# Robust Real-Time Face Detection

PAUL VIOLA
*Microsoft Research, One Microsoft Way, Redmond, WA 98052, USA*
viola@microsoft.com

MICHAEL J. JONES
*Mitsubishi Electric Research Laboratory, 201 Broadway, Cambridge, MA 02139, USA*
mjones@merl.com

**Abstract.**   This paper describes a face detection framework that is capable of processing images extremely rapidly while achieving high detection rates. There are three key contributions. The first is the introduction of a new image representation called the "Integral Image" which allows the features used by our detector to be computed very quickly. The second is a simple and efficient classifier which is built using the AdaBoost learning algorithm (Freund and Schapire, 1995) to select a small number of critical visual features from a very large set of potential features. The third contribution is a method for combining classifiers in a "cascade" which allows background regions of the image to be quickly discarded while spending more computation on promising face-like regions. A set of experiments in the domain of face detection is presented. The system yields face detection performance comparable to the best previous systems (Sung and Poggio, 1998; Rowley et al., 1998; Schneiderman and Kanade, 2000; Roth et al., 2000). Implemented on a conventional desktop, face detection proceeds at 15 frames per second.

## 1.   Introduction

This paper brings together new algorithms and insights to construct a framework for robust and extremely rapid visual detection. Toward this end we have constructed a frontal face detection system which achieves detection and false positive rates which are equivalent to the best published results (Sung and Poggio, 1998; Rowley et al., 1998; Osuna et al., 1997a; Schneiderman and Kanade, 2000; Roth et al., 2000). This face detection system is most clearly distinguished from previous approaches in its ability to detect faces extremely rapidly. Operating on 384 by 288 pixel images, faces are detected at 15 frames per second on a conventional 700 MHz Intel Pentium III. In other face detection systems, auxiliary information, such as image differ-ences in video sequences, or pixel color in color images, have been used to achieve high frame rates. Our system achieves high frame rates working only with the information present in a single grey scale image. These alternative sources of information can also be integrated with our system to achieve even higher frame rates.

There are three main contributions of our face detection framework. We will introduce each of these ideas briefly below and then describe them in detail in subsequent sections.

The first contribution of this paper is a new image representation called an *integral image* that allows for very fast feature evaluation. Motivated in part by the work of Papageorgiou et al. (1998) our detection system does not work directly with image intensities. Like

138    *Viola and Jones*

these authors we use a set of features which are reminiscent of Haar Basis functions (though we will also use related filters which are more complex than Haar filters). In order to compute these features very rapidly at many scales we introduce the integral image representation for images (the integral image is very similar to the summed area table used in computer graphics (Crow, 1984) for texture mapping). The integral image can be computed from an image using a few operations per pixel. Once computed, any one of these Haar-like features can be computed at any scale or location in *constant* time.

The second contribution of this paper is a simple and efficient classifier that is built by selecting a small number of important features from a huge library of potential features using AdaBoost (Freund and Schapire, 1995). Within any image sub-window the total number of Haar-like features is very large, far larger than the number of pixels. In order to ensure fast classification, the learning process must exclude a large majority of the available features, and focus on a small set of critical features. Motivated by the work of Tieu and Viola (2000) feature selection is achieved using the AdaBoost learning algorithm by constraining each weak classifier to depend on only a single feature. As a result each stage of the boosting process, which selects a new weak classifier, can be viewed as a feature selection process. AdaBoost provides an effective learning algorithm and strong bounds on generalization performance (Schapire et al., 1998).

The third major contribution of this paper is a method for combining successively more complex classifiers in a cascade structure which dramatically increases the speed of the detector by focusing attention on promising regions of the image. The notion behind focus of attention approaches is that it is often possible to rapidly determine where in an image a face might occur (Tsotsos et al., 1995; Itti et al., 1998; Amit and Geman, 1999; Fleuret and Geman, 2001). More complex processing is reserved only for these promising regions. The key measure of such an approach is the "false negative" rate of the attentional process. It must be the case that all, or almost all, face instances are selected by the attentional filter.

We will describe a process for training an extremely simple and efficient classifier which can be used as a "supervised" focus of attention operator.[1] A face detection attentional operator can be learned which will filter out over 50% of the image while preserving 99% of the faces (as evaluated over a large dataset). This filter is exceedingly efficient; it can be evaluated in 20 simple operations per location/scale (approximately 60 microprocessor instructions).

Those sub-windows which are not rejected by the initial classifier are processed by a sequence of classifiers, each slightly more complex than the last. If any classifier rejects the sub-window, no further processing is performed. The structure of the cascaded detection process is essentially that of a degenerate decision tree, and as such is related to the work of Fleuret and Geman (2001) and Amit and Geman (1999).

