

# Skin and Face Detection

Linda Shapiro

EE/CSE 576

# What's Coming

1. Review of Bakic flesh detector
2. Fleck and Forsyth flesh detector
3. Details of Rowley face detector
4. Review of the basic AdaBoost algorithm
5. The Viola Jones face detector features
6. The modified AdaBoost algorithm that is used in Viola-Jones face detection

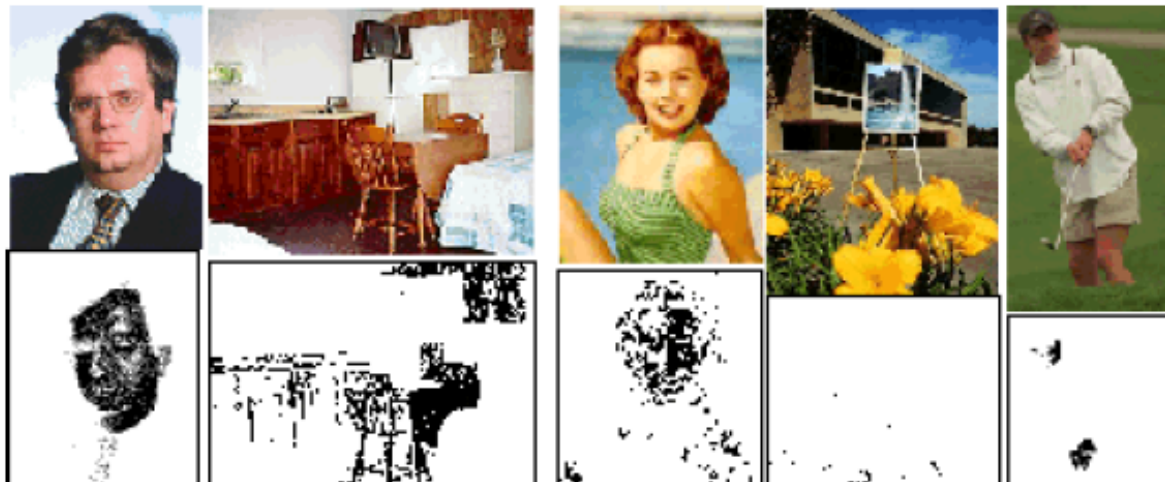
# Object Detection

- Example: Face Detection



([Rowley, Baluja & Kanade, 1998](#))

- Example: Skin Detection



([Jones & Rehg, 1999](#))

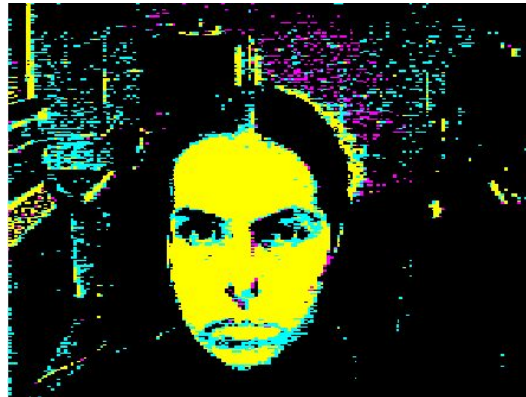
# Review: Bakic Flesh Finder

- Convert pixels to normalized (r,g) space
- Train a binary classifier to recognize pixels in this space as skin or not skin by giving it lots of examples of both classes.
- On test images, have the classifier label the skin pixels.
- Find large connected components.

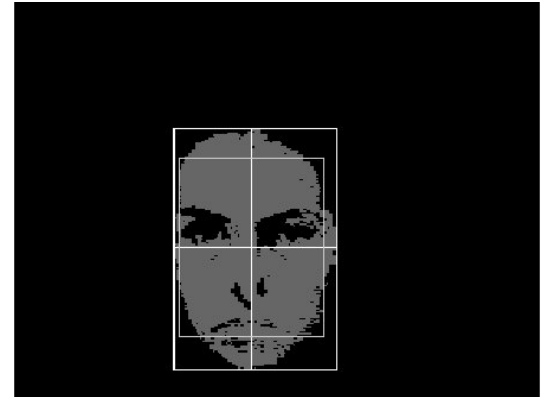
# Finding a face in a video frame



input video frame



pixels classified in  
normalized r-g space



largest connected  
component with aspect  
similar to a face

(all work contributed by Vera Bakic)

# Fleck and Forsyth's Flesh Detector

- Convert RGB to HSI
- Use the intensity component to compute a texture map  
 $\text{texture} = \text{med2} ( | I - \text{med1}(I) | )$
- If a pixel falls into either of the following ranges,  
it's a potential skin pixel

$\text{texture} < 5, 110 < \text{hue} < 150, 20 < \text{saturation} < 60$   
 $\text{texture} < 5, 130 < \text{hue} < 170, 30 < \text{saturation} < 130$

median filters of  
radii 4 and 6

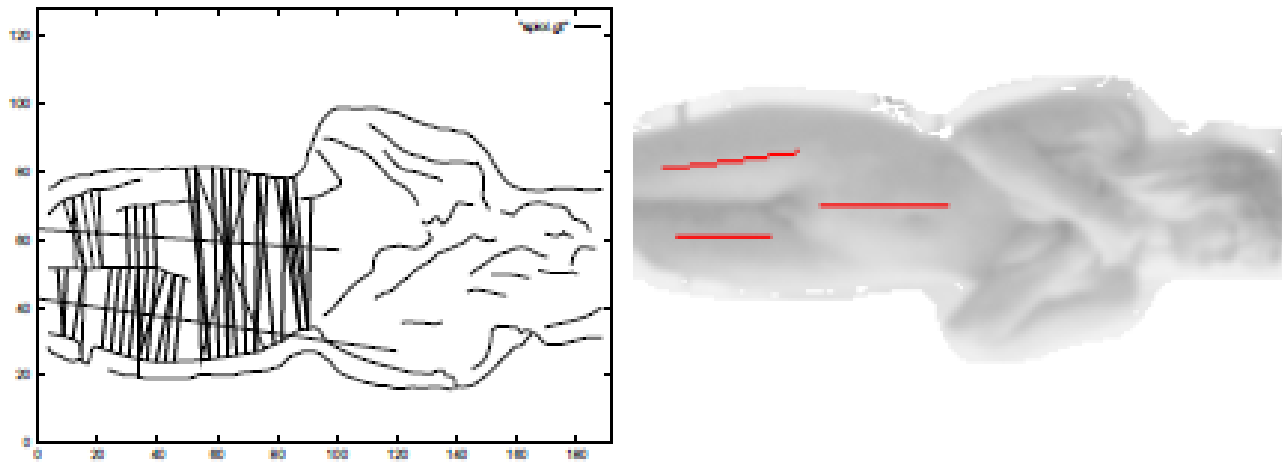
\* Margaret Fleck, David Forsyth, and Chris Bregler (1996)  
"Finding Naked People," 1996 **European Conference on  
Computer Vision** , Volume II, pp. 592-602.

# Algorithm

1. **Skin Filter:** The algorithm first locates images containing large areas whose color and texture is appropriate for skin.
2. **Grouping:** Within these areas, the algorithm finds elongated regions and groups them into possible human limbs and connected groups of limbs, using specialized groupers which incorporate substantial amounts of information about object structure.
3. Images containing sufficiently **large skin-colored groups of possible limbs** are reported as potentially containing naked people.

This algorithm was tested on a database of 4854 images: 565 images of naked people and 4289 control images from a variety of sources. The skin filter identified 448 test images and 485 control images as containing substantial areas of skin. Of these, the grouper identified 241 test images and 182 control images as containing people-like shapes.

# Grouping



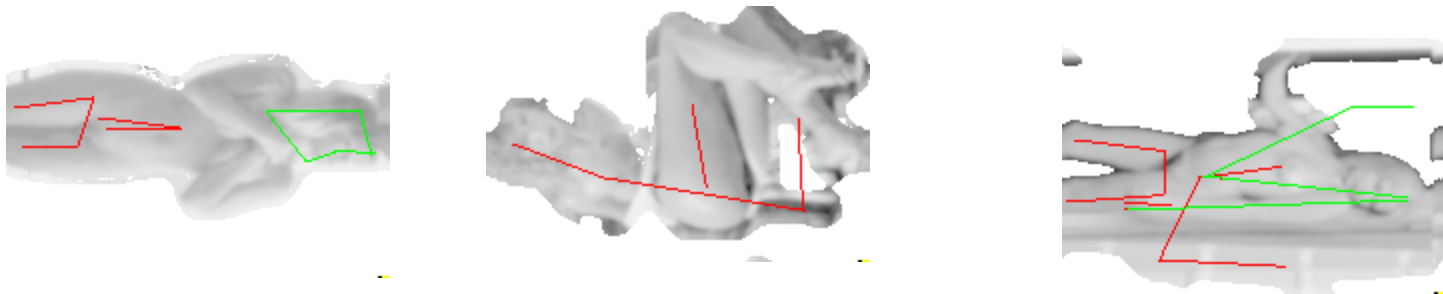
**Fig. 2.** Grouping a spine and two thighs: *Top left* the segment axes that will be grouped into a spine-thigh group, overlaid on the edges, showing the upper bounds on segment length and the their associated symmetries; *Top right* the spine and thigh group assembled from these segments, overlaid on the image.



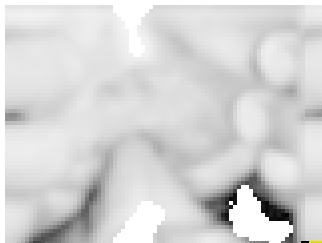
# Results

|                | eliminated by skin filter | eliminated by geometrical analysis | marked as containing naked people |
|----------------|---------------------------|------------------------------------|-----------------------------------|
| test images    | 13.8% (19)                | 34.1% (47)                         | 52.2% (72)                        |
| control images | 92.6% (1297)              | 4.0% (56)                          | 3.4% (48)                         |

**Table 1.** Overall classification performance of the system.



Some True Positives



False Negatives



True Negative

# Object Detection: Rowley's Face Finder

1. convert to gray scale
2. normalize for lighting
3. histogram equalization
4. apply neural net(s)  
trained on 16K images



What data is fed to  
the classifier?

