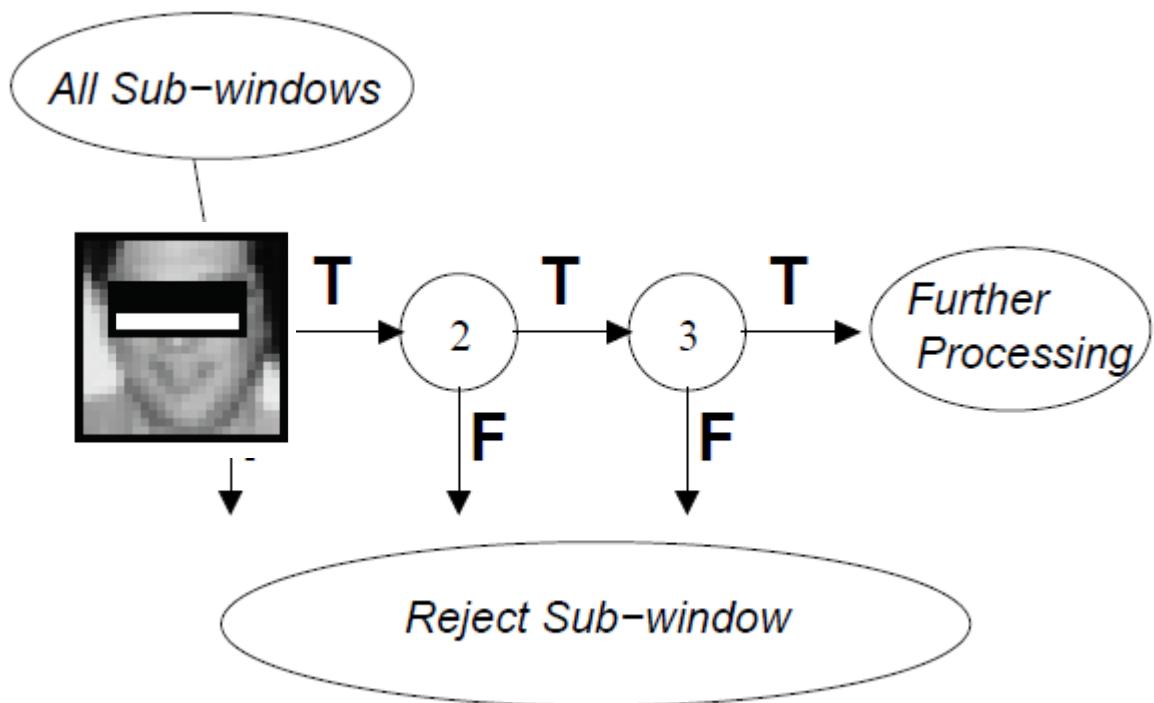
Constructing the classifier

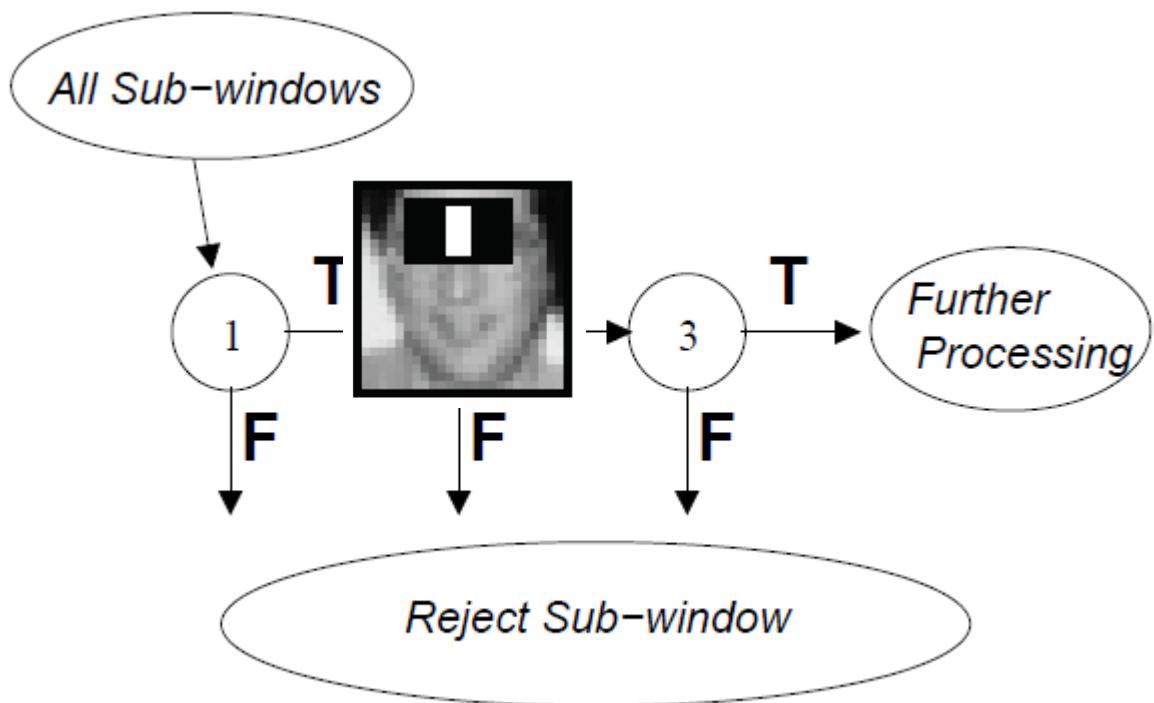
- For each round of boosting:
- Evaluate each rectangle filter on each example
- Sort examples by filter values
- Select best threshold for each filter (min error)
- Use sorting to quickly scan for optimal threshold
- Select best filter/threshold combination
- Weight is a simple function of error rate
- Reweight examples