The complete face detection cascade has 38 classifiers, which total over 80,000 operations. Nevertheless the cascade structure results in extremely rapid average detection times. On a difficult dataset, containing 507 faces and 75 million sub-windows, faces are detected using an average of 270 microprocessor instructions per sub-window. In comparison, this system is about 15 times faster than an implementation of the detection system constructed by Rowley et al. (1998).[2]

An extremely fast face detector will have broad practical applications. These include user interfaces, image databases, and teleconferencing. This increase in speed will enable real-time face detection applications on systems where they were previously infeasible. In applications where rapid frame-rates are not necessary, our system will allow for significant additional post-processing and analysis. In addition our system can be implemented on a wide range of small low power devices, including hand-helds and embedded processors. In our lab we have implemented this face detector on a low power 200 mips *Strong Arm* processor which lacks floating point hardware and have achieved detection at two frames per second.

### 1.1.  Overview

The remaining sections of the paper will discuss the implementation of the detector, related theory, and experiments. Section 2 will detail the form of the features as well as a new scheme for computing them rapidly. Section 3 will discuss the method in which these features are combined to form a classifier. The machine learning method used, a application of AdaBoost, also acts as a feature selection mechanism. While the classifiers that are constructed in this way have good computational and classification performance, they are far too slow for a real-time classifier. Section 4 will describe a method for constructing a cascade of classifiers which

149 together yield an extremely reliable and efficient face
150 detector. Section 5 will describe a number of experi-
151 mental results, including a detailed description of our
152 experimental methodology. Finally Section 6 contains
153 a discussion of this system and its relationship to re-
154 lated systems.

## 155  2.  Features

156 Our face detection procedure classifies images based
157 on the value of simple features. There are many moti-
158 vations for using features rather than the pixels directly.
159 The most common reason is that features can act to en-
160 code ad-hoc domain knowledge that is difficult to learn
161 using a finite quantity of training data. For this system
162 there is also a second critical motivation for features:
163 the feature-based system operates much faster than a
164 pixel-based system.
165     The simple features used are reminiscent of Haar
166 basis functions which have been used by Papageorgiou
167 et al. (1998). More specifically, we use three kinds of
168 features. The value of a *two-rectangle feature* is the
169 difference between the sum of the pixels within two
170 rectangular regions. The regions have the same size
171 and shape and are horizontally or vertically adjacent
172 (see Fig. 1). A *three-rectangle feature* computes the
173 sum within two outside rectangles subtracted from the
174 sum in a center rectangle. Finally a *four-rectangle fea-*
175 *ture* computes the difference between diagonal pairs of
176 rectangles.
177     Given that the base resolution of the detector is
178 $24 \times 24$, the exhaustive set of rectangle features is

quite large, 160,000. Note that unlike the Haar basis, 179
the set of rectangle features is overcomplete.[3] 180

### 2.1.  *Integral Image* 181

Rectangle features can be computed very rapidly using 182
an intermediate representation for the image which we 183
call the integral image.[4] The integral image at location 184
$x, y$ contains the sum of the pixels above and to the left 185
of $x, y$, inclusive: 186

$$ii(x, y) = \sum_{x' \leq x, y' \leq y} i(x', y'),$$

where $ii(x, y)$ is the integral image and $i(x, y)$ is the 187
original image (see Fig. 2). Using the following pair of 188
recurrences: 189

$$s(x, y) = s(x, y - 1) + i(x, y) \qquad (1)$$
$$ii(x, y) = ii(x - 1, y) + s(x, y) \qquad (2)$$

(where $s(x, y)$ is the cumulative row sum, $s(x, -1) =$ 190
$0$, and $ii(-1, y) = 0$) the integral image can be com- 191
puted in one pass over the original image. 192
    Using the integral image any rectangular sum can be 193
computed in four array references (see Fig. 3). Clearly 194
the difference between two rectangular sums can be 195
computed in eight references. Since the two-rectangle 196
features defined above involve adjacent rectangular 197
sums they can be computed in six array references, 198
eight in the case of the three-rectangle features, and 199
nine for four-rectangle features. 200
    One alternative motivation for the integral im- 201
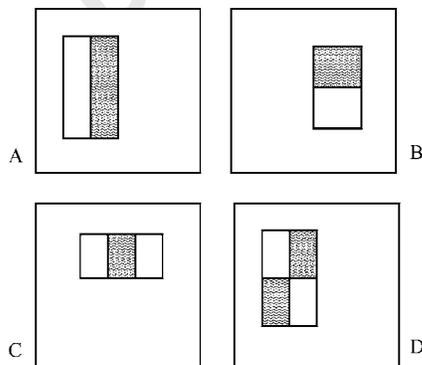age comes from the "boxlets" work of Simard et al. 202



*Figure 1.* Example rectangle features shown relative to the enclos-
ing detection window. The sum of the pixels which lie within the
white rectangles are subtracted from the sum of pixels in the grey
rectangles. Two-rectangle features are shown in (A) and (B). Figure
(C) shows a three-rectangle feature, and (D) a four-rectangle feature.