20 x 20 windows in  
a pyramid structure

# Preprocessing

Oval mask for ignoring background pixels:



Original window:



Best fit linear function:



Lighting corrected window:  
(linear function subtracted)



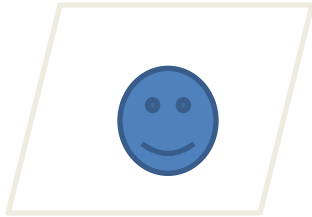
Histogram equalized window:



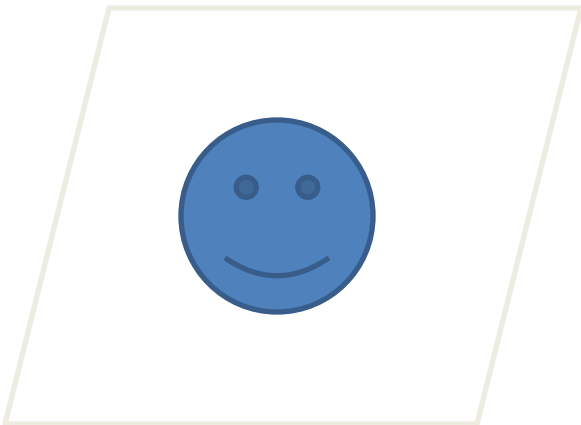
# Image Pyramid Idea



even lower resolution (1/16 of original)



lower resolution image (1/4 of original)



original image (full size)

# Training the Neural Network

## Positive Face Examples

- Nearly 1051 face examples collected from face databases at CMU, Harvard, and WWW
- Faces of various sizes, positions, orientations, intensities
- Eyes, tip of nose, corners and center of mouth labeled manually and used to normalize each face to the same scale, orientation, and position

Result: set of 20 X 20 face training samples

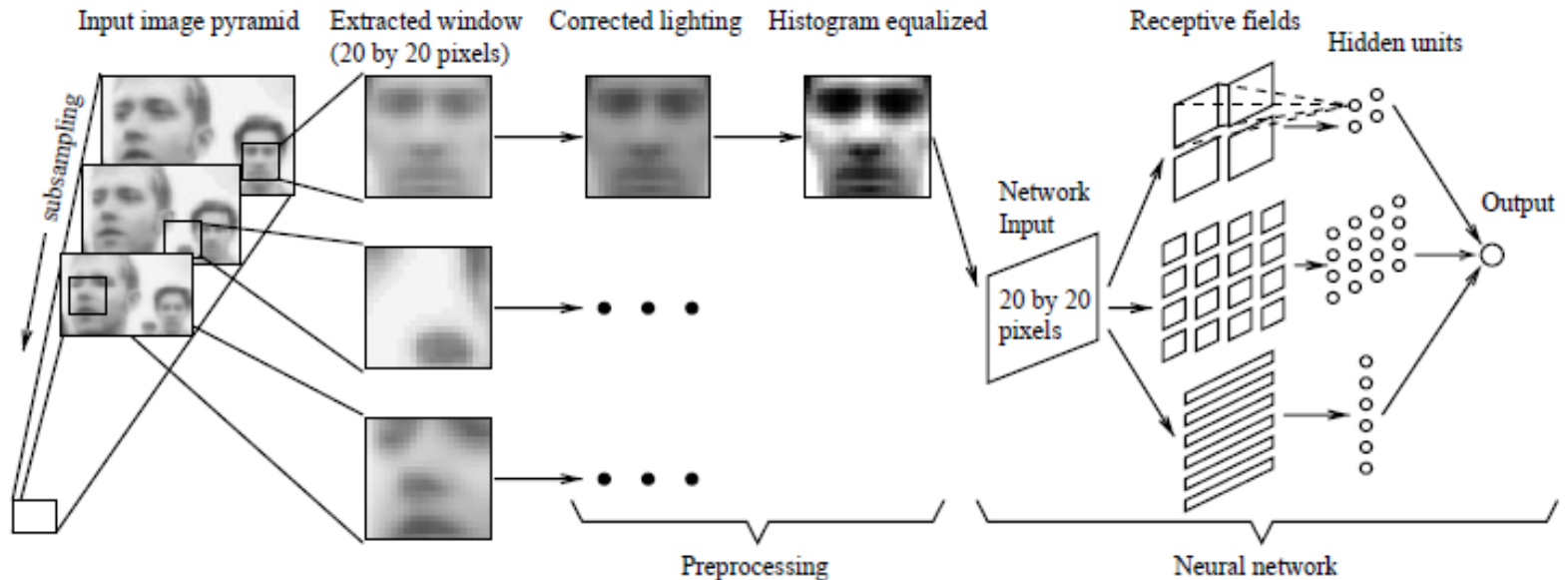
# Training the Neural Network

## Negative Face Examples

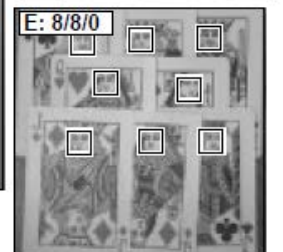
- Generate 1000 random nonface images and apply the preprocessing
- Train a neural network on these plus the face images
- Run the system on real scenes that contain no faces
- Collect the false positives
- Randomly select 250 of these and apply preprocessing
- Label them as negative and add to the training set

# Overall Algorithm

Rowley, Baluja, and Kanade: *Neural Network-Based Face Detection* (PAMI, January 1998) 17

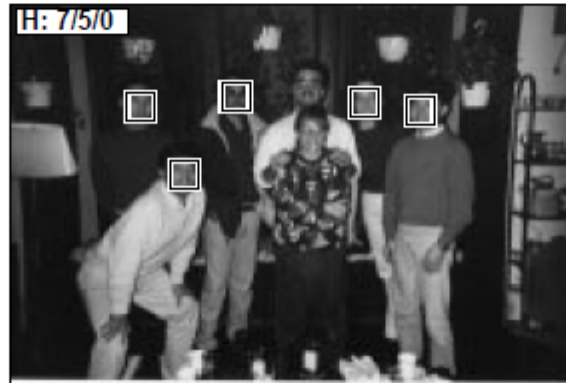
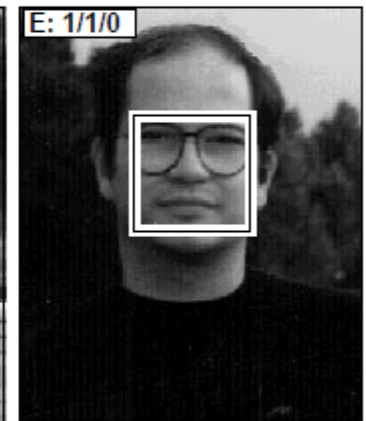
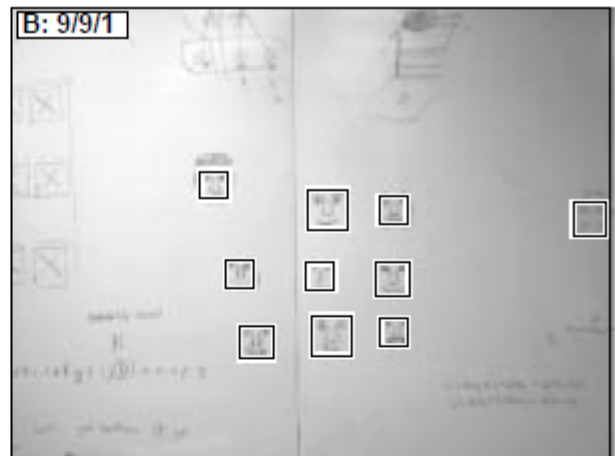


# More Pictures





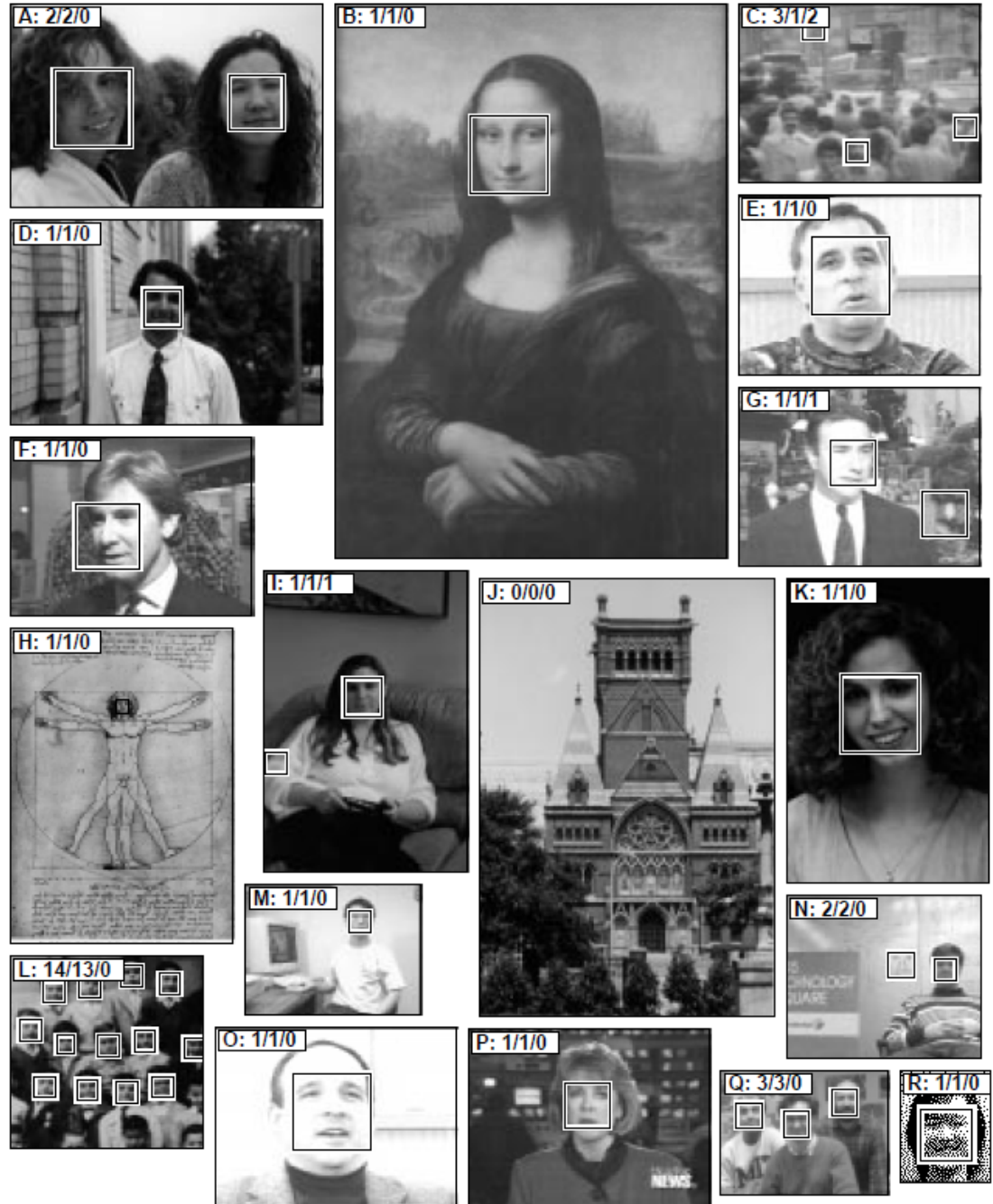
# Even More



# And More

Accuracy: detected 80-90% on different image sets with an “acceptable number” of false positives

Fast Version: 2-4 seconds per image (in 1998)