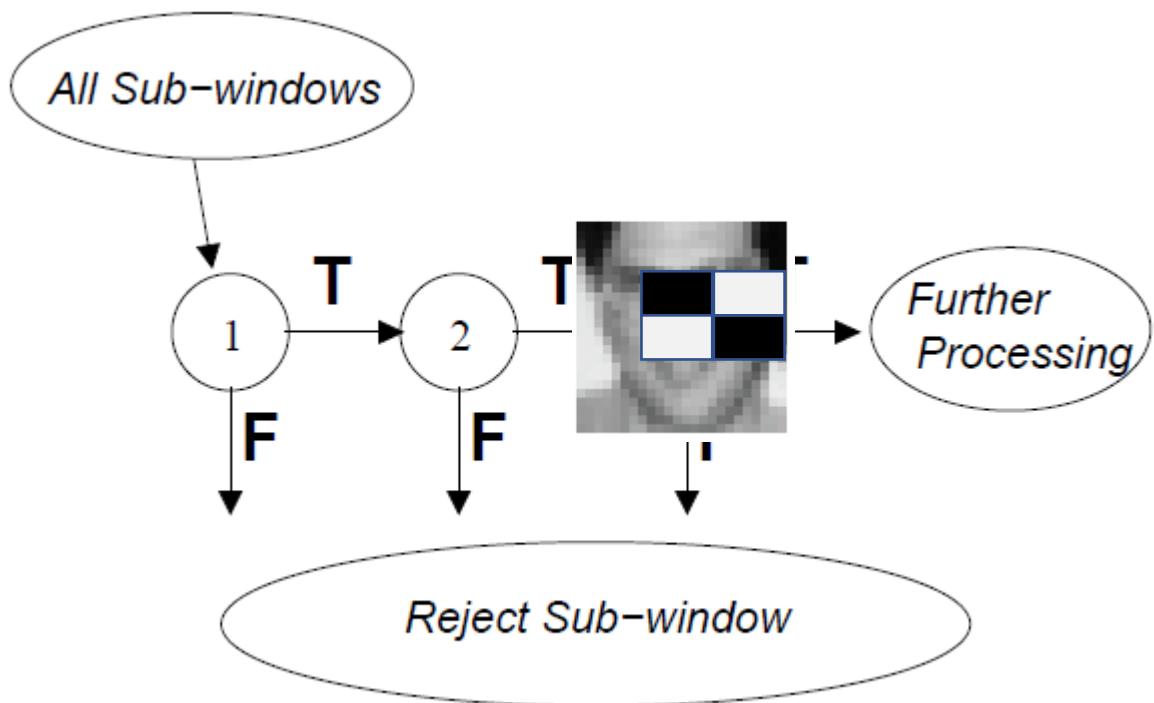


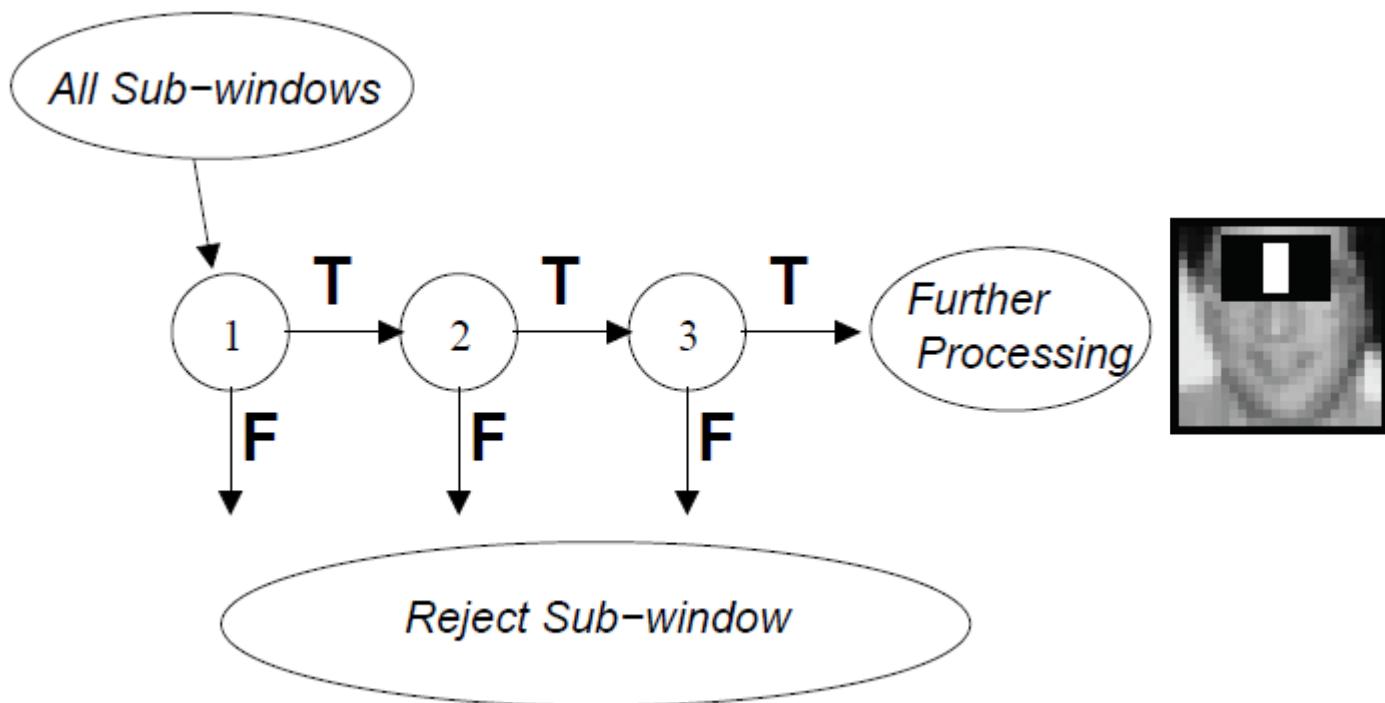


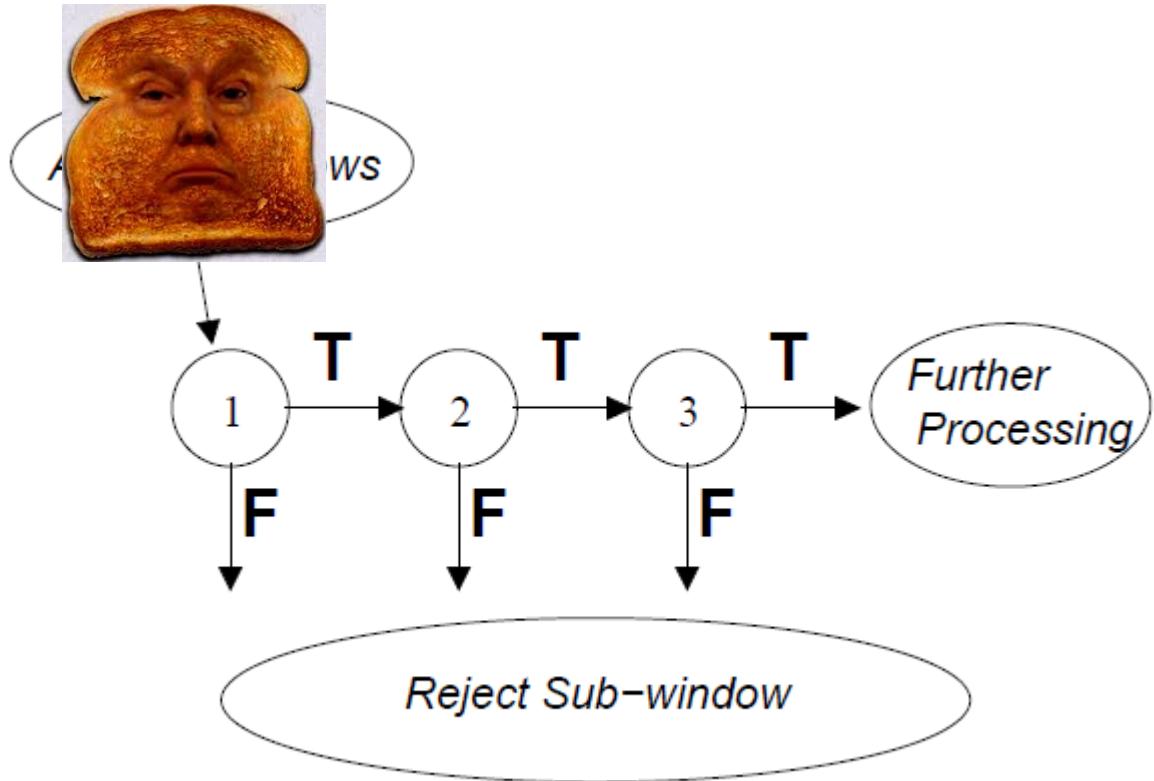


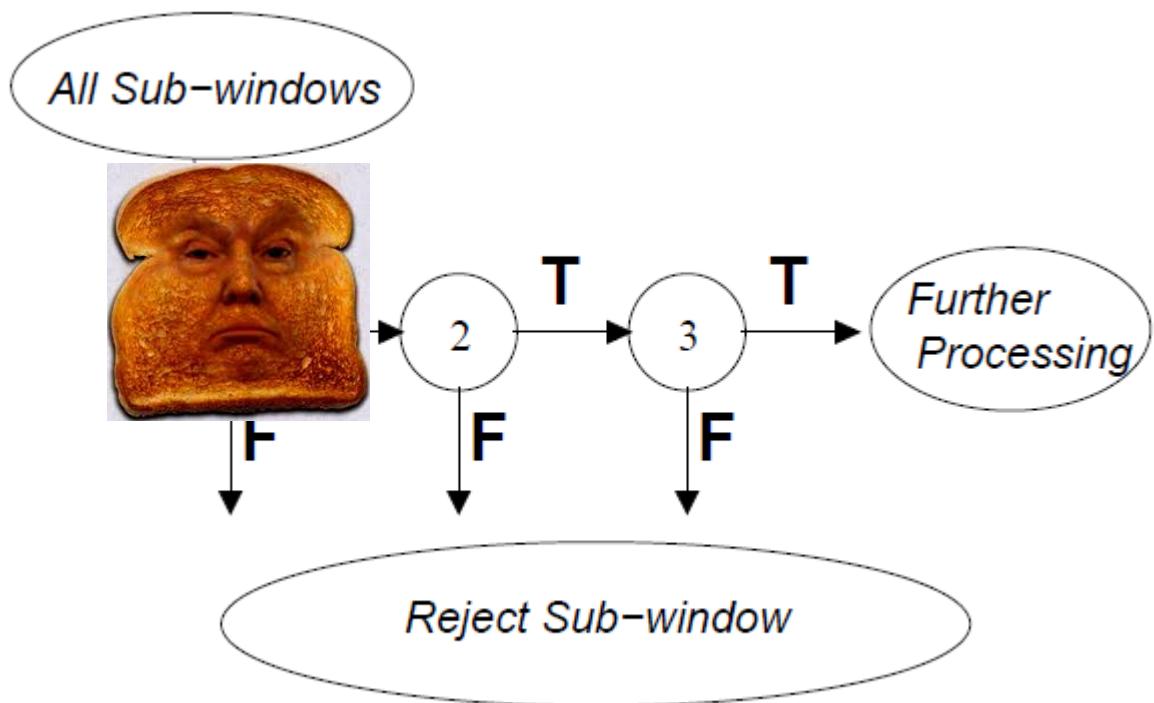


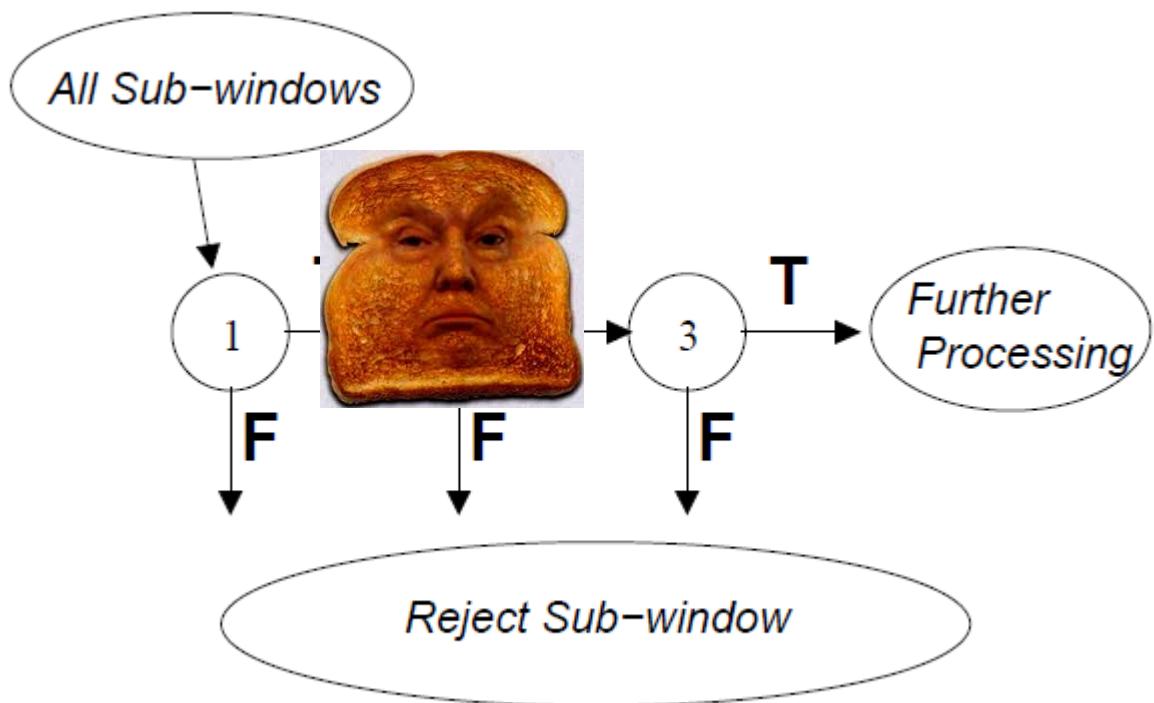


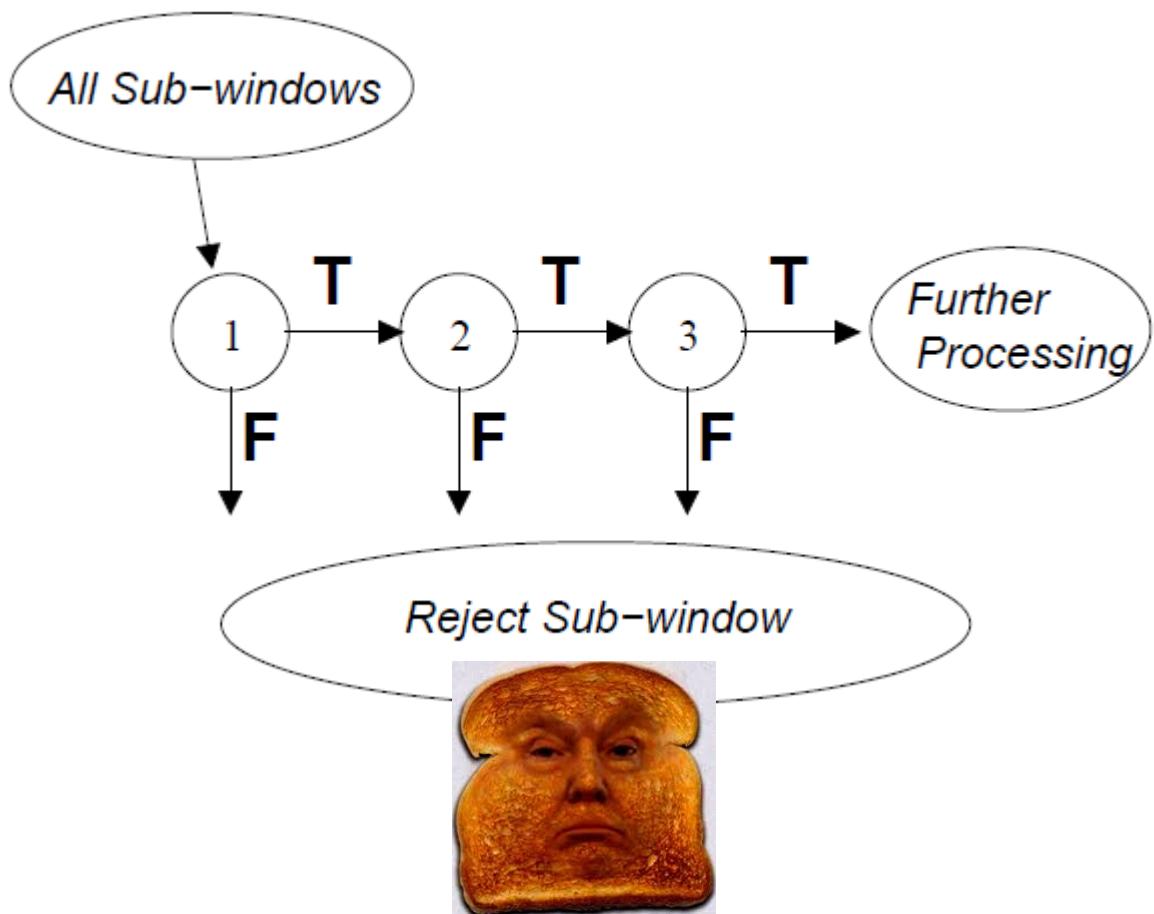




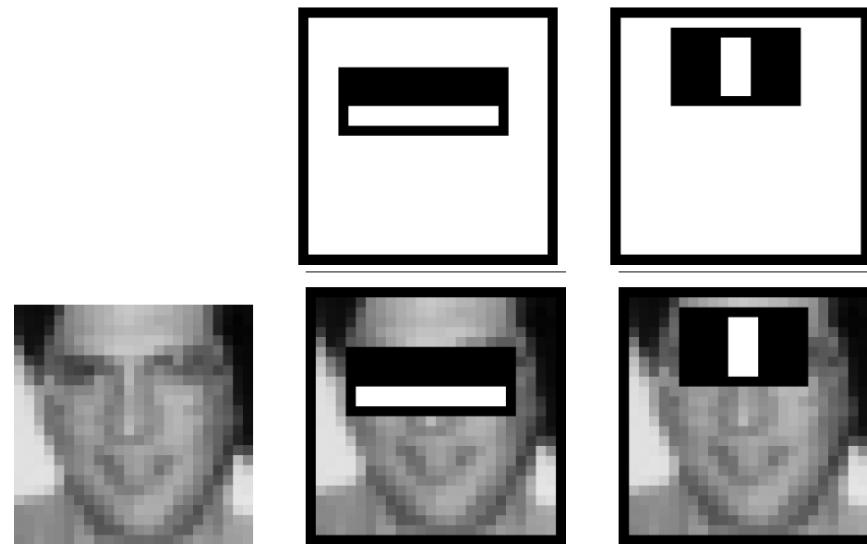








Viola-Jones Face Detector: Results



First two features
selected

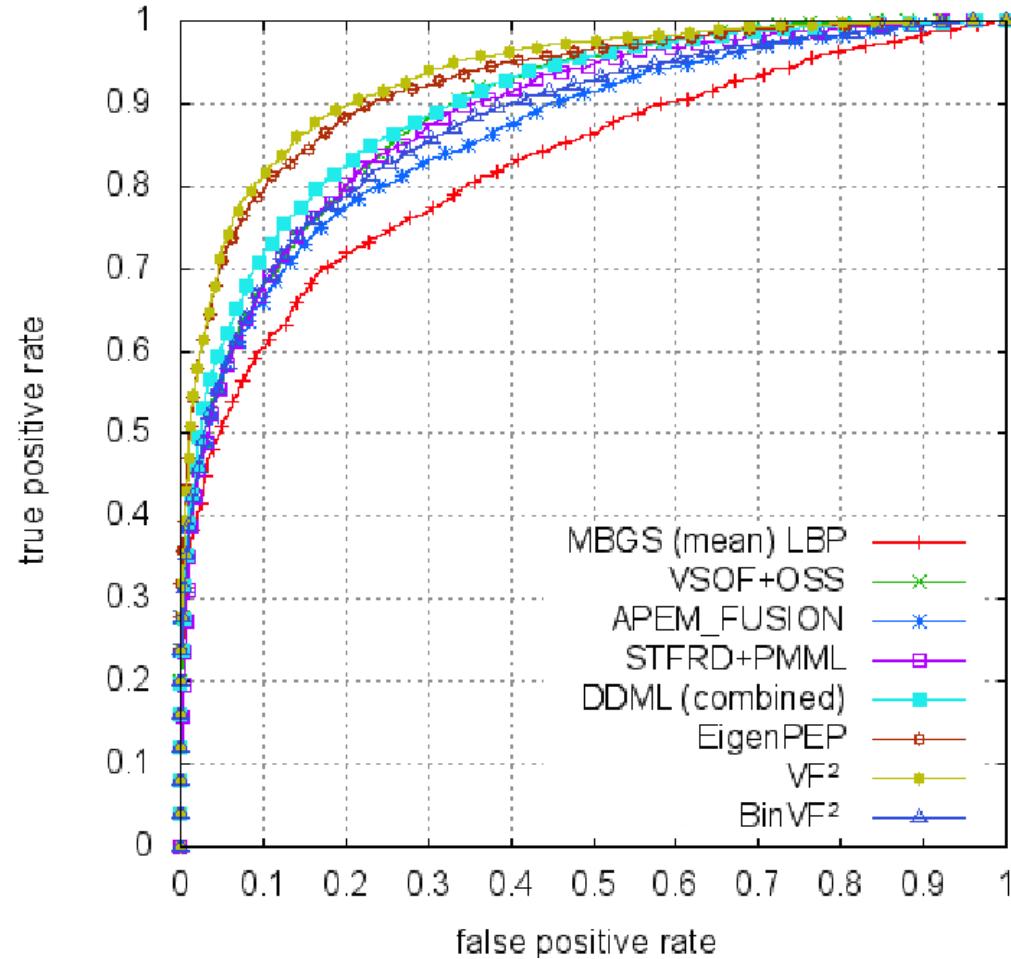
Possible outcomes of a classifier

	is face	is NOT face
says face	TP	FP
Says NOT face	FN	TN

TP : true positive
FP : false positive
TN : true negative
FN : false negative

$$TP + FN = 1.0$$

$$FP + TN = 1.0$$



ROC curve for face recognition.

A **receiver operating characteristic curve**, or **ROC curve**, is a [graphical plot](#) that illustrates the diagnostic ability of a [binary classifier](#) system as its discrimination threshold is varied.

The ROC curve is created by plotting the [true positive rate](#) (TPR) against the [false positive rate](#) (FPR) at various threshold settings. [Wiki](#)

Assuming tests $i = 1, \dots, n$ are independent:

Probability of true face passing test $i = \text{TP}_i$

Probability of true face passing all tests = $\text{TP}_{\text{all}} = \prod_i \text{TP}_i$

Probability of false face passing test $i = \text{FP}_i$

Probability of false face passing all tests = $\text{FP}_{\text{all}} = \prod_i \text{FP}_i$

Assuming tests $i = 1, \dots, n$ are independent:

Probability of true face passing test $i = \text{TP}_i$

Probability of true face passing all tests = $\text{TP}_{\text{all}} = \prod_i \text{TP}_i$

Probability of false face passing test $i = \text{FP}_i$

Probability of false face passing all tests = $\text{FP}_{\text{all}} = \prod_i \text{FP}_i$

We want TP_{all} to be large, and FP_{all} small.

For example, suppose that

$$\text{TP}_i = 0.999, \quad \text{FN}_i = 0.001$$

$$\text{FP}_i = 0.95, \quad \text{TN}_i = 0.05$$

and suppose there are 100 tests. Then

$$\text{TP}_{\text{all}} = 0.999^{100} = 0.904$$

$$\text{FP}_{\text{all}} = 0.95^{100} = 0.006$$

Test is relatively cheap on negative examples.

Expected number of tests applied to a non-face is

$$E[\text{number tests}] = \frac{1}{\text{TN}_i}$$

which is 20 if $\text{FP}_i = 0.95$

Proof:

$$\begin{aligned} E[\text{number tests}] &= 1 \times \text{TN} + 2 \times \text{FP} \times \text{TN} \\ &\quad + 3 \times \text{FP}^2 \times \text{TN} + 4 \times \text{FP}^3 \times \text{TN} + \dots \\ &= \sum_{i=1}^{\infty} i \times \text{FP}^{i-1} \times \text{TN} \\ &= aa \end{aligned} \tag{1}$$

We multiply FP on both sides of Eq. 1, we can get

$$\frac{aa}{\text{TN}} \times \text{FP} = \sum_{i=1}^{\infty} i \times \text{FP}^i \tag{2}$$

$$\frac{aa}{\text{TN}}(1 - \text{FP}) = 1 + \text{FP} + \text{FP}^2 + \text{FP}^3 + \dots \tag{3}$$

$$aa = \frac{1}{1 - \text{FP}} = \frac{1}{\text{TN}} \tag{4}$$

Given $\text{TN} = 0.05$, then $E[\text{number tests}] = 20$.