# Learning from weighted data

| sample  | class | weight |
|---------|-------|--------|
| 1.5 2.6 | I     | 1/4    |
| 2.3 8.9 | II    | 3/4    |

| sample  | class |
|---------|-------|
| 1.5 2.6 | I     |
| 2.3 8.9 | II    |
| 2.3 8.9 | II    |
| 2.3 8.9 | II    |

## Consider a weighted dataset

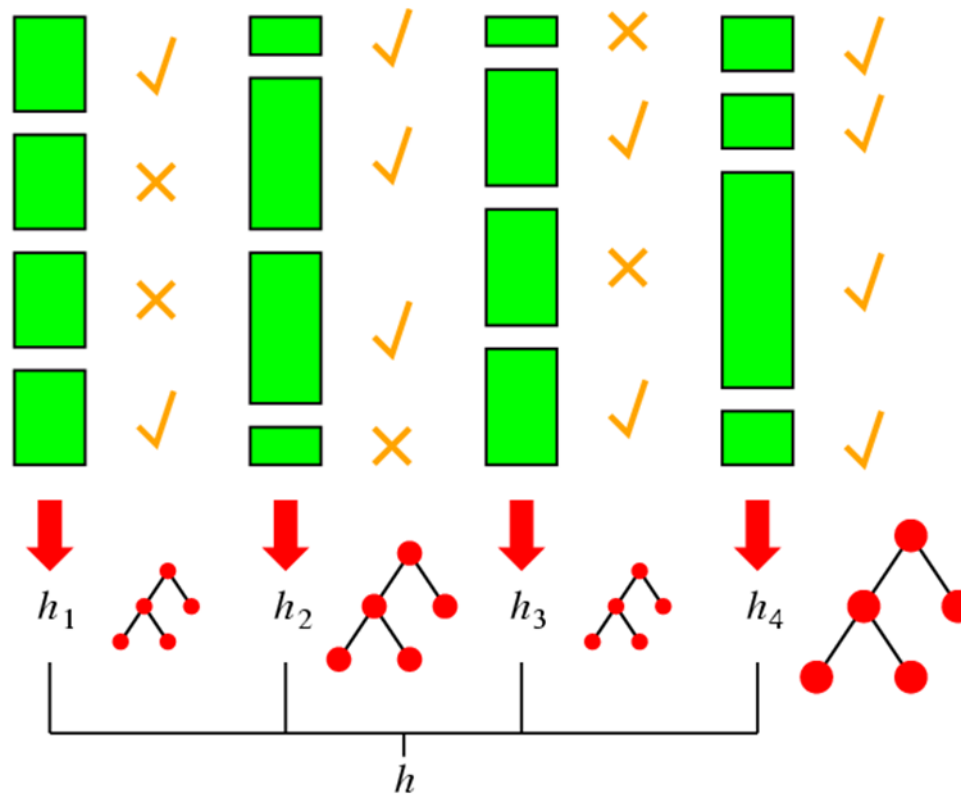
- $W_i$  – **weight** of  $i$ th training example  $(\mathbf{x}^i, y^i)$
- Interpretations:
  - $i$ th training example counts as if it occurred  $W_i$  % times
  - If I were to “resample” data, I would get **more samples** of “heavier” data points

**Now, always do weighted calculations:**

# Basic AdaBoost Overview

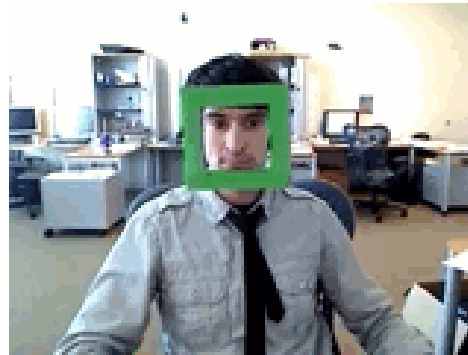
- **Input** is a set of training examples  $(X_i, y_i)$   $i = 1$  to  $m$ .
- We train a **sequence of weak classifiers**, such as decision trees, neural nets or SVMs. Weak because not as strong as the final classifier.
- The training examples will have **weights**, initially all equal.
- At each step, we use the current weights, train a new classifier, and use its performance on the training data to produce **new weights** for the next step (normalized).
- But we **keep ALL** the weak classifiers.
- When it's time for testing on a new feature vector, we will **combine the results from all of the weak classifiers**.

# Idea of Boosting (from AI text)



But Viola Jones will make the weights of the correct samples smaller instead of making the weights of the incorrect ones bigger.

# Face detection



State-of-the-art face detection demo  
(Courtesy [Boris Babenko](#))

# Face detection and recognition



Detection



Recognition

“Sally”



# Face detection

Where are the faces?





# Face Detection

## What kind of features?

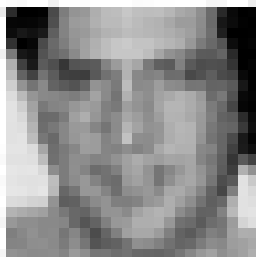
- Rowley: 32 x 32 subimages at different levels of a pyramid
- **Viola/Jones: rectangular features**

## What kind of classifiers?

- Rowley: neural nets
- **Viola/Jones: Adaboost over simple one-node decision trees (stumps)**

# Image Features

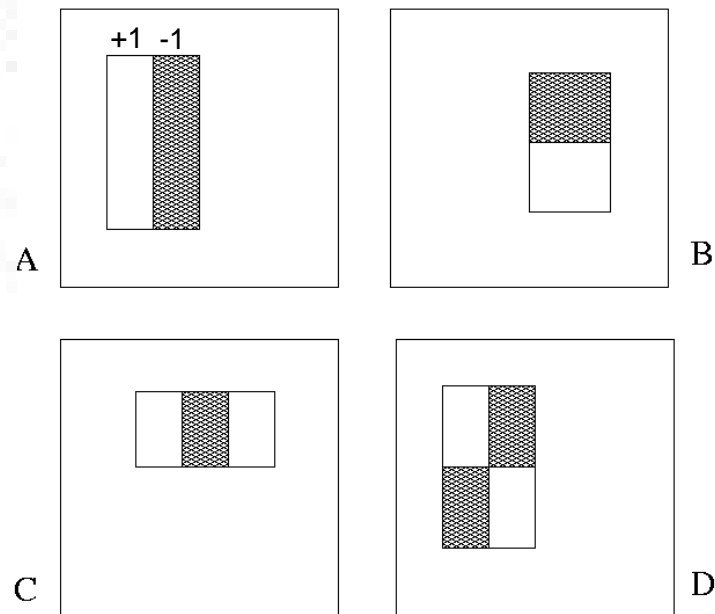
## “Rectangle filters”



People call them Haar-like features, since similar to 2D Haar wavelets.

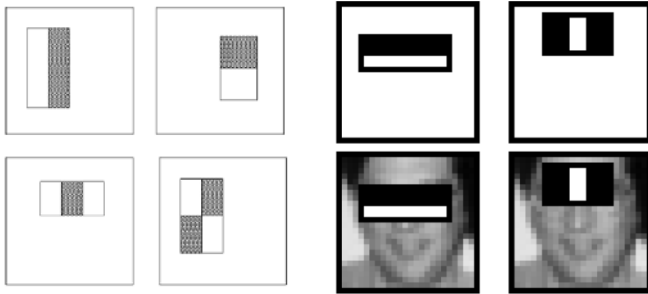
*Value =*

$$\sum (\text{pixels in white area}) - \sum (\text{pixels in black area})$$



# Feature extraction

## “Rectangular” filters



- Filters can be different sizes.
- Filters can be anywhere in the box being analyzed.
- Feature output is very simple convolution.
- Requires sums of large boxes.

Efficiently computable  
with **integral image**: any  
sum can be computed  
in constant time

Avoid scaling images  
scale features directly  
for same cost

# Recall: Sums of rectangular regions

---

How do we compute the sum of the pixels in the red box?

After some pre-computation, this can be done in constant time for any box.

This “trick” is commonly used for computing Haar wavelets (a fundamental building block of many object recognition approaches.)

|     |     |     |     |     |     |     |     |     |     |     |     |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 243 | 239 | 240 | 225 | 206 | 185 | 188 | 218 | 211 | 206 | 216 | 225 |
| 242 | 239 | 218 | 110 | 67  | 31  | 34  | 152 | 213 | 206 | 208 | 221 |
| 243 | 242 | 123 | 58  | 94  | 82  | 132 | 77  | 108 | 208 | 208 | 215 |
| 235 | 217 | 115 | 212 | 243 | 236 | 247 | 139 | 91  | 209 | 208 | 211 |
| 233 | 208 | 131 | 222 | 219 | 226 | 196 | 114 | 74  | 208 | 213 | 214 |
| 232 | 217 | 131 | 116 | 77  | 150 | 69  | 56  | 52  | 201 | 228 | 223 |
| 232 | 232 | 182 | 186 | 184 | 179 | 159 | 123 | 93  | 232 | 235 | 235 |
| 232 | 236 | 201 | 154 | 216 | 133 | 129 | 81  | 175 | 252 | 241 | 240 |
| 235 | 238 | 230 | 128 | 172 | 138 | 65  | 63  | 234 | 249 | 241 | 245 |
| 237 | 236 | 247 | 143 | 59  | 78  | 10  | 94  | 255 | 248 | 247 | 251 |
| 234 | 237 | 245 | 193 | 55  | 33  | 115 | 144 | 213 | 255 | 253 | 251 |
| 248 | 245 | 161 | 128 | 149 | 109 | 138 | 65  | 47  | 156 | 239 | 255 |
| 190 | 107 | 39  | 102 | 94  | 73  | 114 | 58  | 17  | 7   | 51  | 137 |
| 23  | 32  | 33  | 148 | 168 | 203 | 179 | 43  | 27  | 17  | 12  | 8   |
| 17  | 26  | 12  | 160 | 255 | 255 | 109 | 22  | 26  | 19  | 35  | 24  |

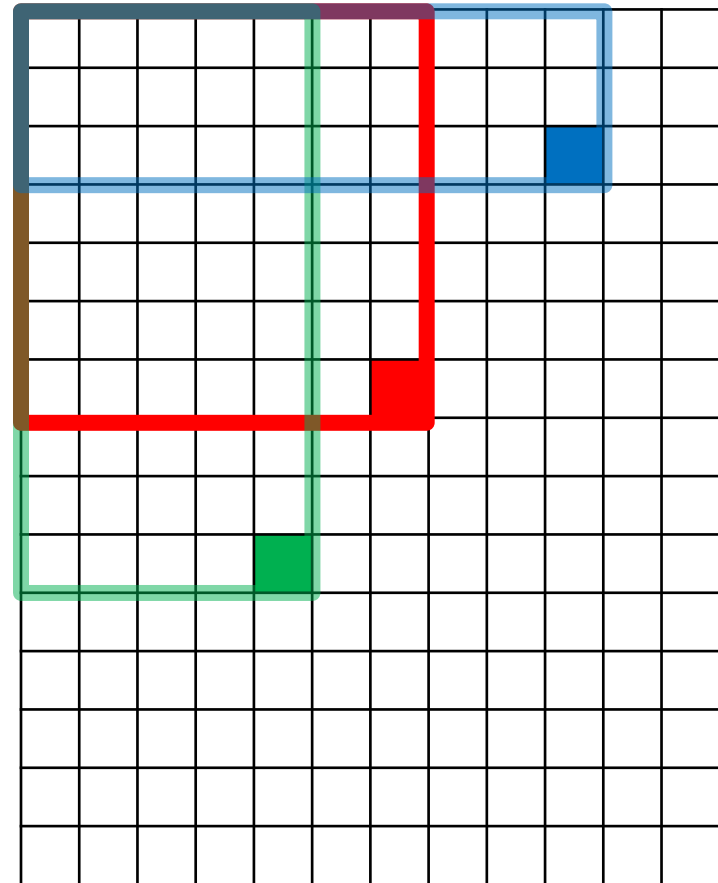
# Sums of rectangular regions

---

The trick is to compute an “integral image.” Every pixel is the sum of its neighbors to the upper left.