- Given example images $(x_1, y_1), \dots, (x_n, y_n)$ where $y_i = 0, 1$ for negative and positive examples respectively.
- Initialize weights $w_{1,i} = \frac{1}{2m}, \frac{1}{2l}$ for $y_i = 0, 1$ respectively, where m and l are the number of negatives and positives respectively.
- For $t = 1, \dots, T$:

1. Normalize the weights,

$$w_{t,i} \leftarrow \frac{w_{t,i}}{\sum_{j=1}^n w_{t,j}}$$

so that w_t is a probability distribution.

2. For each feature, j , train a classifier h_j which is restricted to using a single feature. The error is evaluated with respect to w_t , $\epsilon_j = \sum_i w_i |h_j(x_i) - y_i|$.
3. Choose the classifier, h_t , with the lowest error ϵ_t .
4. Update the weights:

$$w_{t+1,i} = w_{t,i} \beta_t^{1-e_i}$$

where $e_i = 0$ if example x_i is classified correctly, $e_i = 1$ otherwise, and $\beta_t = \frac{\epsilon_t}{1-\epsilon_t}$.

- The final strong classifier is:

$$h(x) = \begin{cases} 1 & \sum_{t=1}^T \alpha_t h_t(x) \geq \frac{1}{2} \sum_{t=1}^T \alpha_t \\ 0 & \text{otherwise} \end{cases}$$

where $\alpha_t = \log \frac{1}{\beta_t}$

Table 1: The boosting algorithm for learning a query online. T hypotheses are constructed each using a single feature. The final hypothesis is a weighted linear combination of the T hypotheses where the weights are inversely proportional to the training errors.

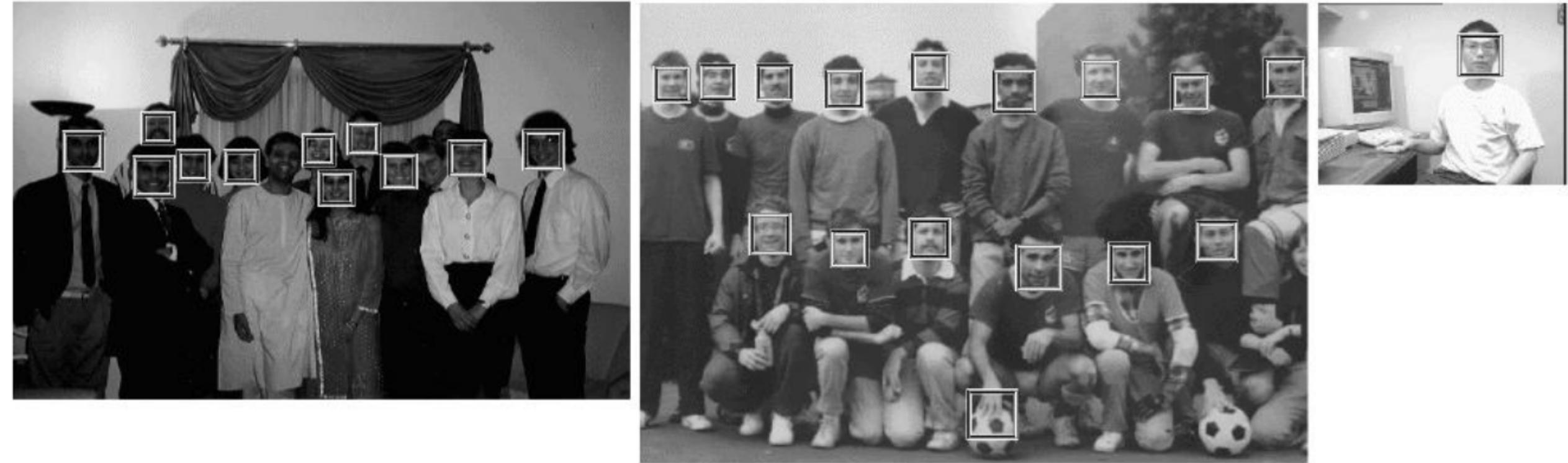
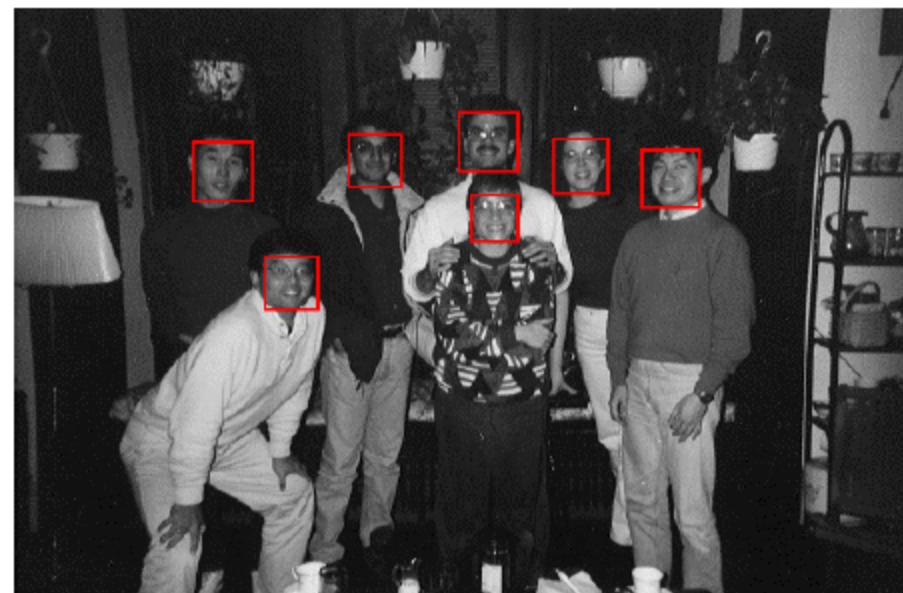
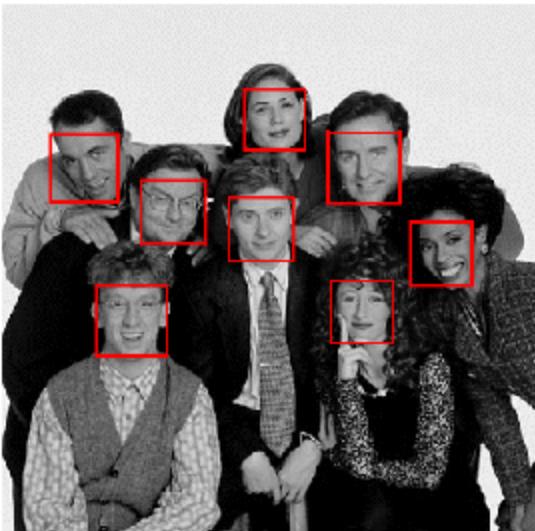
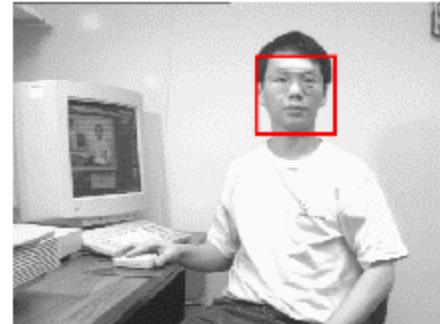
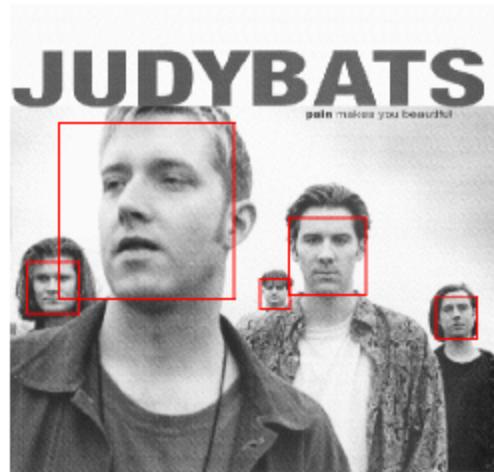
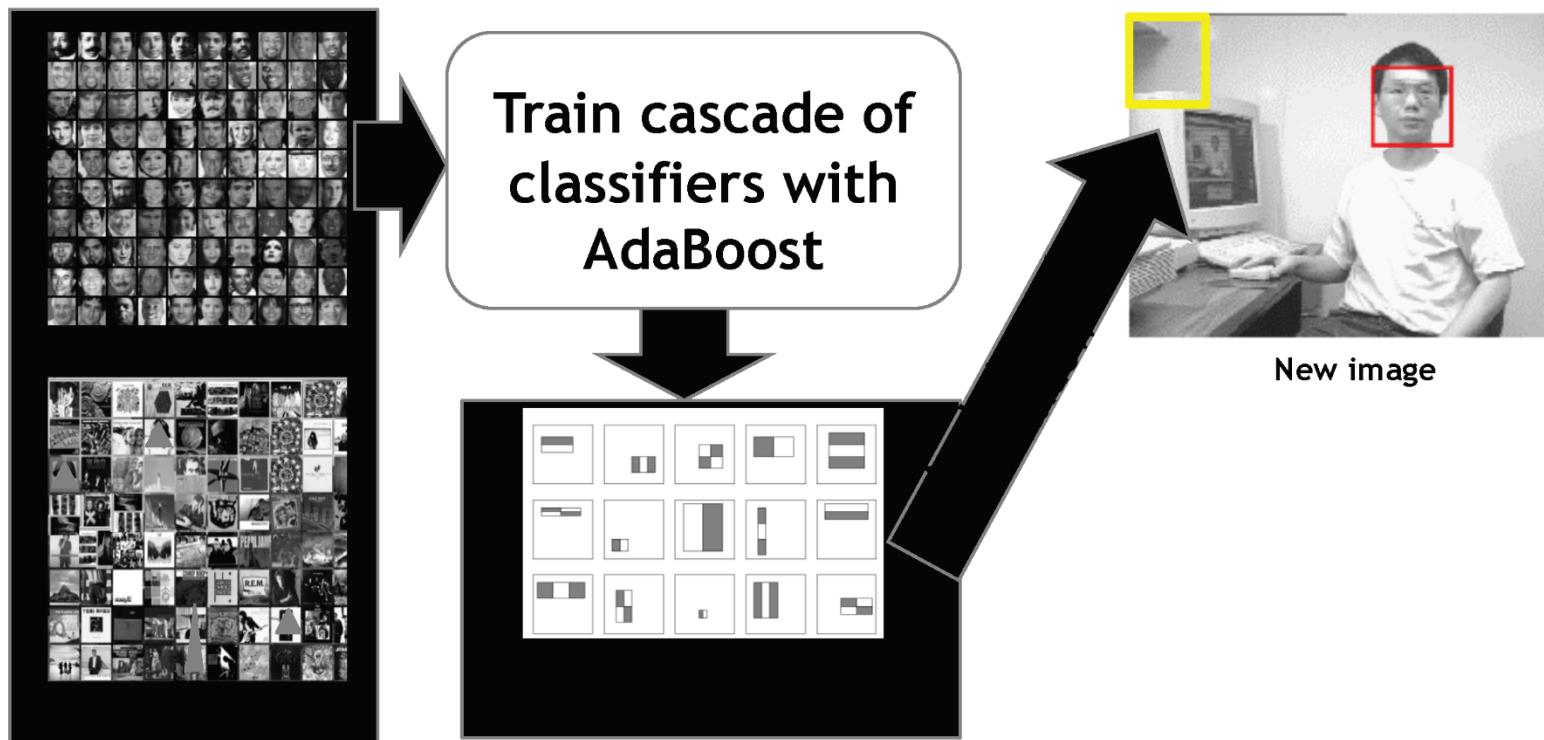


Figure 7: Output of our face detector on a number of test images from the MIT+CMU test set.

Viola-Jones Face Detector: Results



Viola-Jones Face Detector: Summary



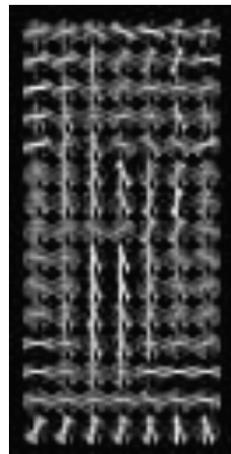
- Train with 5K positives, 350M negatives
- Real-time detector using 38 layer cascade
- 6061 features in final layer
- [Implementation available in OpenCV:
<http://www.intel.com/technology/computing/opencv/>]

Summary (Viola-Jones)

- Fastest known real-time face detector for gray images (in 2001)
- Three contributions with broad applicability:
 - ❖ Cascaded classifier yields rapid classification
 - ❖ AdaBoost as an extremely efficient feature selector
 - ❖ Rectangle Features + Integral Image can be used for rapid image analysis

Linear SVM and HOG for pedestrian detection

Person detection with HoG's & linear SVM's

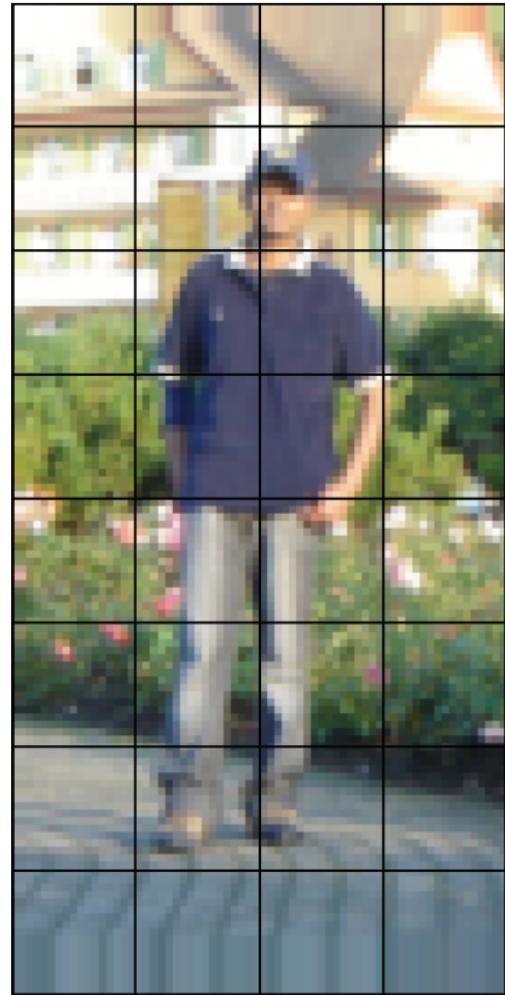


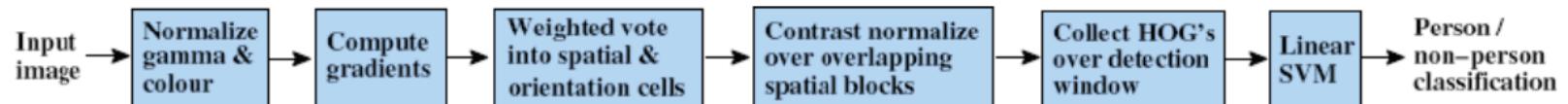
- Histogram of oriented gradients (HoG): Map each grid cell in the input window to a histogram counting the gradients per orientation.
- Train a linear SVM using training set of pedestrian vs. non-pedestrian windows.