Sequentially compute using:

$$I(x, y) = I(x, y) + \\ I(x - 1, y) + I(x, y - 1) - \\ I(x - 1, y - 1)$$



# Sums of rectangular regions

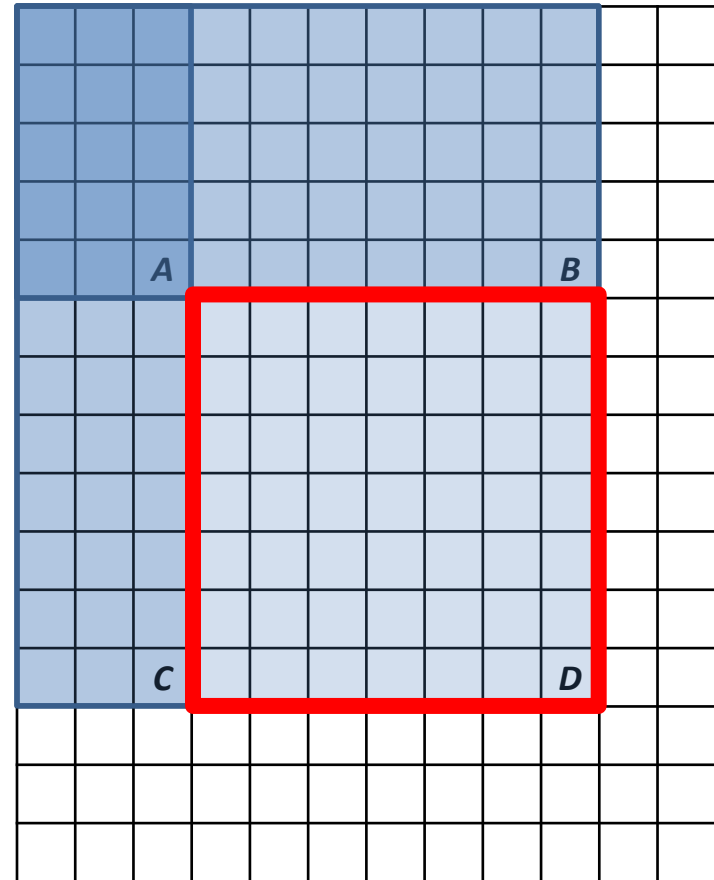
---

Solution is found using:

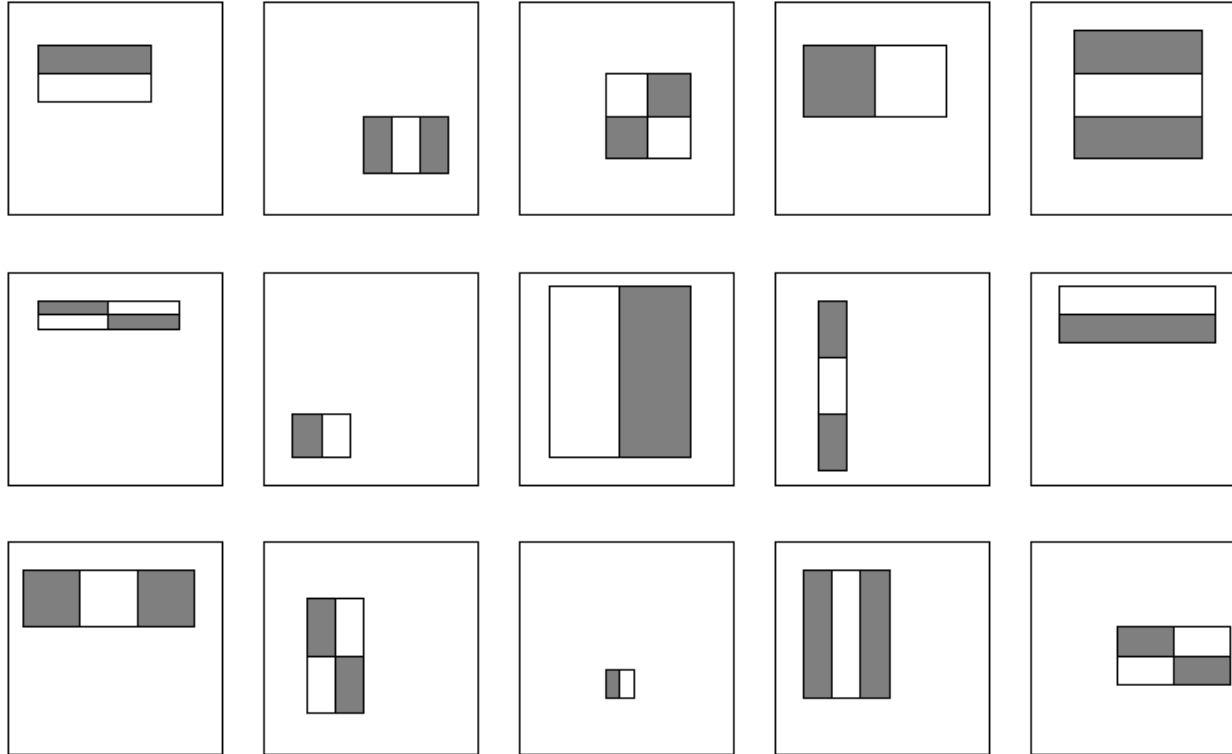
$$A + D - B - C$$

What if the position of the box lies  
between pixels?

Use bilinear interpolation.



# Large library of filters



Considering all possible filter parameters:  
position, scale,  
and type:

**160,000+**  
possible features  
associated with  
each 24 x 24  
window

Use **AdaBoost** both to select the informative features and to form the classifier

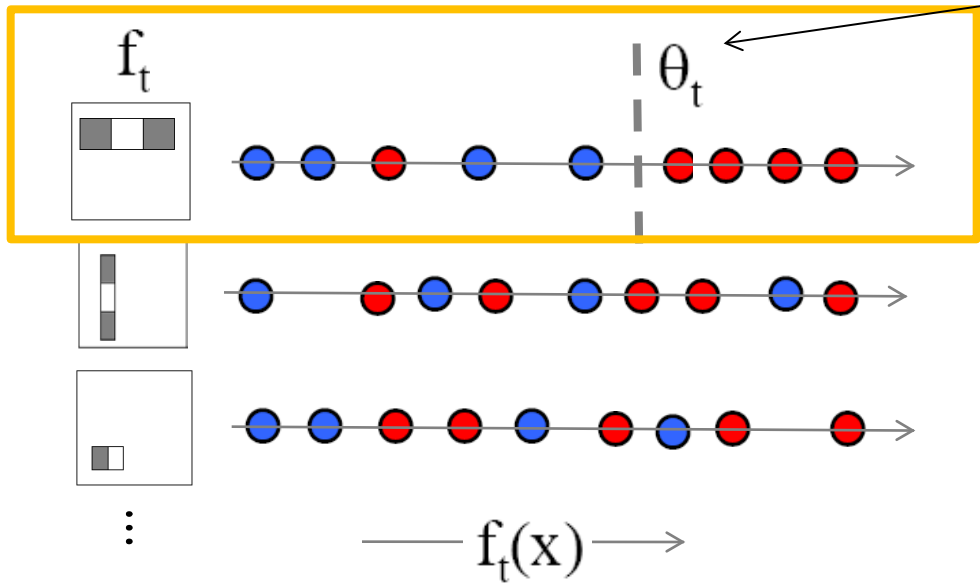
# Feature selection

- For a 24x24 detection region, the number of possible rectangle features is ~160,000!
- At test time, it is impractical to evaluate the entire feature set
- Can we create a good classifier using just a small subset of all possible features?
- How to select such a subset?



# AdaBoost for feature+classifier selection

- Want to select the single rectangle feature and threshold that best separates **positive** (faces) and **negative** (non-faces) training examples, in terms of *weighted* error.



$\theta_t$  is a threshold for classifier  $h_t$   
 Resulting **weak classifier**:

$$h_t(x) = \begin{cases} +1 & \text{if } f_t(x) > \theta_t \\ 0 & \text{otherwise} \end{cases}$$

For next round, reweight the examples according to errors, choose another filter/threshold combo.

Outputs of a possible rectangle feature on training examples x (faces and non faces)

# Weak Classifiers

- Each weak classifier works on exactly **one rectangle feature**.
- Each weak classifier has 3 associated variables
  1. its threshold  $\theta$
  2. its polarity  $p$
  3. its weight  $\alpha$
- The polarity can be 0 or 1 (in our code)
- The weak classifier computes its one feature  $f$ 
  - When the polarity is 1, we want  $f > \theta$  for face
  - When the polarity is 0, we want  $f < \theta$  for face
- The weight will be used in the final classification by AdaBoost.

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

$\beta_t = \varepsilon_t / (1 - \varepsilon_t)$ : the training error of the classifier  $h_t$

- Final classifier is combination of the weak ones, weighted according to error they had.
- The code computes a **SCORE** based on the difference of the two above summations.

# AdaBoost Algorithm modified by Viola Jones

- 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,

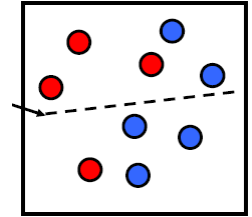
$$w_{t,i} \leftarrow \frac{w_{t,i}}{\sum_{j=1}^n w_{t,j}}$$

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  $w_t$ ,  $\epsilon_j = \sum_i w_i |h_j(x_i) - y_i|$ . **sum over training samples**
3. Choose the classifier,  $h_t$ , with the lowest error  $\epsilon_t$ .
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  $\beta_t = \frac{\epsilon_t}{1-\epsilon_t}$ .



**NOTE:** Our code uses equal weights for all samples

$\{x_1, \dots, x_n\}$

For  $T$  rounds: meaning we will construct  $T$  weak classifiers

← Normalize weights

← Find the best threshold and polarity for each feature, and return error  $\epsilon$ .