Person detection with HoGs & linear SVMs



- Histograms of Oriented Gradients for Human Detection. [Navneet Dalal](#), [Bill Triggs](#), International Conference on Computer Vision & Pattern Recognition - June 2005
- <http://lear.inrialpes.fr/pubs/2005/DT05/>





-1	0	1
----	---	---

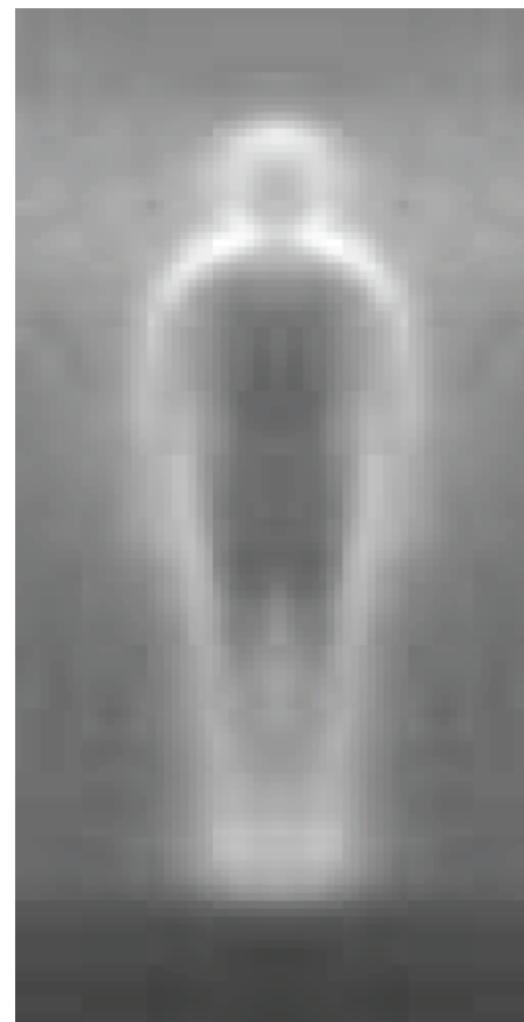
centered

-1	1
----	---

uncentered

1	-8	0	8	-1
---	----	---	---	----

cubic-corrected

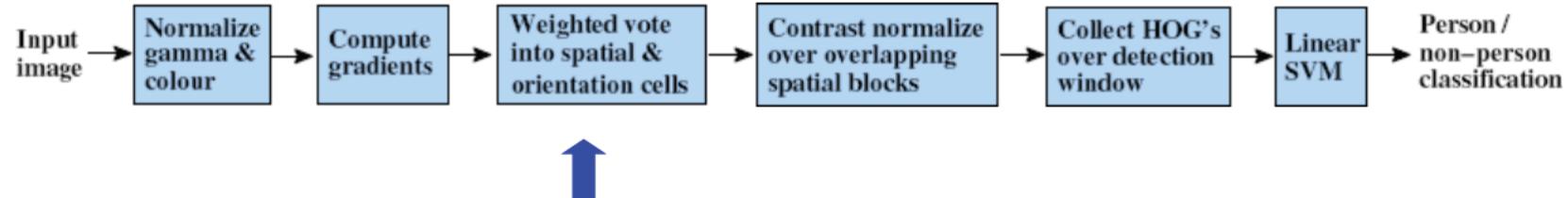


0	1
-1	0

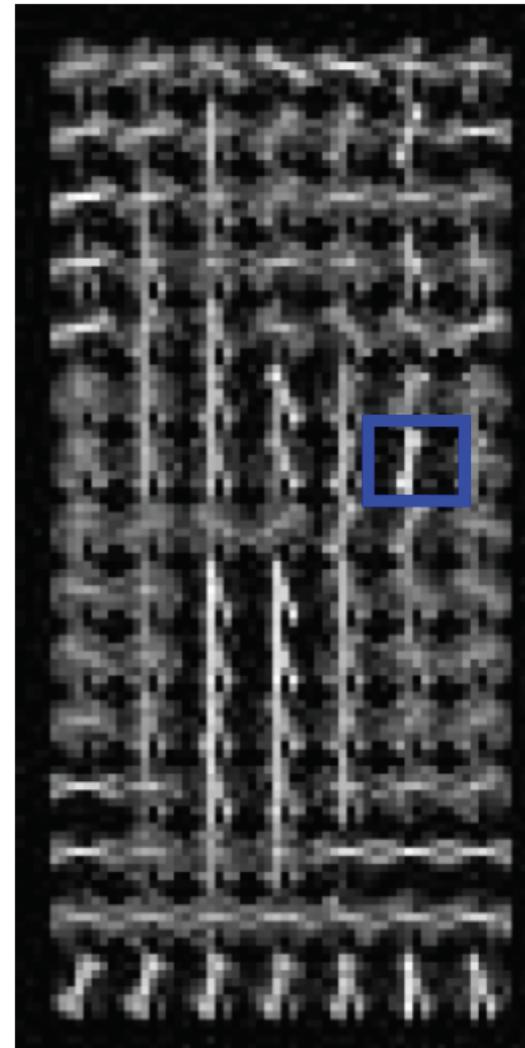
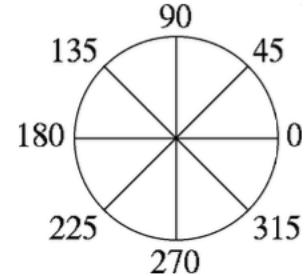
diagonal

-1	0	1
-2	0	2
-1	0	1

Sobel



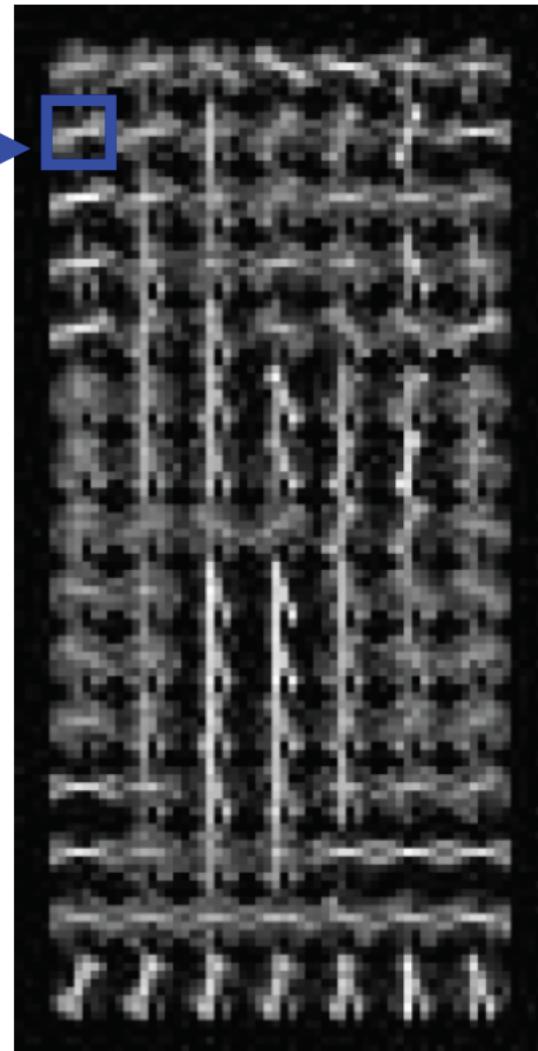
- Histogram of gradient orientations
 - Orientation





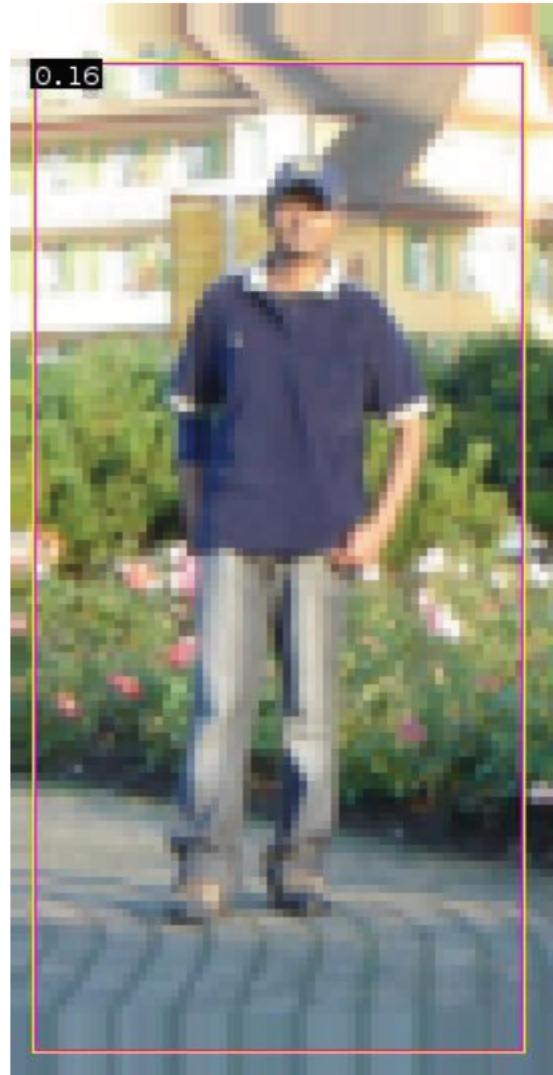
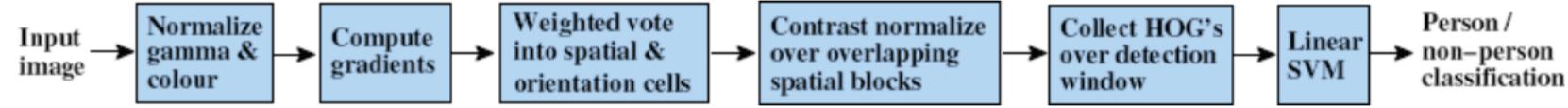
8 orientations

X =



15x7 cells

$$\in R^{840}$$

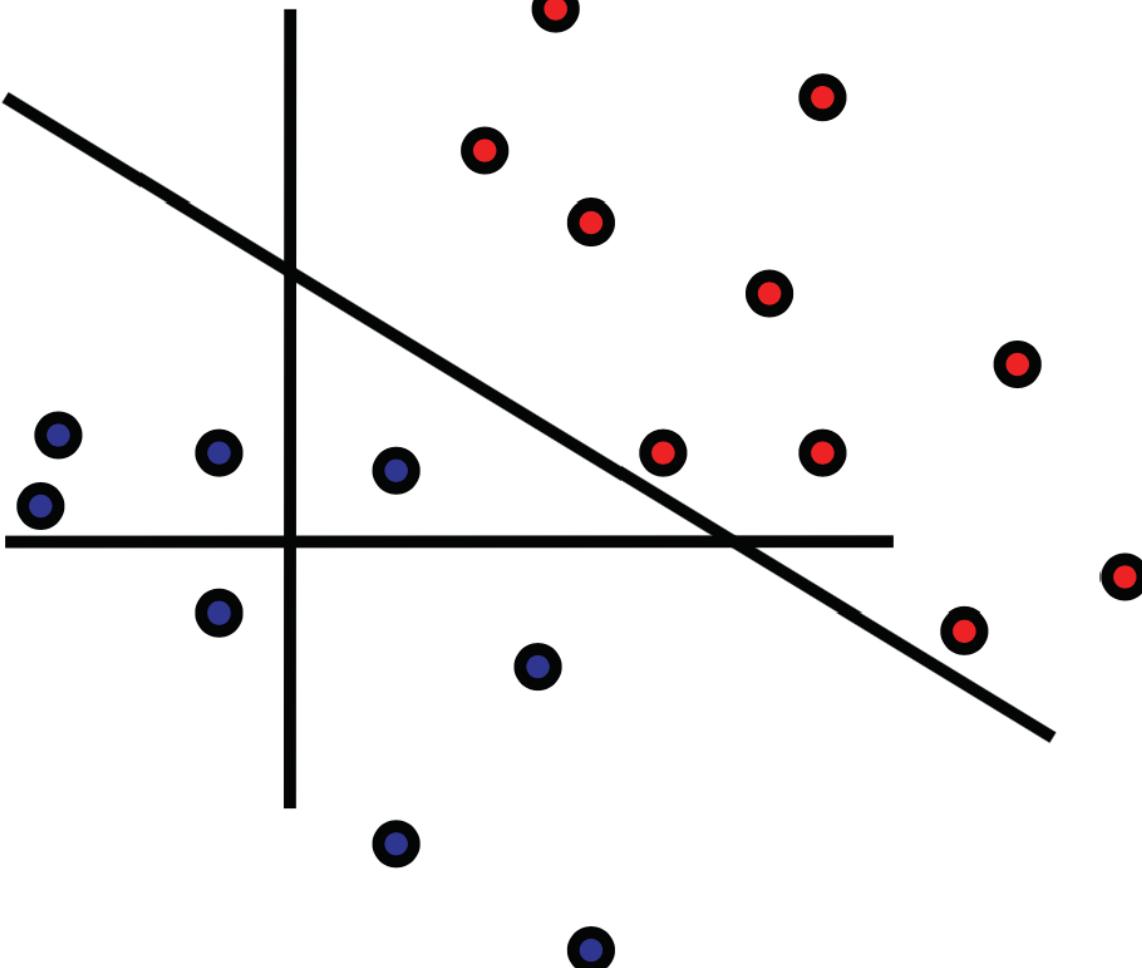


$$0.16 = w^T x - b$$

$$\text{sign}(0.16) = 1$$

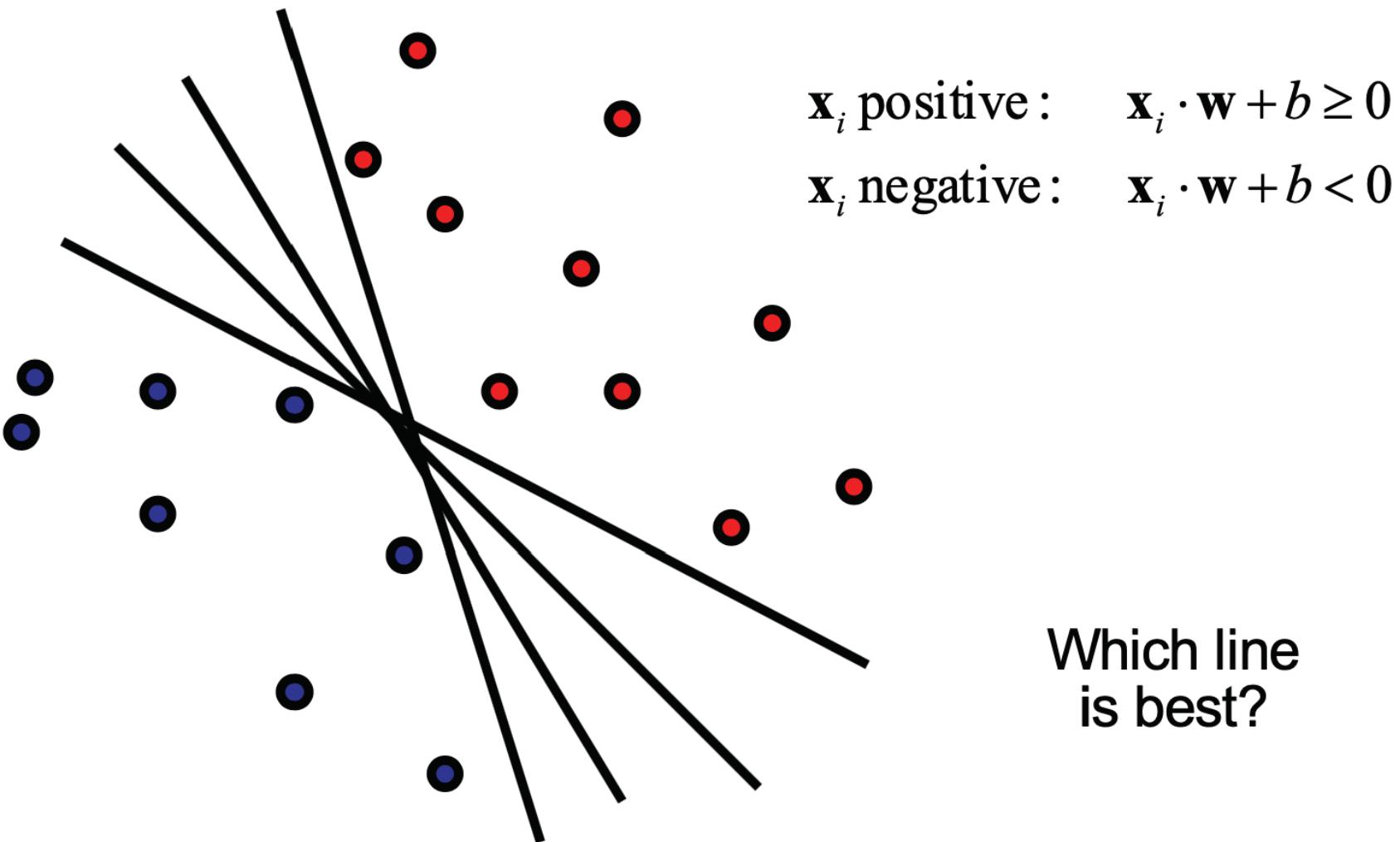
\Rightarrow pedestrian

Linear classifiers

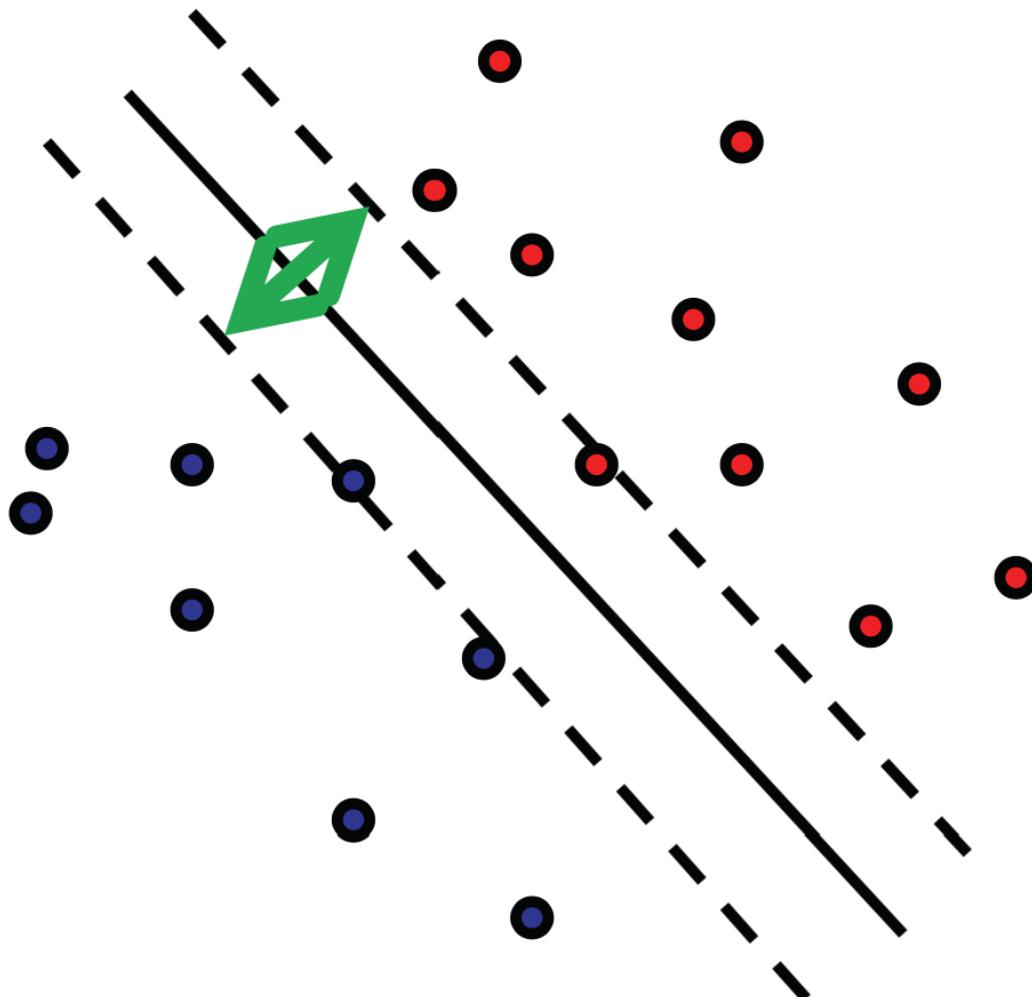


Linear classifiers

- Find linear function to separate positive and negative examples



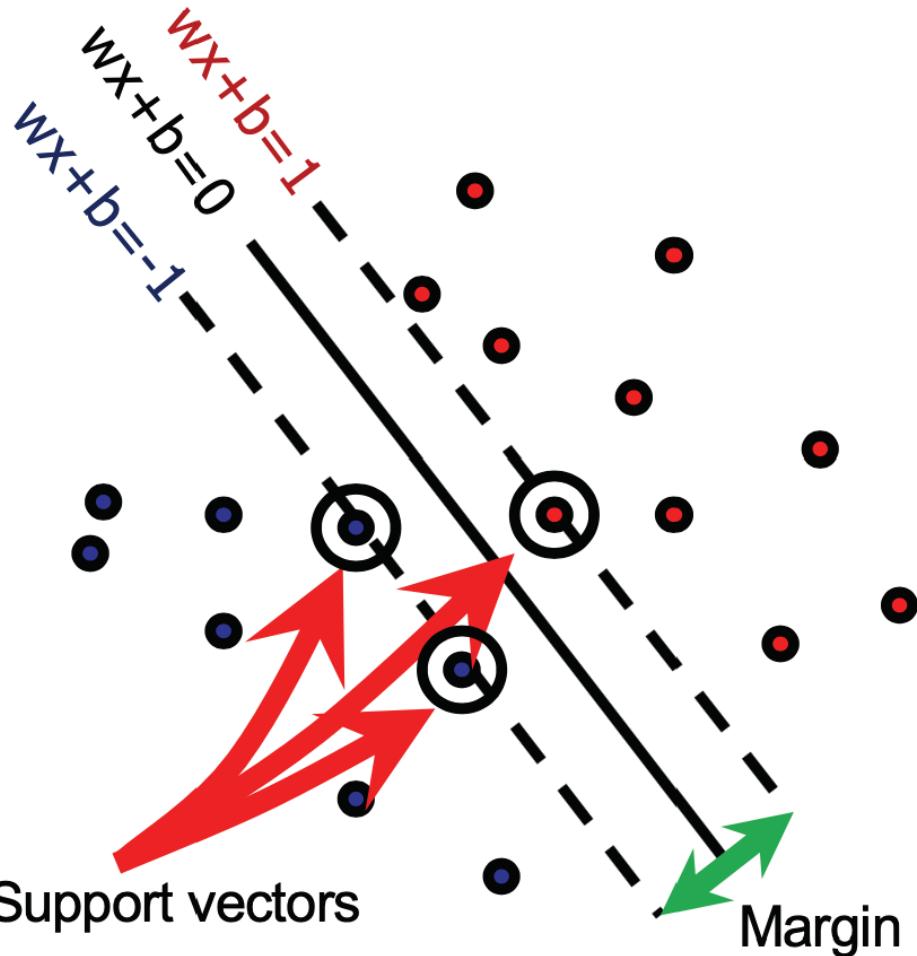
Support Vector Machines (SVMs)



- Discriminative classifier based on *optimal separating line* (for 2d case)
- Maximize the *margin* between the positive and negative training examples

Support vector machines

- Want line that maximizes the margin.



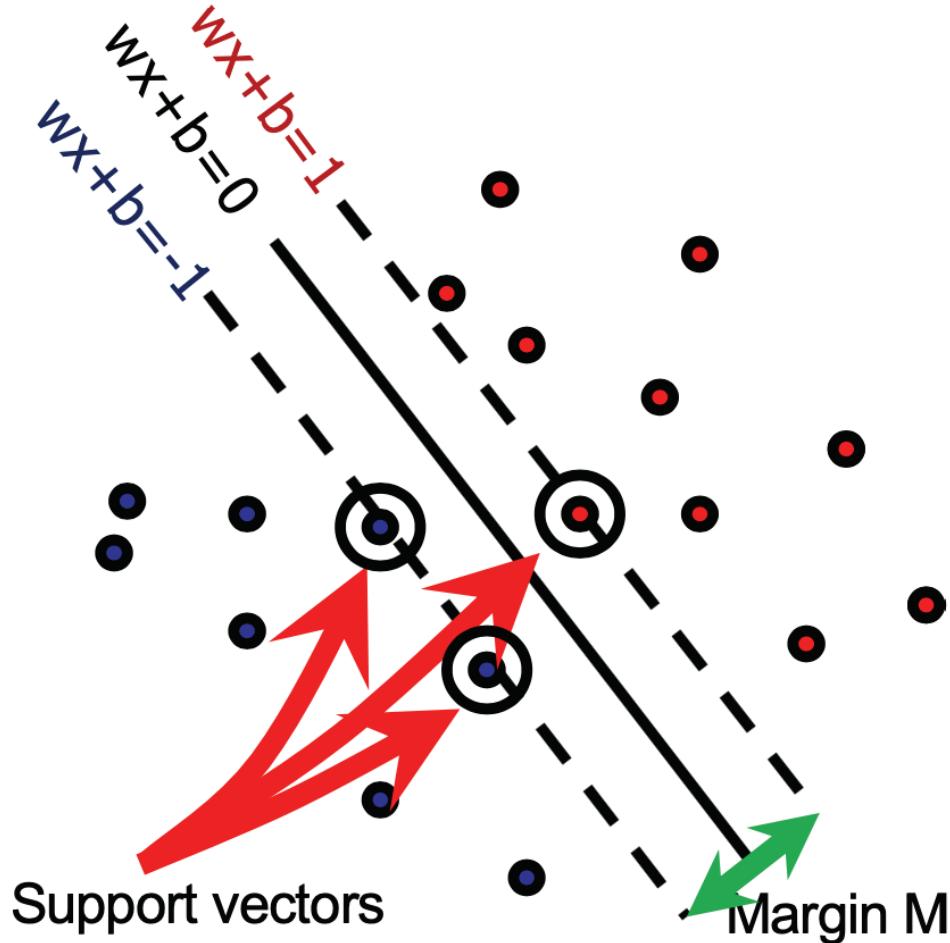
\mathbf{x}_i positive ($y_i = 1$): $\mathbf{x}_i \cdot \mathbf{w} + b \geq 1$

\mathbf{x}_i negative ($y_i = -1$): $\mathbf{x}_i \cdot \mathbf{w} + b \leq -1$

For support vectors, $\mathbf{x}_i \cdot \mathbf{w} + b = \pm 1$

Support vector machines

- Want line that maximizes the margin.



\mathbf{x}_i positive ($y_i = 1$): $\mathbf{x}_i \cdot \mathbf{w} + b \geq 1$

\mathbf{x}_i negative ($y_i = -1$): $\mathbf{x}_i \cdot \mathbf{w} + b \leq -1$

For support, vectors, $\mathbf{x}_i \cdot \mathbf{w} + b = \pm 1$

Distance between point
and line:

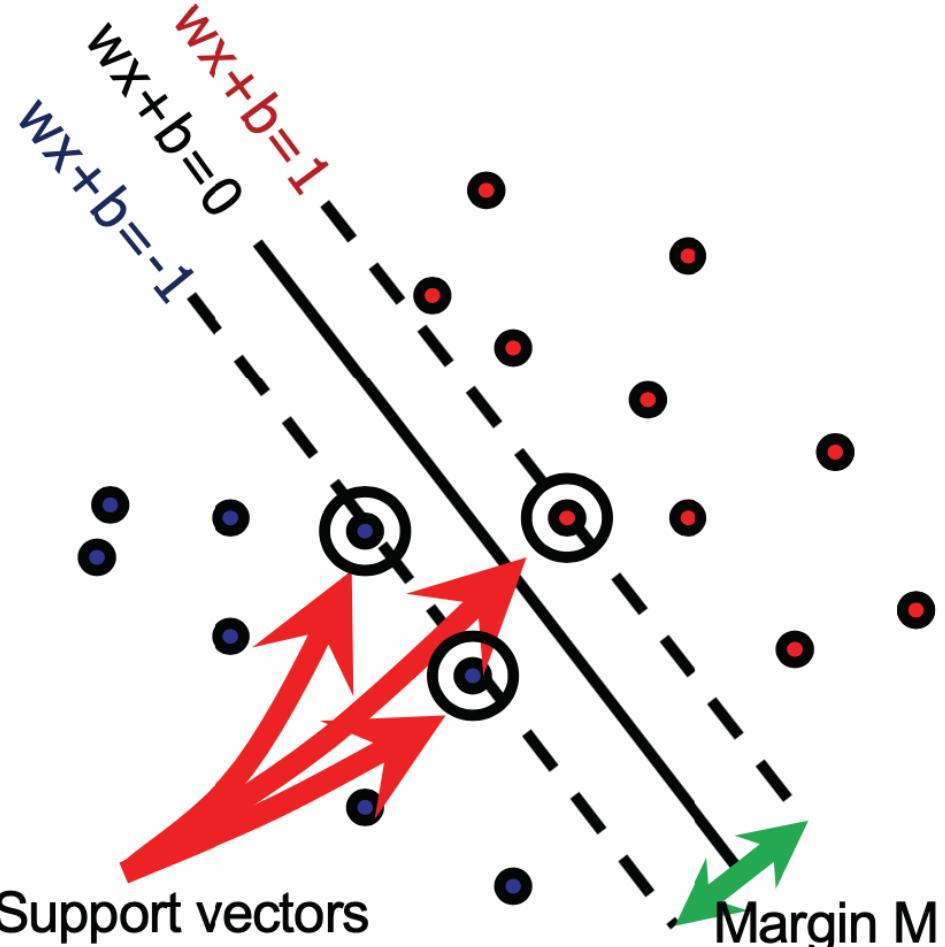
$$\frac{|\mathbf{x}_i \cdot \mathbf{w} + b|}{\|\mathbf{w}\|}$$

For support vectors:

$$\frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|} = \frac{\pm 1}{\|\mathbf{w}\|} \quad M = \left| \frac{1}{\|\mathbf{w}\|} - \frac{-1}{\|\mathbf{w}\|} \right| = \frac{2}{\|\mathbf{w}\|}$$

Support vector machines

- Want line that maximizes the margin.