← Re-weight the examples:  
Incorrectly classified -> more weight  
Correctly classified -> less weight

# Updating the Weights

Suppose a weak classifier  $i$  has error  $e_i$ .

The weight alpha for this classifier is

$$\alpha = \ln((1-e_i)/e_i)$$

The updating formula for the weight  $w_i$  for classifier  $i$  is given as

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

where  $e_i = 0$  if example  $x_i$  is classified correctly else 1.

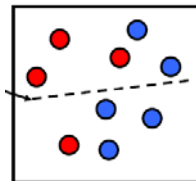
And  $\beta_t = \exp(-\alpha_t)$  which is  $e_i/(1-e_i)$

After the weights are updated, they are normalized by the sum of all of them.

# Recall

- Classification
  - Decision Trees and Forests
  - Neural Nets
  - SVMs
  - Boosting
  - ....

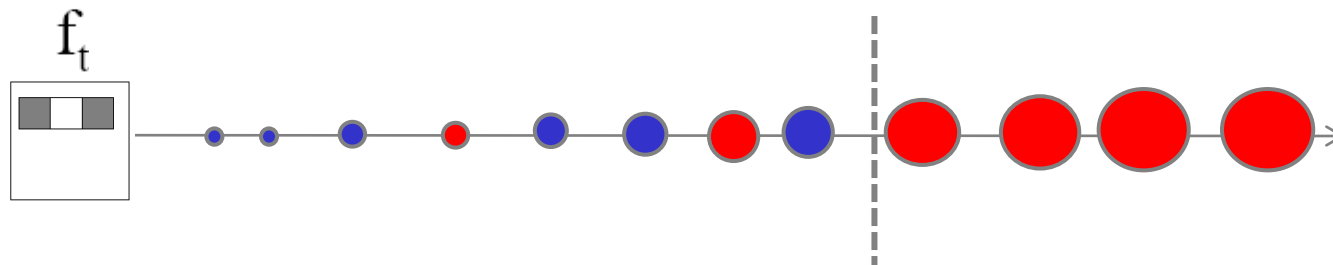
- Face Detection
  - Simple Features
  - Integral Images
  - Boosting



# Picking the threshold for the best classifier: the idea

Efficient single pass approach:

The features are actually **sorted** in the code according to numeric value!



At each sample compute:

$$e = \min ( S + (T - S), S + (T - S) )$$

Find the minimum value of  $e$ , and use the value of the corresponding sample as the threshold.

$S$  = sum of **weights of** samples with feature value below the current sample

$T$  = total sum of all samples

$S$  and  $T$  are for faces;  $S$  and  $T$  are for background.

# Picking the threshold for the best classifier: the details for coding

The **features** for the training samples are actually **sorted** in the code according to numeric value!

Algorithm:

1. find **AFS**, the sum of the weights of all the **face** samples
2. find **ABG**, the sum of the weights of all the **background** samples
3. set to zero **FS**, the sum of the weights of face samples so far
4. set to zero **BG**, the sum of the weights of background samples so far
5. go through each sample *s* in a loop **IN THE SORTED ORDER**

At each sample, add weight to FS or BG and compute:

$$e = \min (\text{BG} + (\text{AFS} - \text{FS}), \text{FS} + (\text{ABG} - \text{BG}))$$

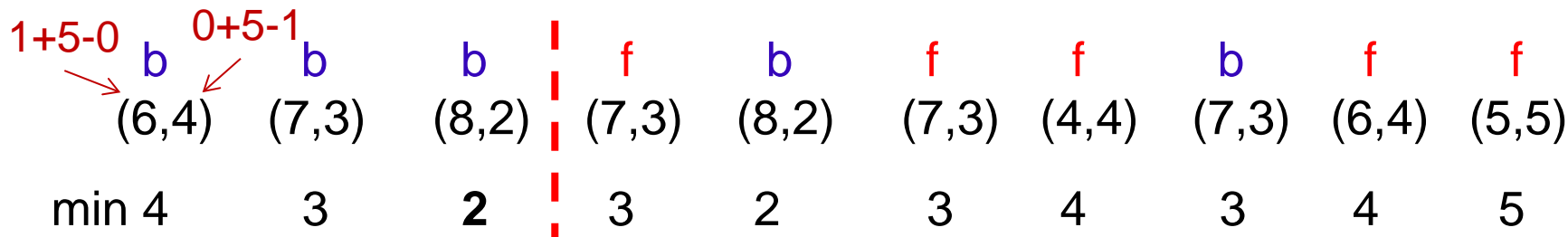
Find the minimum value of  $e$ , and use the feature value of the corresponding sample as the threshold.



# What's going on?

$$\text{error} = \min \left( \frac{\text{BG} + (\text{AFS} - \text{FS})}{\text{left}}, \frac{\text{FS} + (\text{ABG} - \text{BG})}{\text{right}} \right)$$

- Let's pretend the weights on the samples are all 1's.
- The samples are arranged in ascending order by feature value and we know which ones are faces (f) and background (b).
- **Left** is the number of background patches so far plus the number of faces **yet to be encountered**.
- **Right** is the number of faces so far plus the number of background patches **yet to be encountered**.



# Setting the Polarity

$$\text{error} = \min \left( \underset{\text{left}}{\text{BG} + (\text{AFS} - \text{FS})}, \underset{\text{right}}{\text{FS} + (\text{ABG} - \text{BG})} \right)$$

- **left** is the number of background patches so far plus the number of faces yet to be encountered.
  - **right** is the number of faces so far plus the number of background patches yet to be encountered.
- 
- When  $\text{left} < \text{right}$ , set polarity to 0
  - Else set polarity to 1

# Threshold and Polarity Example

$$\text{error} = \min (\text{BG} + (\text{AFS} - \text{FS}), \text{FS} + (\text{ABG} - \text{BG}))$$

|          |     |     |     |     |     |                  |
|----------|-----|-----|-----|-----|-----|------------------|
| samples  | 0   | 1   | 2   | 3   | 4   |                  |
| labels   | F   | B   | F   | B   | B   | initialize       |
| features | 6   | 3   | 10  | 2   | 1   | AFS = 0          |
| weight   | 1/5 | 1/5 | 1/5 | 1/5 | 1/5 | ABG = 0          |
| index    | 4   | 3   | 1   | 0   | 2   | besterr = 999999 |

AFS becomes sum of face sample weights = 2/5; ABG = 3/5

**step 0:** idx = 4; FS stays 0; BG = 1/5  
 error =  $\min(1/5 + (2/5 - 0), 0 + (3/5 - 1/5)) = 2/5$   
 besterr = 2/5; bestpolarity = 1; bestthreshold=1

**step 1:** idx = 3; FS stays 0; BG = 2/5  
 error =  $\min(2/5 + (2/5 - 0), 0 + (3/5 - 2/5)) = 1/5$   
 besterr = 1/5; bestpolarity = 1; bestthreshold=2

# Threshold and Polarity Example

$$\text{error} = \min (\text{BG} + (\text{AFS} - \text{FS}), \text{FS} + (\text{ABG} - \text{BG}))$$

|          |     |     |     |     |     |
|----------|-----|-----|-----|-----|-----|
| samples  | 0   | 1   | 2   | 3   | 4   |
| labels   | F   | B   | F   | B   | B   |
| features | 6   | 3   | 10  | 2   | 1   |
| weight   | 1/5 | 1/5 | 1/5 | 1/5 | 1/5 |
| index    | 4   | 3   | 1   | 0   | 2   |

initialize

AFS = 0

ABG = 0

besterr = 999999

**step 2:** idx = 1; FS stays 0; BG = 3/5  
error =  $\min(3/5 + (2/5 - 0), 0 + (3/5 - 3/5)) = 0$   
besterr = 0; bestpolarity = 1; bestthreshold=3

**step 3:** idx = 0; FS = 1/5; BG = 3/5  
error =  $\min(3/5 + (2/5 - 1/5), 1/5 + (3/5 - 3/5)) = 1/5$   
NO CHANGE

**step 4:** idx = 2; FS = 2/5; BG = 3/5  
error =  $\min(3/5 + (2/5 - 2/5), 2/5 + (3/5 - 3/5)) = 2/5$   
NO CHANGE

**RESULT**

|   |   |   |  |   |    |
|---|---|---|--|---|----|
| 1 | 2 | 3 |  | 6 | 10 |
|---|---|---|--|---|----|

$\theta > 3$

# Measuring classification performance

- Confusion matrix
- Accuracy
  - $(TP+TN)/(TP+TN+FP+FN)$
- True Positive Rate=Recall
  - $TP/(TP+FN)$
- False Positive Rate
  - $FP/(FP+TN)$
- Precision
  - $TP/(TP+FP)$
- F1 Score
  - $2*Recall*Precision/(Recall+Precision)$

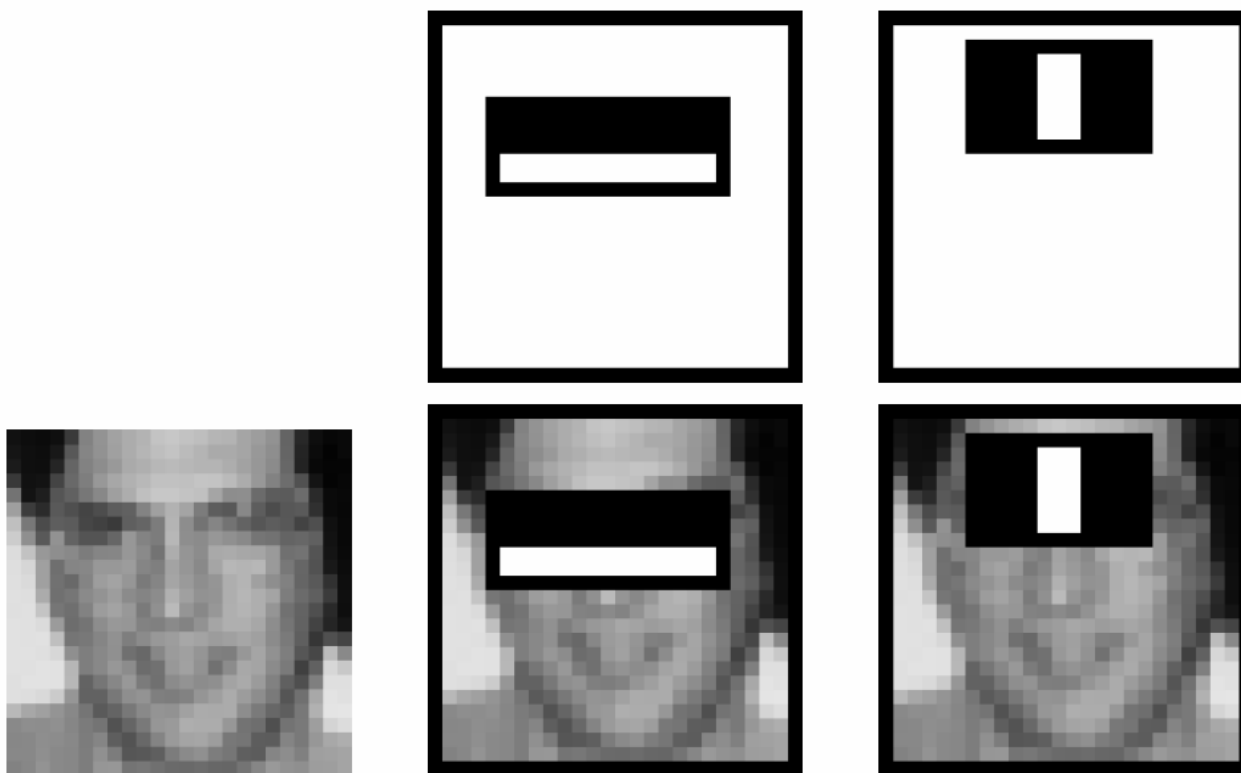
|              |        | Predicted class |        |        |
|--------------|--------|-----------------|--------|--------|
|              |        | Class1          | Class2 | Class3 |
| Actual class | Class1 | 40              | 1      | 6      |
|              | Class2 | 3               | 25     | 7      |
|              | Class3 | 4               | 9      | 10     |

|        |          | Predicted      |                |
|--------|----------|----------------|----------------|
|        |          | Positive       | Negative       |
| Actual | Positive | True Positive  | False Negative |
|        | Negative | False Positive | True Negative  |

# Boosting for face detection

---

- First two features selected by boosting:

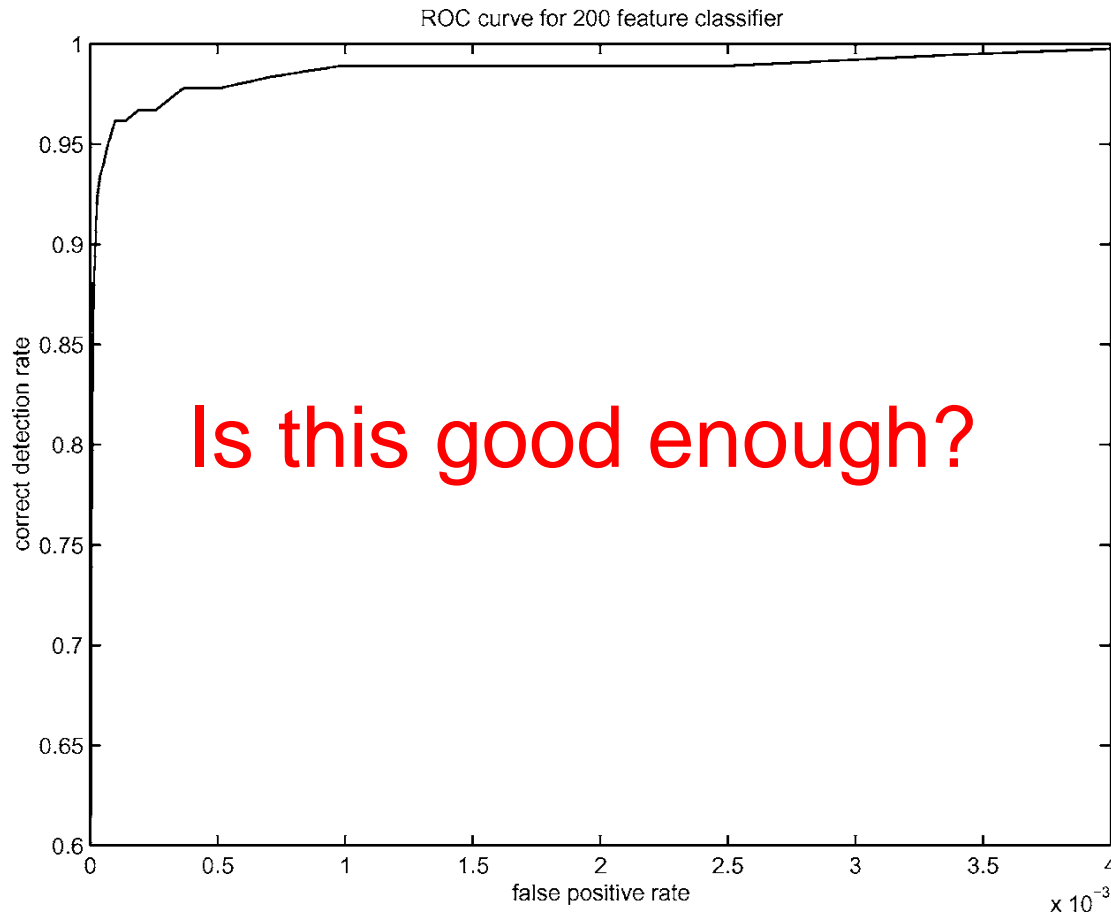


This feature combination can yield 100% detection rate and 50% false positive rate

# Boosting for face detection

---

- A 200-feature classifier can yield 95% detection rate and a false positive rate of 1 in 14084

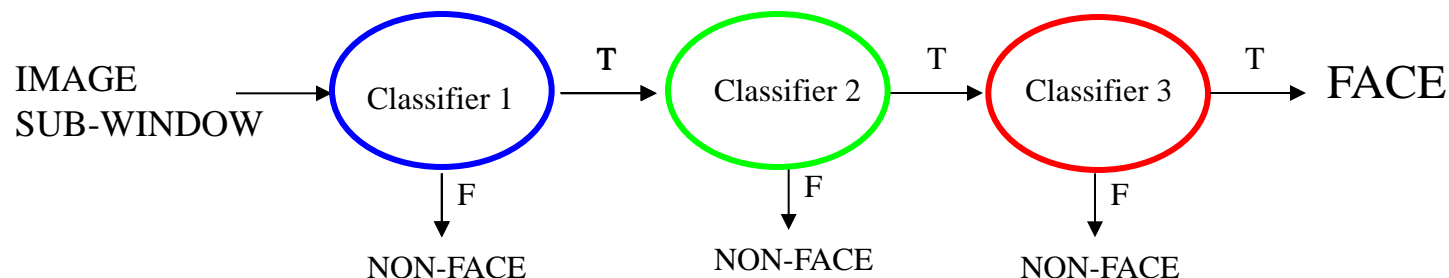


Receiver operating characteristic (ROC) curve

# Attentional cascade (from Viola-Jones)

---

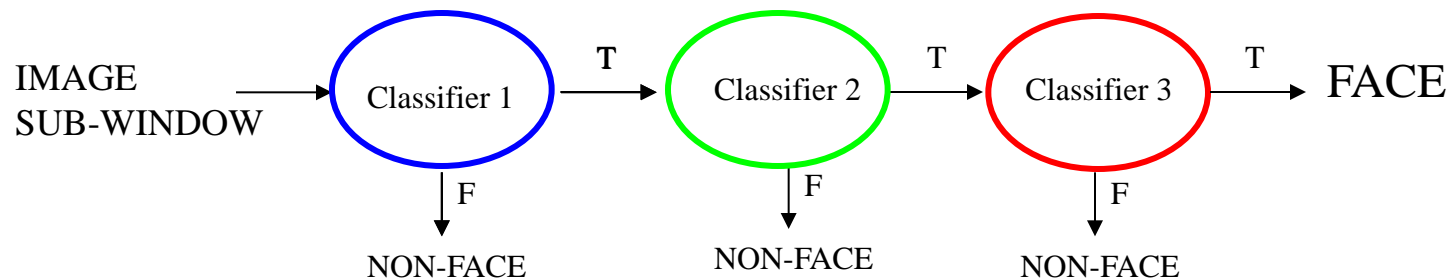
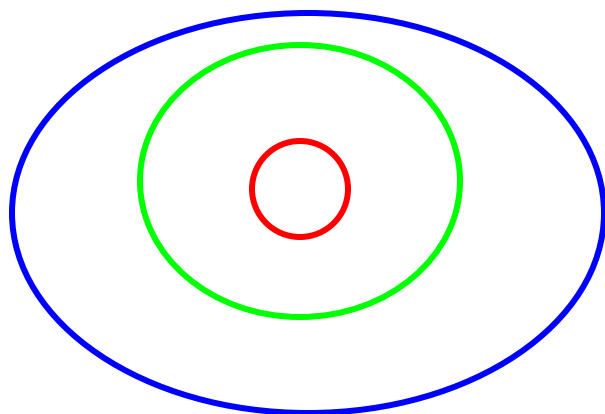
- We start with **simple classifiers** which reject many of the negative sub-windows while detecting almost all positive sub-windows
- Positive response from the first classifier triggers the evaluation of a second (more complex) classifier, and so on
- **A negative outcome at any point leads to the immediate rejection of the sub-window**



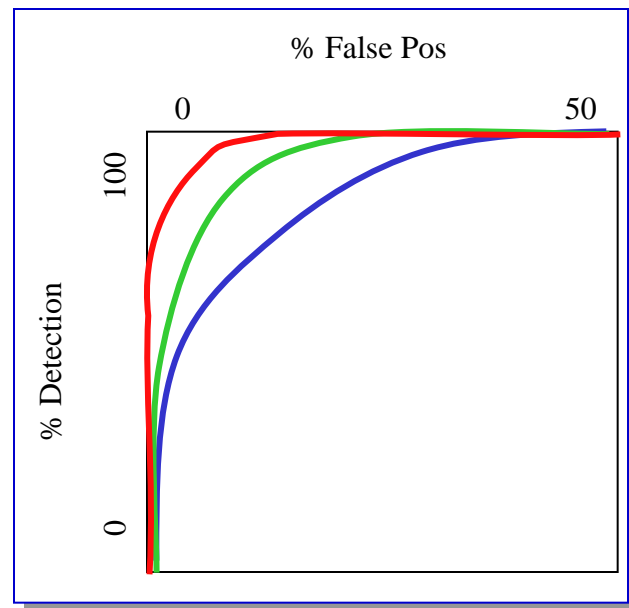


# Attentional cascade

- Chain of classifiers that are progressively more complex and have lower false positive rates:



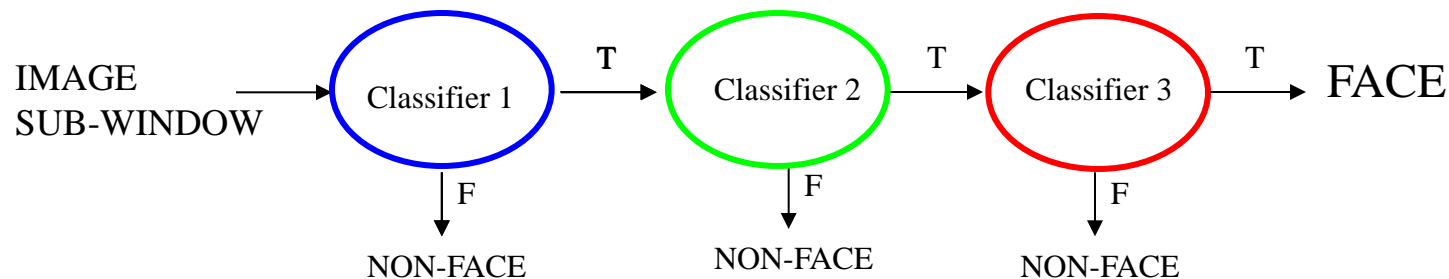
Receiver operating characteristic



# Attentional cascade

---

- The detection rate and the false positive rate of the cascade are found by multiplying the respective rates of the individual stages
- A detection rate of 0.9 and a false positive rate on the order of  $10^{-6}$  can be achieved by a 10-stage cascade if each stage has a detection rate of 0.99 ( $0.99^{10} \approx 0.9$ ) and a false positive rate of about 0.30 ( $0.3^{10} \approx 6 \times 10^{-6}$ )