\mathbf{x}_i positive ($y_i = 1$): $\mathbf{x}_i \cdot \mathbf{w} + b \geq 1$

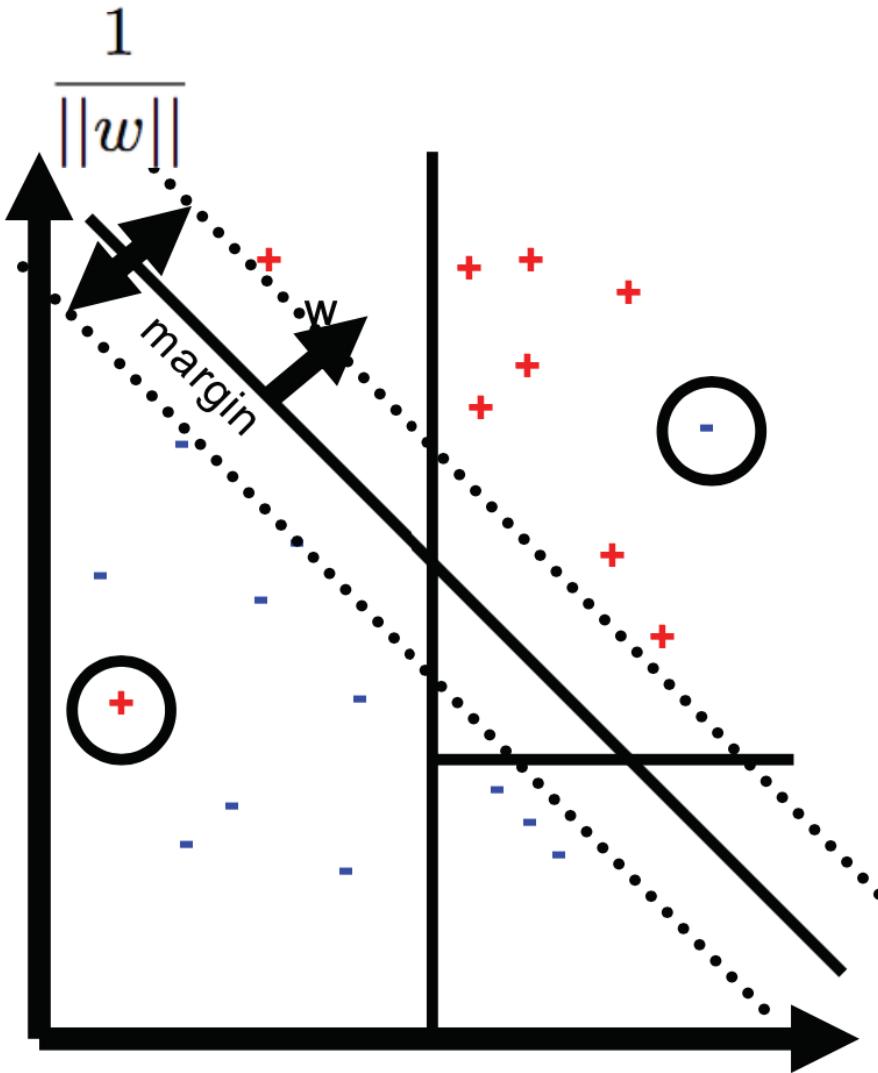
\mathbf{x}_i negative ($y_i = -1$): $\mathbf{x}_i \cdot \mathbf{w} + b \leq -1$

For support vectors, $\mathbf{x}_i \cdot \mathbf{w} + b = \pm 1$

Distance between point and line:
$$\frac{|\mathbf{x}_i \cdot \mathbf{w} + b|}{\|\mathbf{w}\|}$$

Therefore, the margin is $2 / \|\mathbf{w}\|$

Support vector machines



- Simple decision
- Good classification
- Good generalization

$$\min_{w, \xi} \frac{\|w\|^2}{2} + C \sum_j \xi_j$$

$$y_j w^T x_j \geq 1 - \xi_j$$

$$\xi_j \geq 0$$

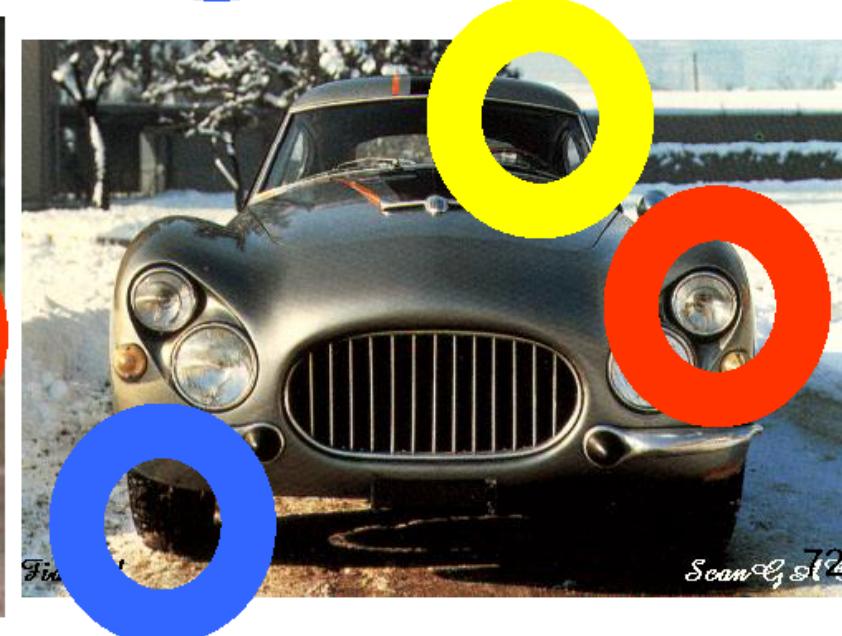
How do I solve the SVM problem?

- It's a convex optimization problem
 - Can solve in Matlab (don't)
- Download from the web
 - SMO: Sequential Minimal Optimization
 - SVM-Light <http://svmlight.joachims.org/>
 - LibSVM <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>
 - LibLinear <http://www.csie.ntu.edu.tw/~cjlin/liblinear/>
 - SVM-Perf <http://svmlight.joachims.org/>
 - Pegasos <http://ttic.uchicago.edu/~shai/>



Bag of Words model for object detection

Different appearance, similar parts



Object

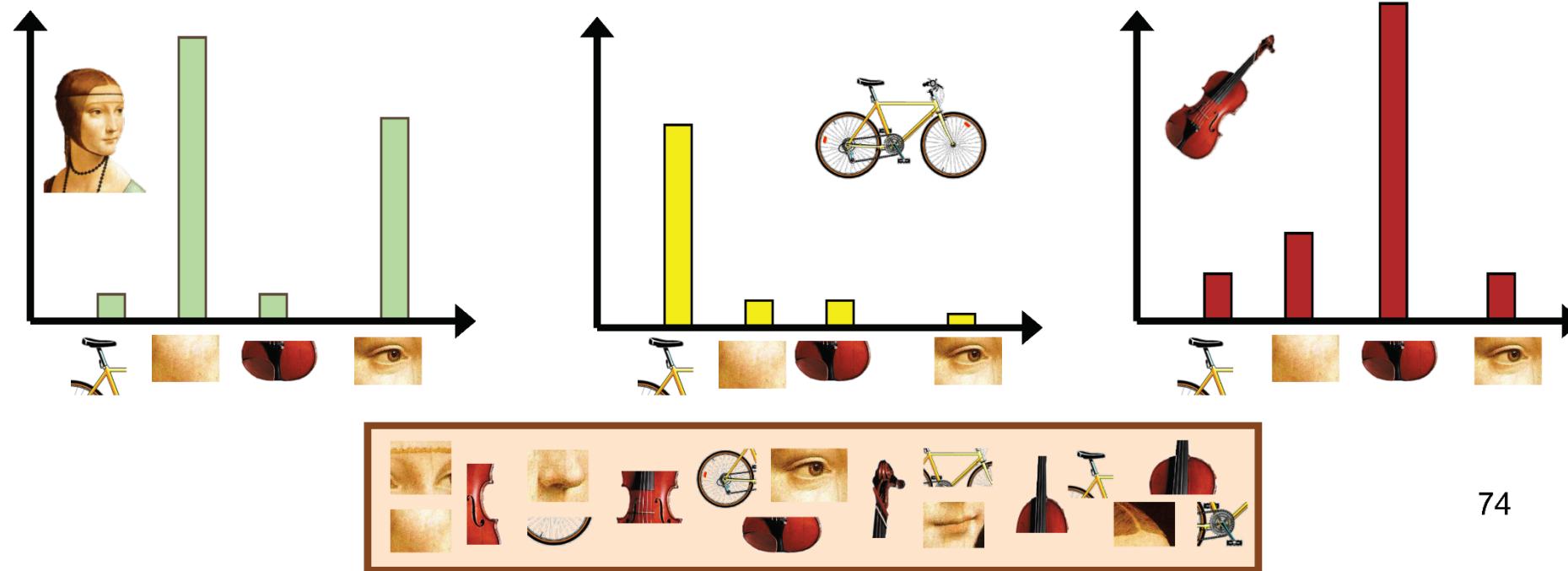


Bag of ‘words’

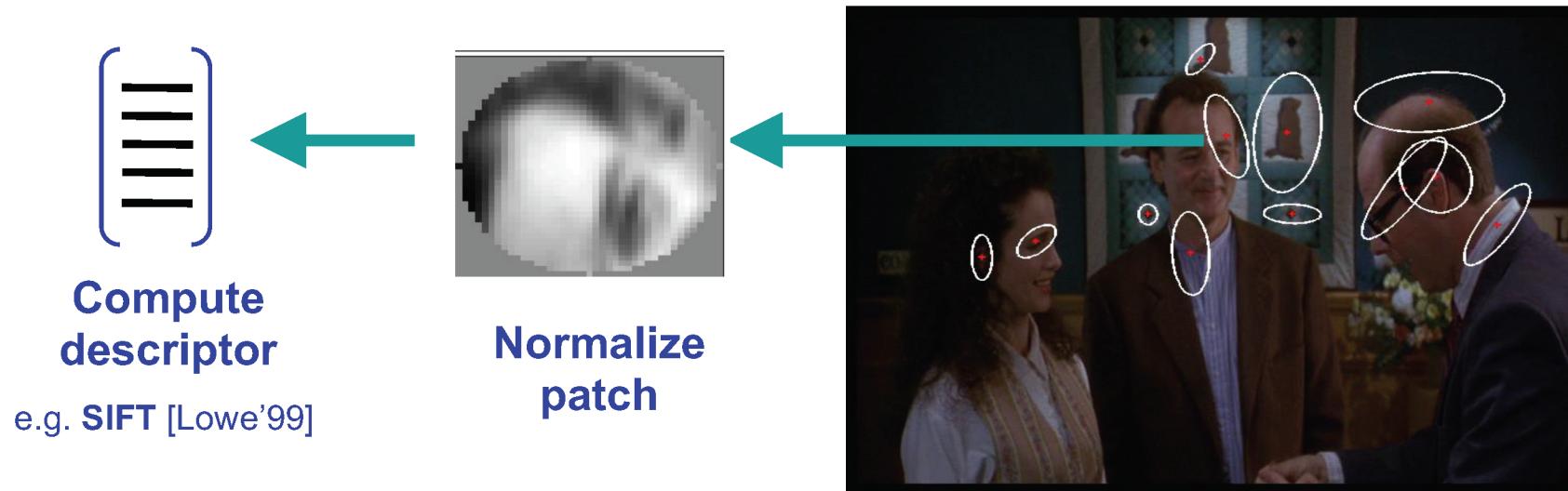


Bag of Words

- Independent features
- Histogram representation



1. Feature detection and representation



Detect patches

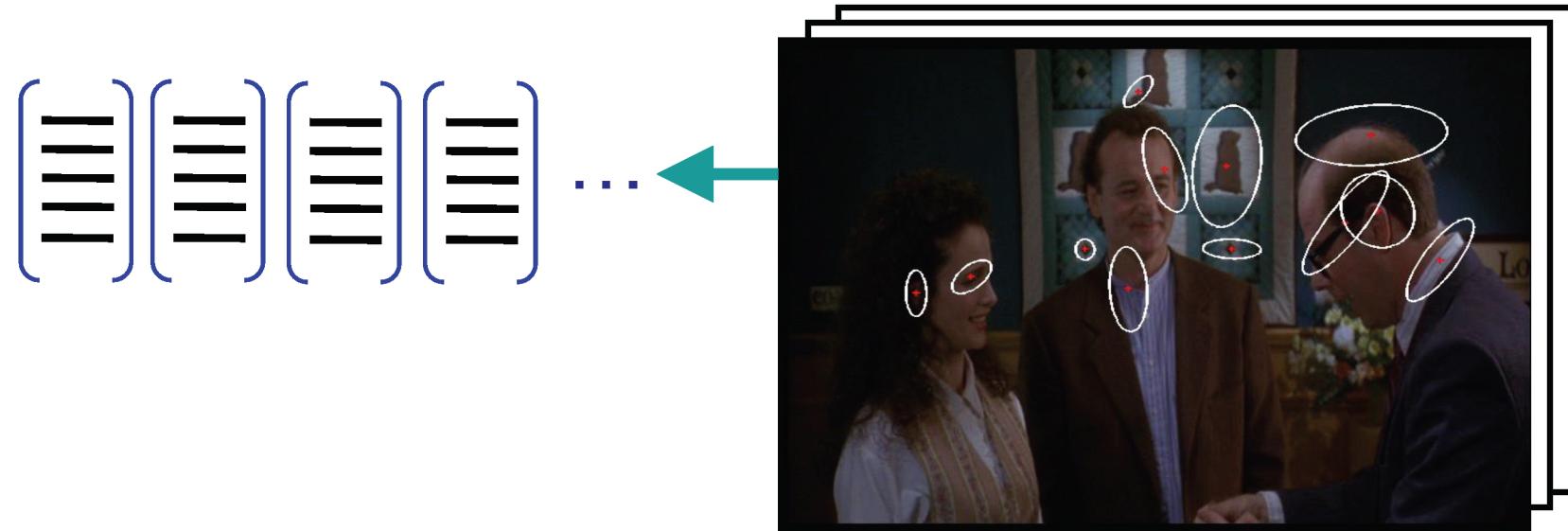
[Mikojaczyk and Schmid '02]

[Mata, Chum, Urban & Pajdla, '02]