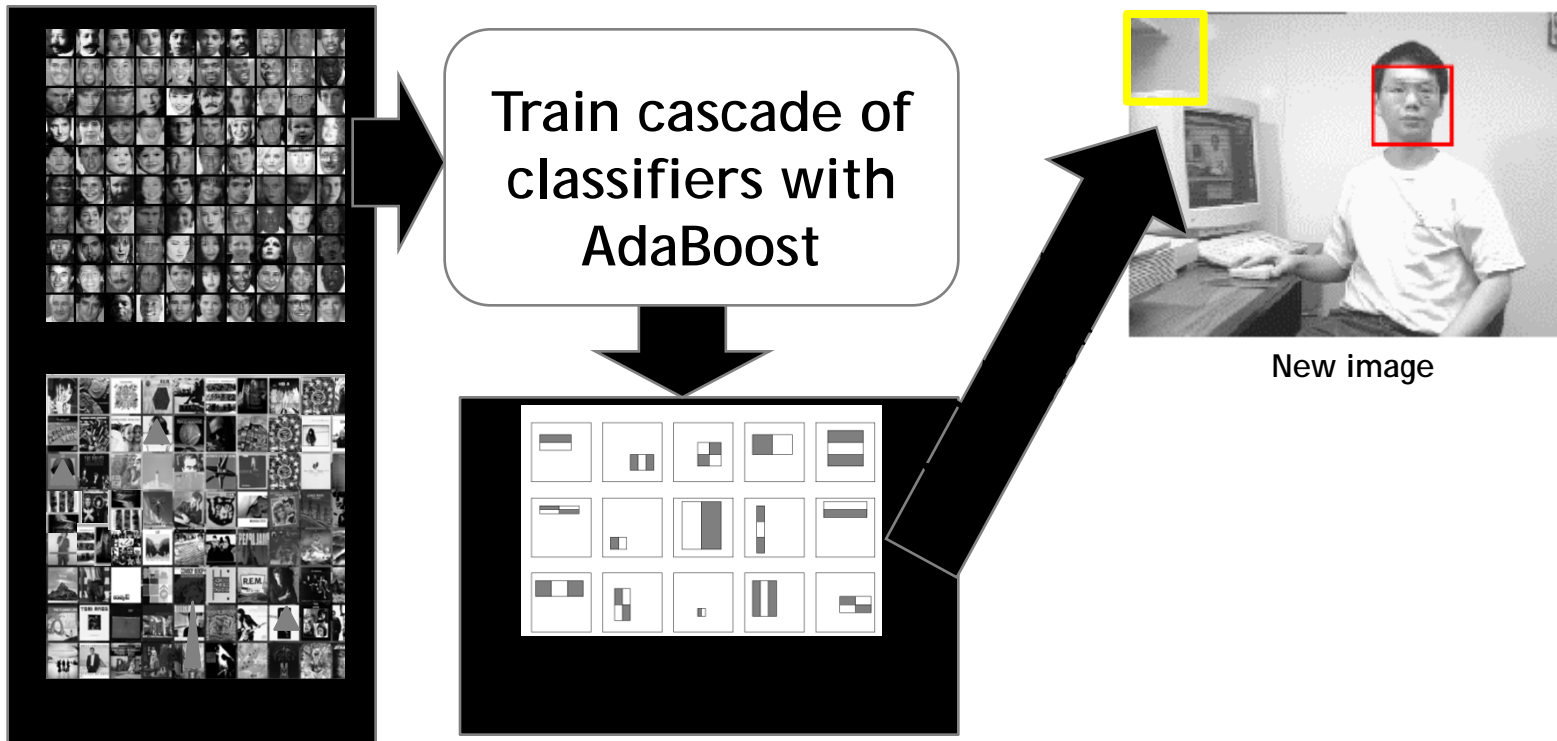
# Training the cascade

---

- Set target detection and false positive rates for each stage
- Keep adding features to the current stage until its target rates have been met
  - Need to lower AdaBoost threshold to maximize detection (as opposed to minimizing total classification error)
  - Test on a *validation set*
- If the overall false positive rate is not low enough, then add another stage
- Use false positives from current stage as the negative training examples for the next stage

# 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/>]

# The implemented system

---

- Training Data
  - 5000 faces
    - All frontal, rescaled to 24x24 pixels
  - 300 million non-faces
    - 9500 non-face images
  - Faces are normalized
    - Scale, translation
- Many variations
  - Across individuals
  - Illumination
  - Pose



# System performance

---

- Training time: “weeks” on 466 MHz Sun workstation
- 38 layers, total of 6061 features
- Average of 10 features evaluated per window on test set
- “On a 700 Mhz Pentium III processor, the face detector can process a 384 by 288 pixel image in about .067 seconds”
  - 15 Hz
  - 15 times faster than previous detector of comparable accuracy (Rowley et al., 1998)

# Non-maximal suppression (NMS)

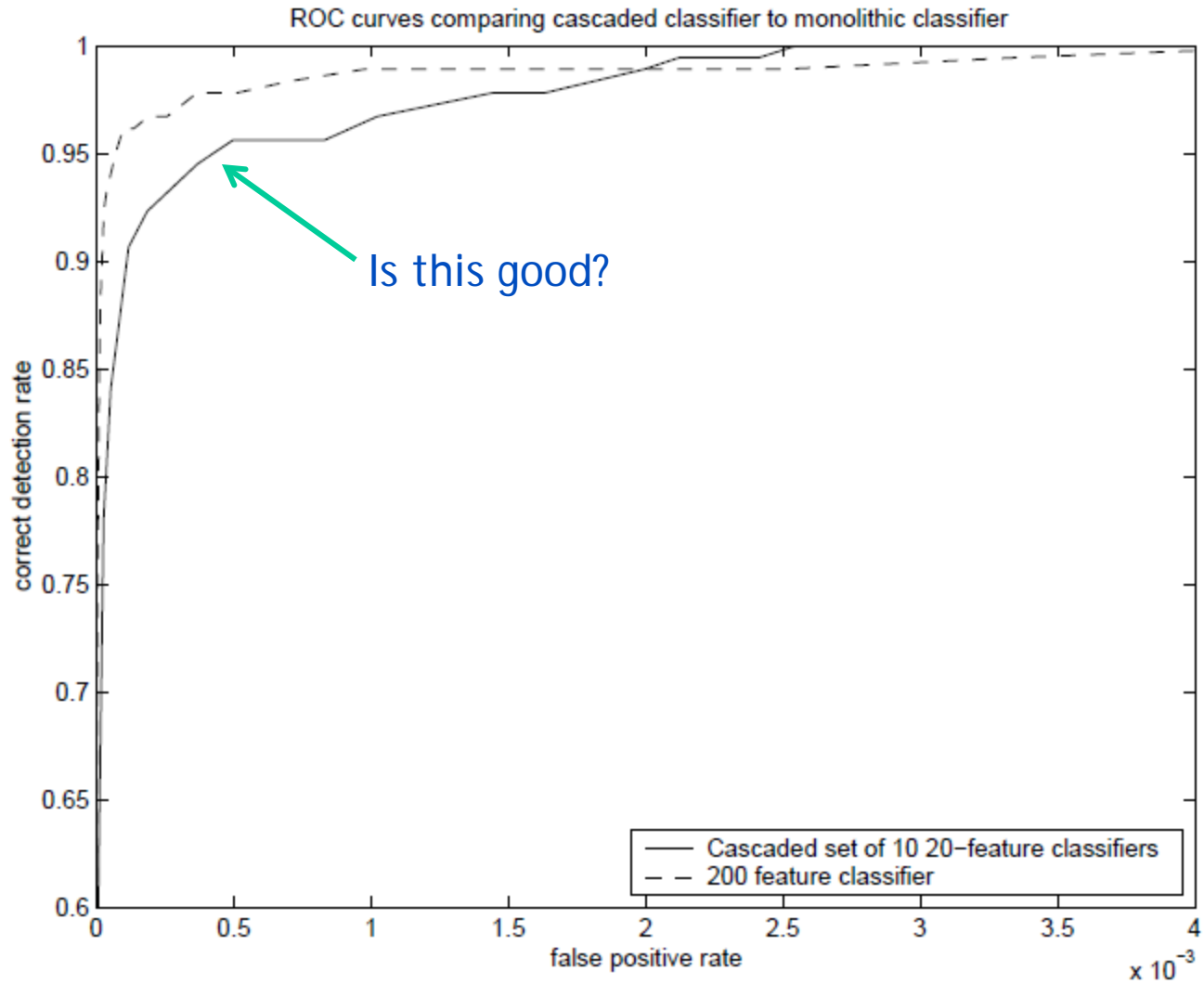


Many detections above threshold.