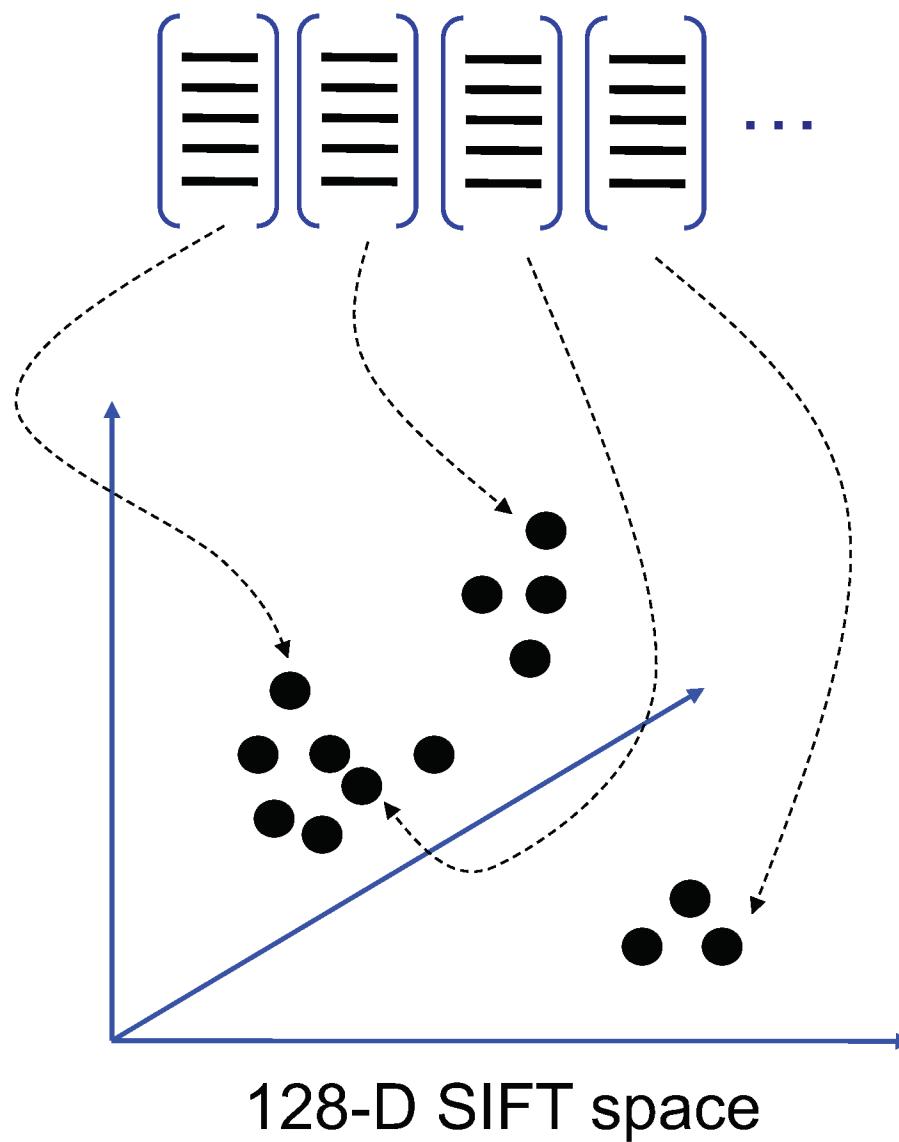
[Sivic & Zisserman, '03]

Local interest operator
or
Regular grid

1. Feature detection and representation



2. Codewords dictionary formation



2. Codewords dictionary formation

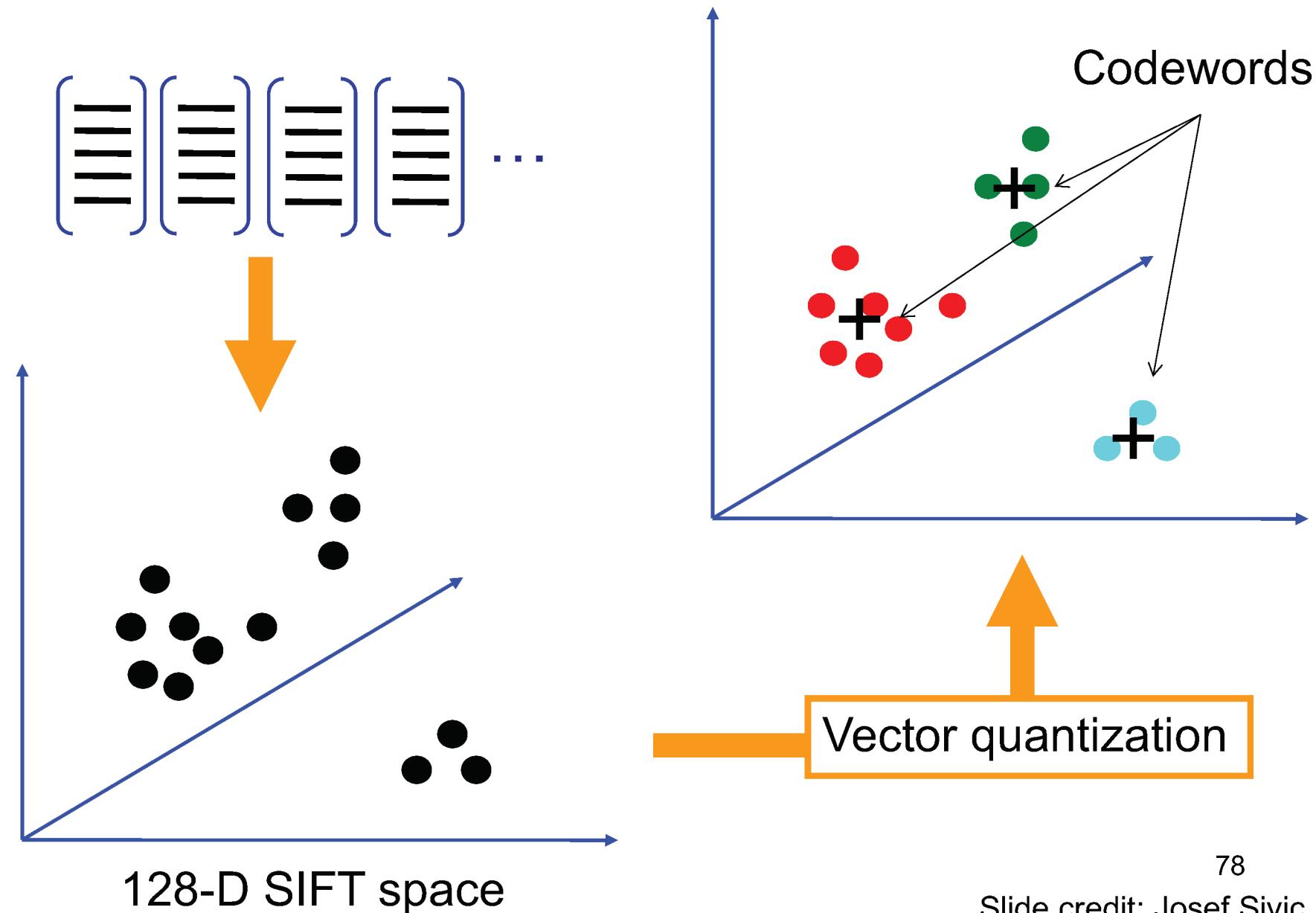


Image patch examples of codewords

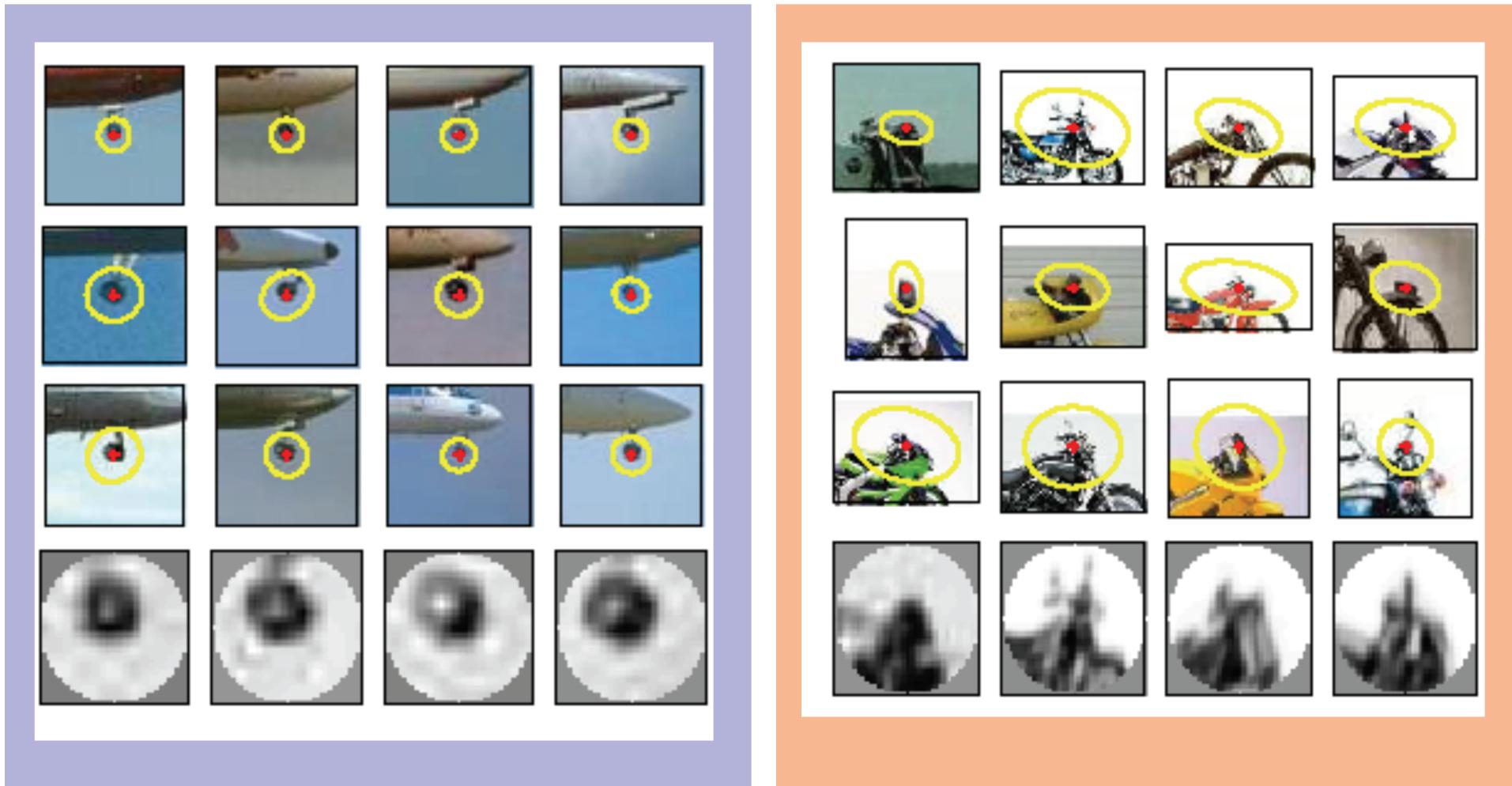
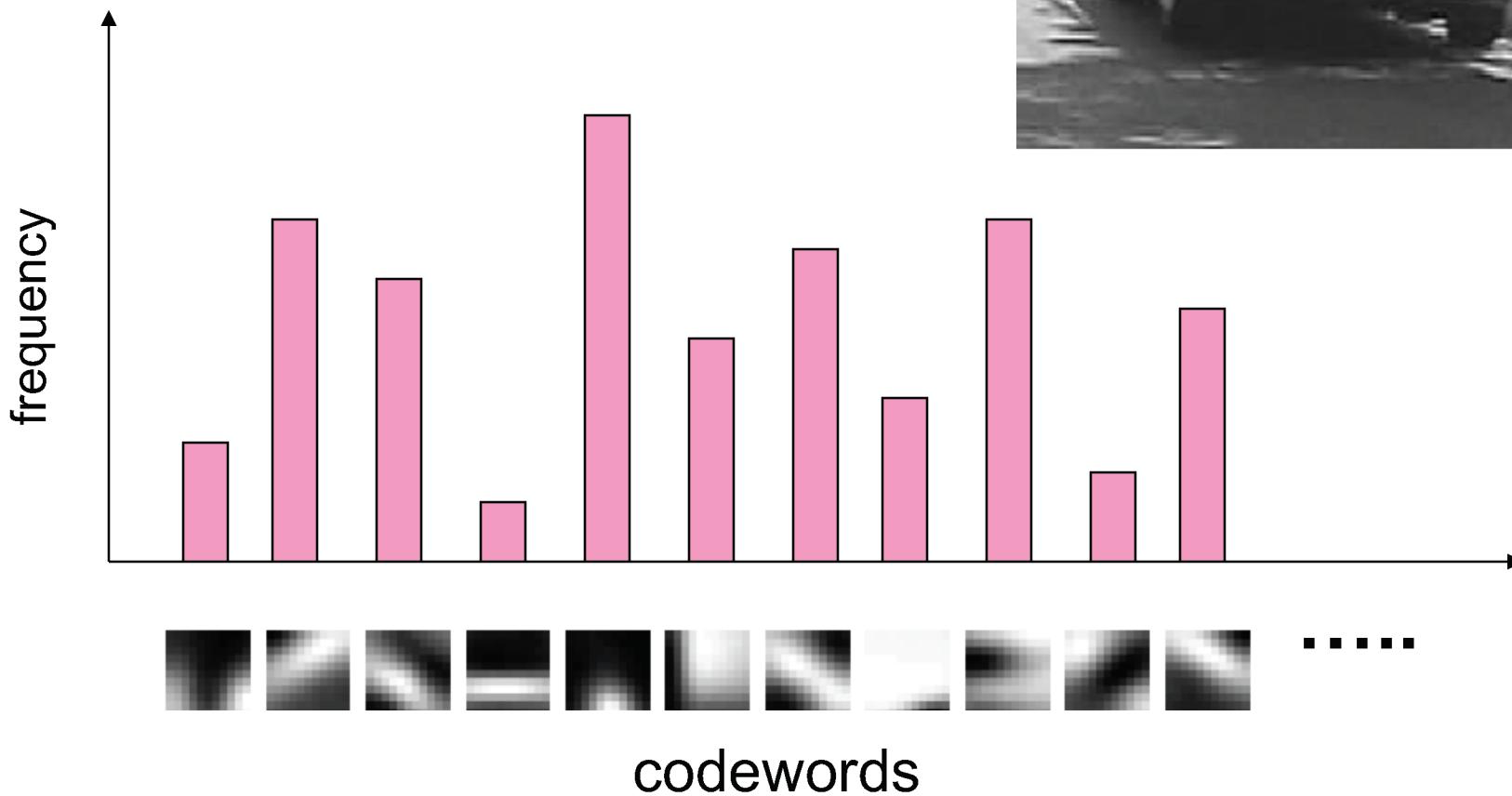


Image representation

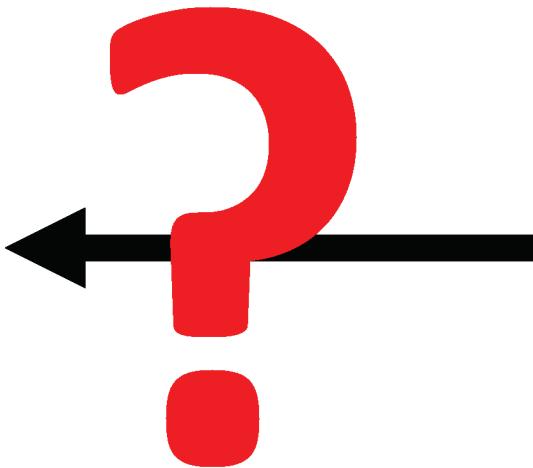
Histogram of features
assigned to each cluster



Uses of BoW representation

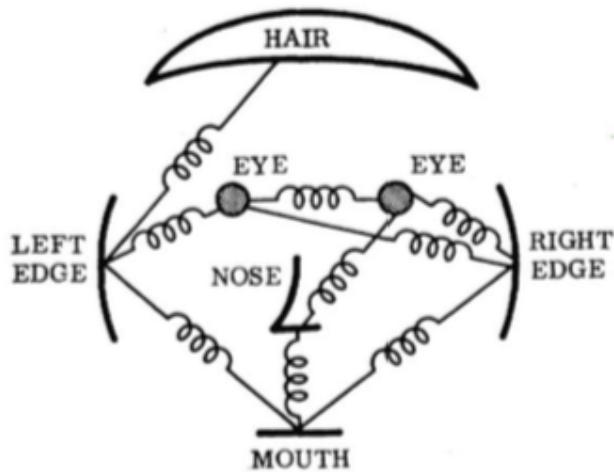
- Treat as feature vector for standard classifier
 - e.g k-nearest neighbors, support vector machine
- Cluster BoW vectors over image collection
 - Discover visual themes

What about spatial info?



Deformable Part-based Models

- Felzenszwalb et al. PAMI'10
- Winner of the PASCAL detection challenge (2008,2009)
- Objects are decomposed into parts and spatial relations among parts



Fischler and Elschlager '73

Deformable Part-based Models

- Score of hypothesis

The diagram illustrates the decomposition of a score function into two components: a "data term" and a "spatial prior".

“data term”

$$\text{score}(p_0, \dots, p_n) = \sum_{i=0}^n F_i \cdot \phi(H, p_i)$$

A red arrow points from the term F_i to the word "filters" below it.

“spatial prior”

$$- \sum_{i=1}^n d_i \cdot (dx_i^2, dy_i^2)$$

A blue arrow points from the term d_i to the words "displacements" below it.

Below the "spatial prior" section, the words "deformation parameters" are written in blue.

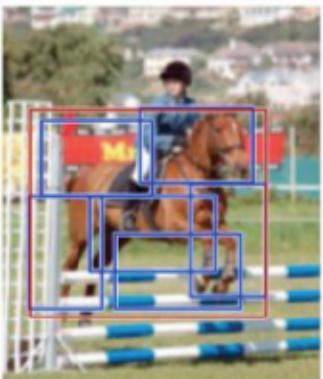
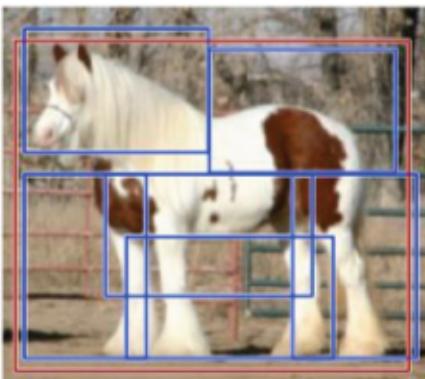


Score is sum of filter scores minus deformation costs

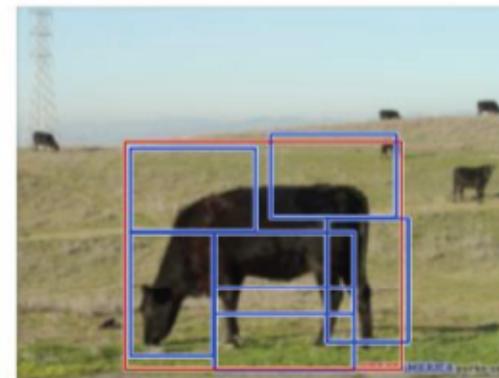
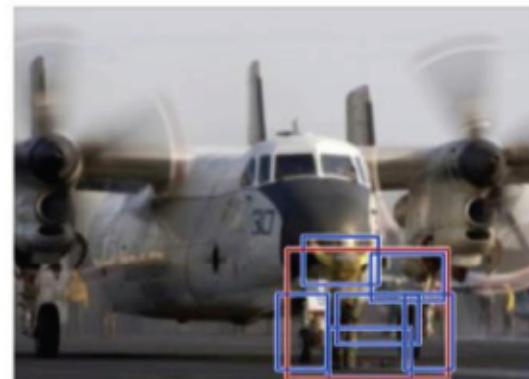
Deformable Part-based Models: Results

- Horse detections

high scoring true positives

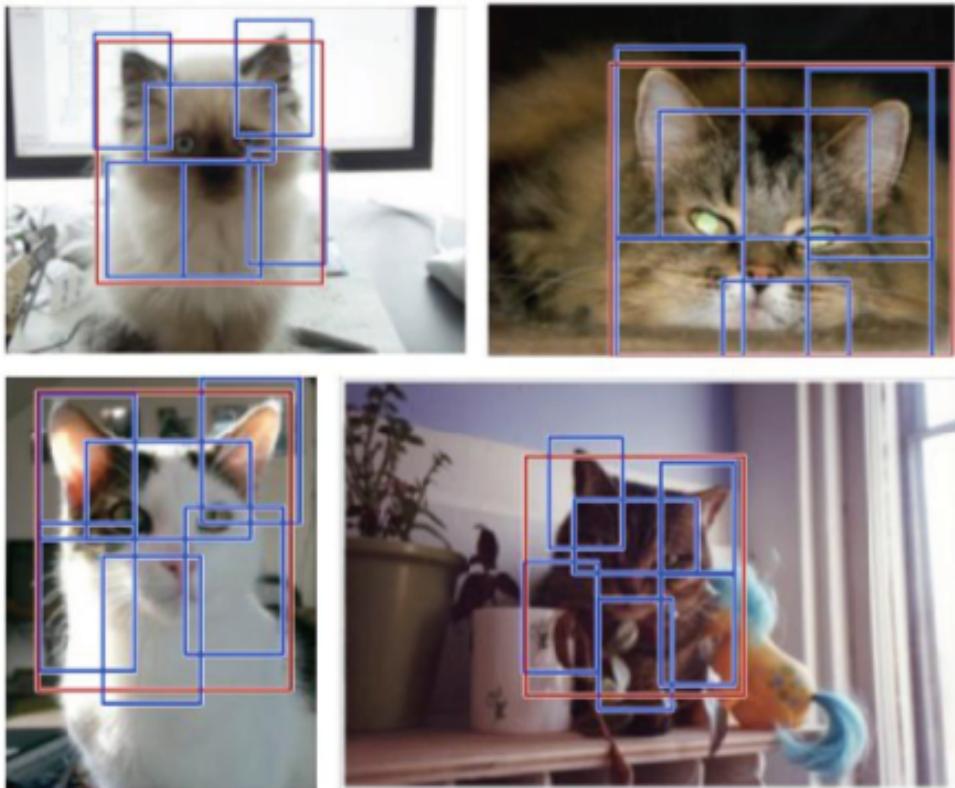


high scoring false positives

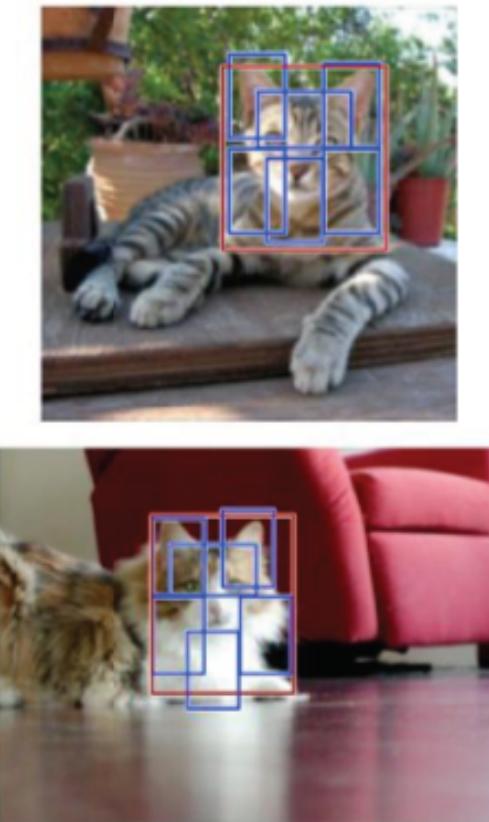


- Cat detections

high scoring true positives

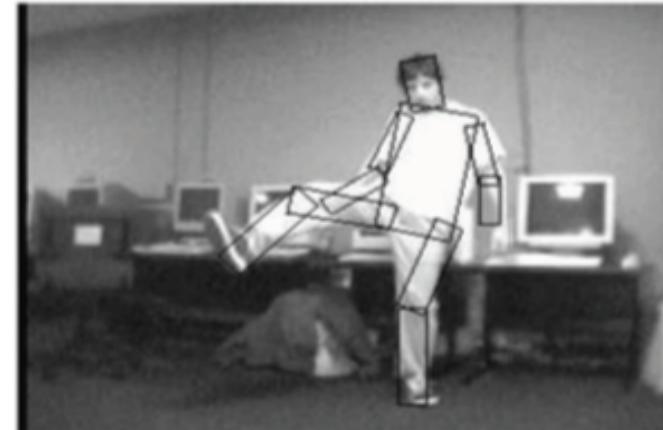
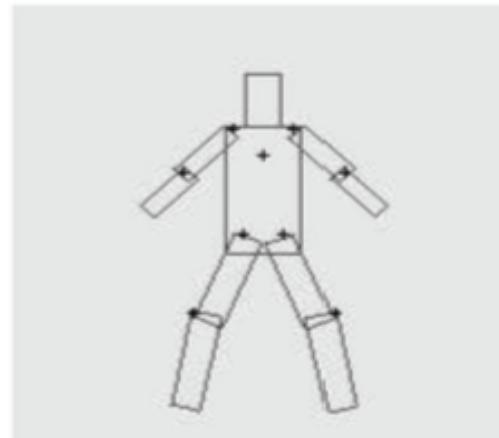


high scoring false positives
(not enough overlap)



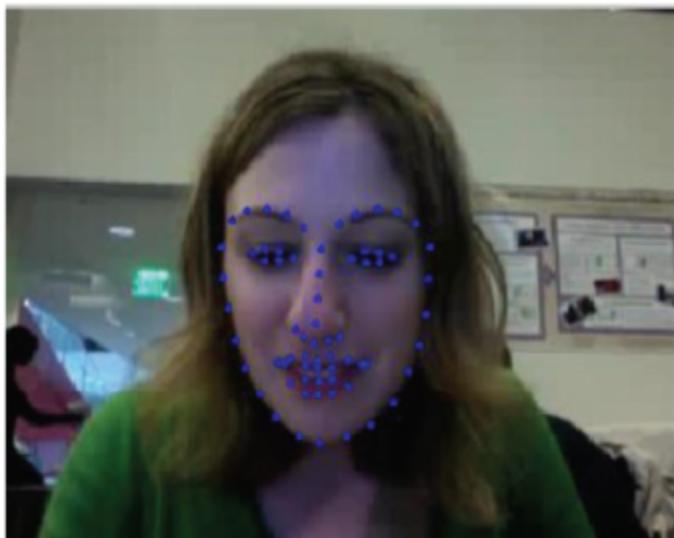
Part-based representation

- Tree model → Articulated objects



Part-based models: Pose estimation

- Pose estimation in video (Ramanan et al, 2007; Yang et al, 2015)



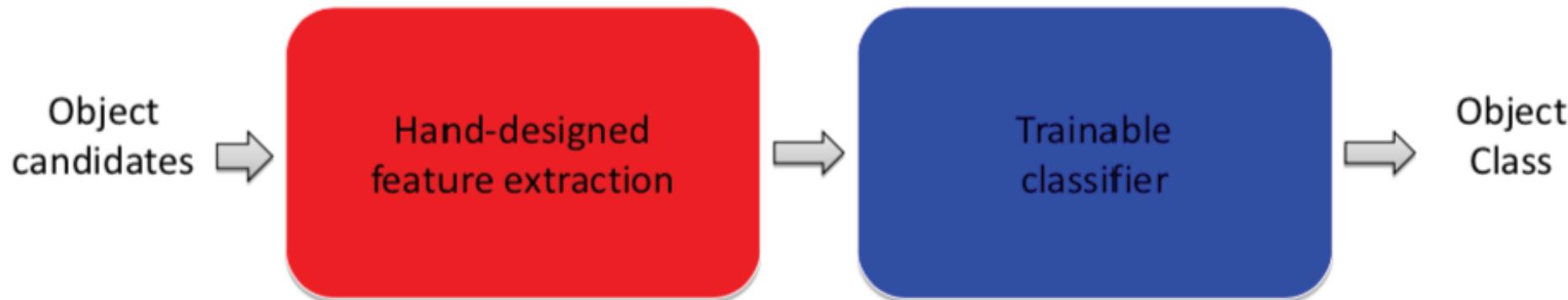
Face capturing



Dancing

DPMs: Limitations

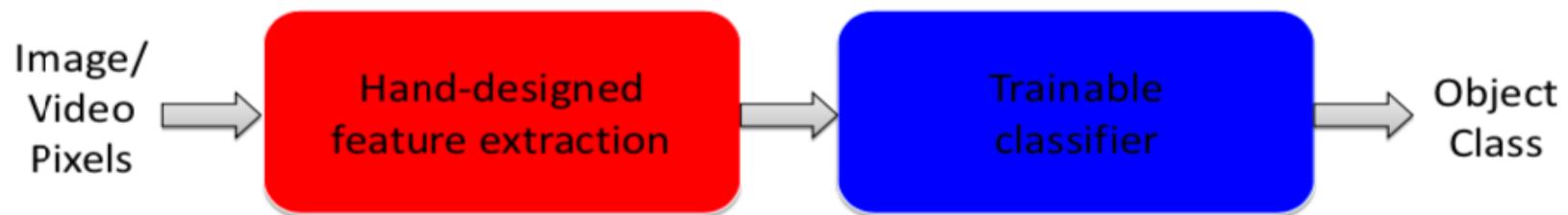
- Manually designed feature (HOG)
- Pre-defined object-part structure
- Shallow structure: limited representation capacity



Move into “deep learning” regime

Next: “deep” architectures

Traditional recognition: “Shallow” architecture



Deep learning: “Deep” architecture



Readings

- M. Turk and A. Pentland, [Face Recognition using Eigenfaces](#), CVPR 1991
- M. Turk and A. Pentland, Eigenfaces for recognition, Journal of Cognitive Neuroscience, 1991
- P. Viola and M. Jones, Robust Real-time Object Detection, IJCV 2001
- N. Dalal and B. Triggs, Histograms of Oriented Gradients for Human Detection, CVPR 2005
- P. Felzenszwalb, R. Girshick, D. McAllester and D. Ramanan, Object Detection with Discriminatively Trained Part Based Models, TPAMI 2010