# Non-maximal suppression (NMS)

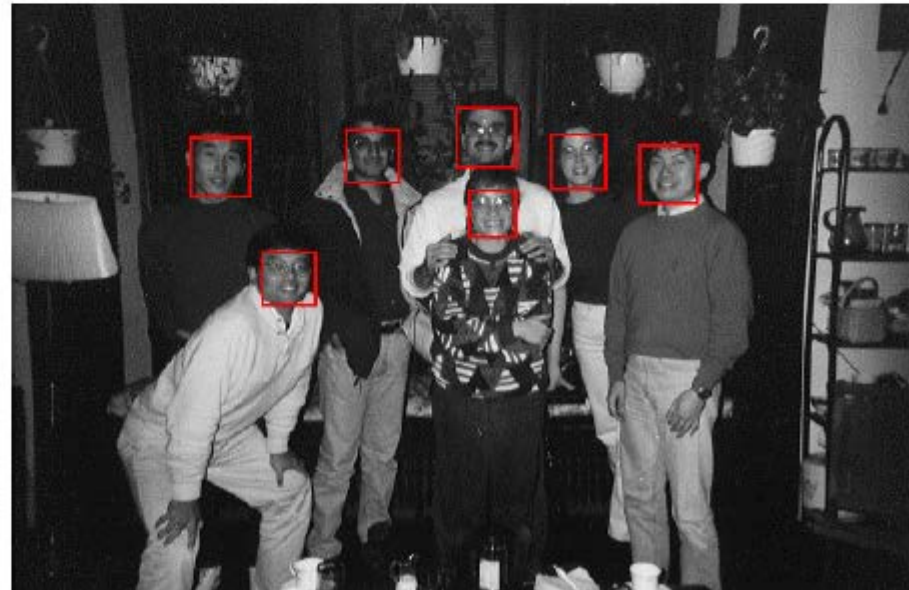
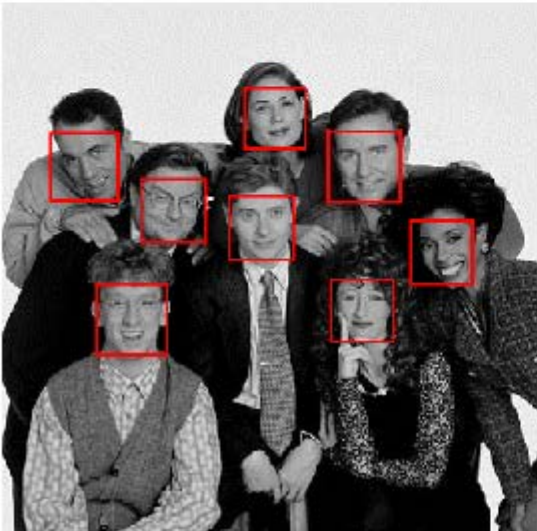
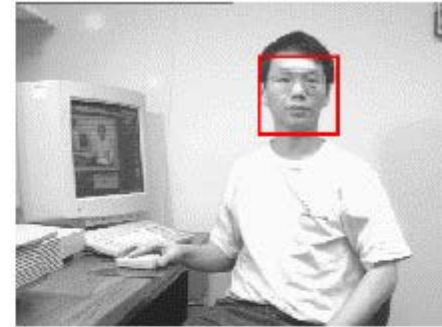
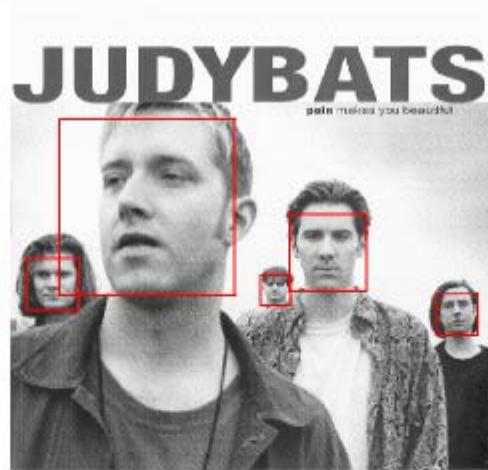




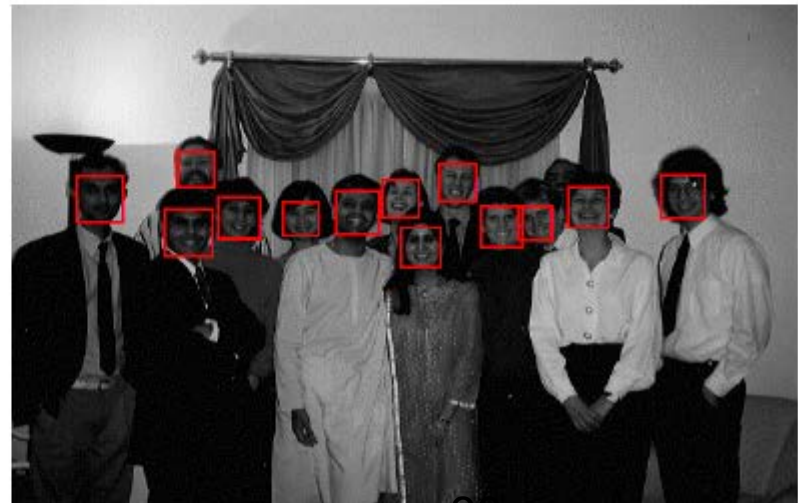
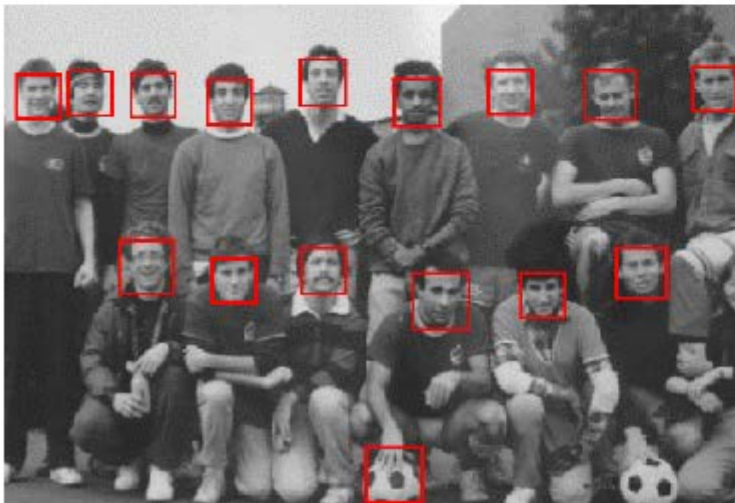
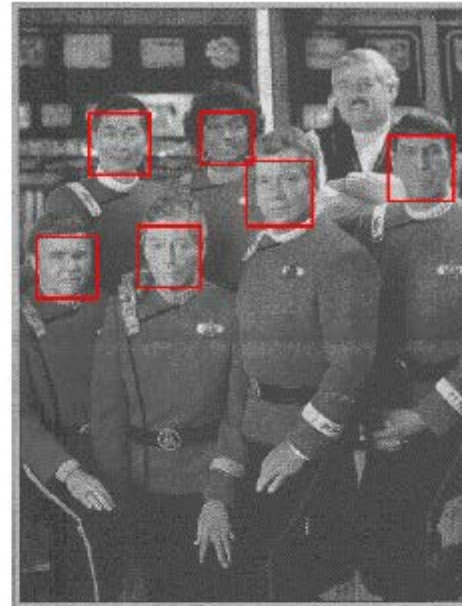
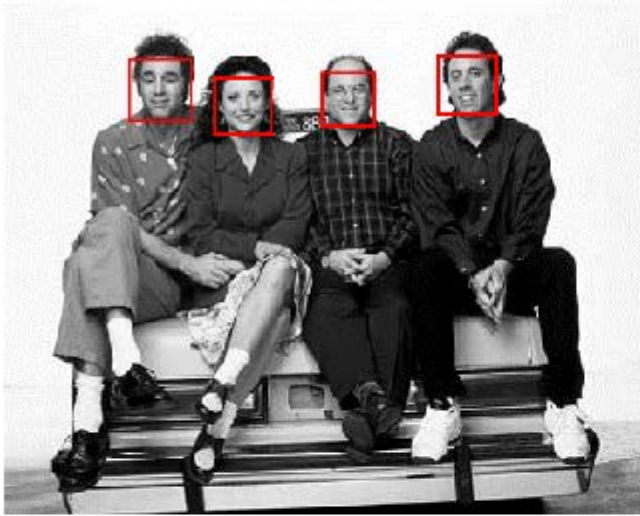


Similar accuracy, but 10x faster

# Viola-Jones Face Detector: Results



# Viola-Jones Face Detector: Results

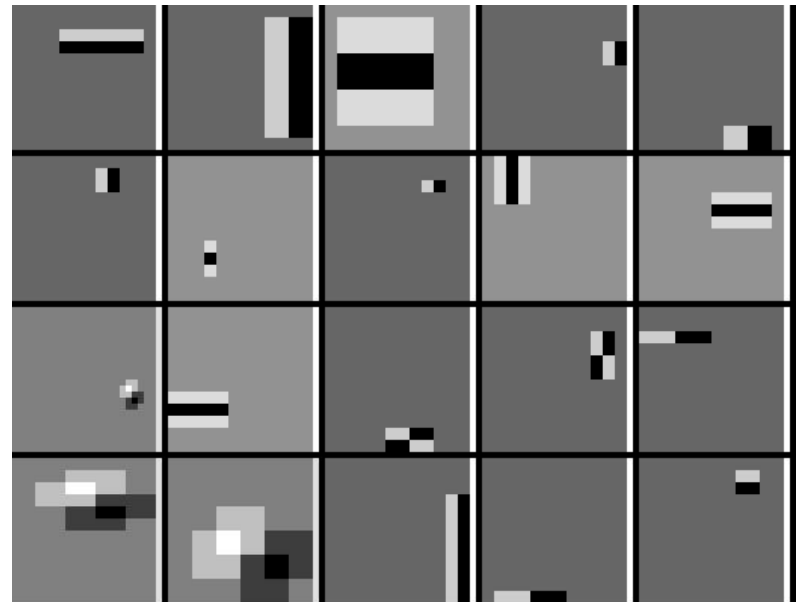


# Viola-Jones Face Detector: Results



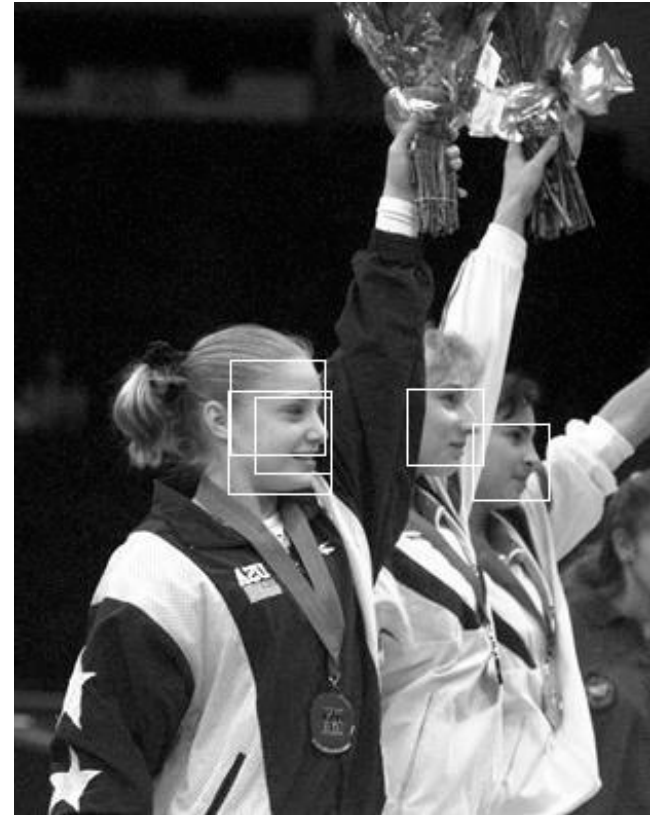
# Detecting profile faces?

Detecting profile faces requires training separate detector with profile examples.





# Viola-Jones Face Detector: Results



# Summary: Viola/Jones detector

- Rectangle features
- Integral images for fast computation
- Boosting for feature selection
- Attentional cascade for fast rejection of negative windows