*Figure 2.* The value of the integral image at point $(x, y)$ is the sum
of all the pixels above and to the left.
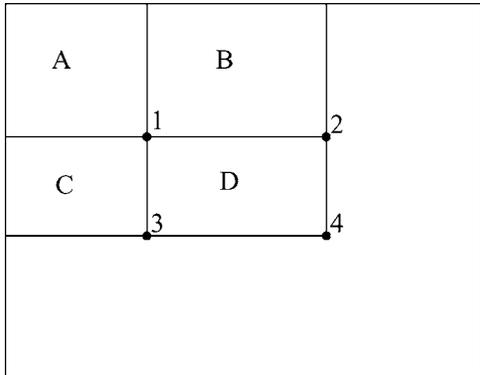
140    *Viola and Jones*



*Figure 3*.    The sum of the pixels within rectangle *D* can be computed with four array references. The value of the integral image at location 1 is the sum of the pixels in rectangle *A*. The value at location 2 is $A + B$, at location 3 is $A + C$, and at location 4 is $A + B + C + D$. The sum within *D* can be computed as $4 + 1 - (2 + 3)$.

203  (1999). The authors point out that in the case of linear
204  operations (e.g. $f \cdot g$), any invertible linear operation
205  can be applied to $f$ or $g$ if its inverse is applied to the
206  result. For example in the case of convolution, if the
207  derivative operator is applied both to the image and the
208  kernel the result must then be double integrated:

$$f * g = \int \int (f' * g').$$

209  The authors go on to show that convolution can be
210  significantly accelerated if the derivatives of $f$ and $g$
211  are sparse (or can be made so). A similar insight is that
212  an invertible linear operation can be applied to $f$ if its
213  inverse is applied to $g$:

$$(f'') * \left( \int \int g \right) = f * g.$$

214  Viewed in this framework computation of the rect-
215  angle sum can be expressed as a dot product, $i \cdot r$, where
216  $i$ is the image and $r$ is the box car image (with value
217  1 within the rectangle of interest and 0 outside). This
218  operation can be rewritten

$$i \cdot r = \left( \int \int i \right) \cdot r''.$$

219  The integral image is in fact the double integral of the
220  image (first along rows and then along columns). The
221  second derivative of the rectangle (first in row and then
222  in column) yields four delta functions at the corners of

223  the rectangle. Evaluation of the second dot product is
224  accomplished with four array accesses.

### 2.2.    *Feature Discussion*

226  Rectangle features are somewhat primitive when
227  compared with alternatives such as steerable filters
228  (Freeman and Adelson, 1991; Greenspan et al., 1994).
229  Steerable filters, and their relatives, are excellent for the
230  detailed analysis of boundaries, image compression,
231  and texture analysis. While rectangle features are also
232  sensitive to the presence of edges, bars, and other sim-
233  ple image structure, they are quite coarse. Unlike steer-
234  able filters, the only orientations available are vertical,
235  horizontal and diagonal. Since orthogonality is not cen-
236  tral to this feature set, we choose to generate a very
237  large and varied set of rectangle features. Typically the
238  representation is about 400 times overcomplete. This
239  overcomplete set provides features of arbitrary aspect
240  ratio and of finely sampled location. Empirically it ap-
241  pears as though the set of rectangle features provide
242  a rich image representation which supports effective
243  learning. The extreme computational efficiency of rect-
244  angle features provides ample compensation for their
245  limitations.

246  In order to appreciate the computational advantage
247  of the integral image technique, consider a more con-
248  ventional approach in which a pyramid of images is
249  computed. Like most face detection systems, our de-
250  tector scans the input at many scales; starting at the
251  base scale in which faces are detected at a size of
252  $24 \times 24$ pixels, a 384 by 288 pixel image is scanned
253  at 12 scales each a factor of 1.25 larger than the last.
254  The conventional approach is to compute a pyramid of
255  12 images, each 1.25 times smaller than the previous
256  image. A fixed scale detector is then scanned across
257  each of these images. Computation of the pyramid,
258  while straightforward, requires significant time. Imple-
259  mented efficiently on conventional hardware (using bi-
260  linear interpolation to scale each level of the pyramid) it
261  takes around .05 seconds to compute a 12 level pyramid
262  of this size (on an Intel PIII 700 MHz processor).[5]

263  In contrast we have defined a meaningful set of rect-
264  angle features, which have the property that a single
265  feature can be evaluated at any scale and location in a
266  few operations. We will show in Section 4 that effec-
267  tive face detectors can be constructed with as few as *two*
268  rectangle features. Given the computational efficiency
269  of these features, the face detection process can be com-
270  pleted for an entire image at every scale at 15 frames per

271 second, about the same time required to evaluate the 12
272 level image pyramid alone. Any procedure which re-
273 quires a pyramid of this type will necessarily run slower
274 than our detector.

## 3. Learning Classification Functions

276 Given a feature set and a training set of positive and
277 negative images, any number of machine learning ap-
278 proaches could be used to learn a classification func-
279 tion. Sung and Poggio use a mixture of Gaussian model
280 (Sung and Poggio, 1998). Rowley et al. (1998) use a
281 small set of simple image features and a neural net-
282 work. Osuna et al. (1997b) used a support vector ma-
283 chine. More recently Roth et al. (2000) have proposed
284 a new and unusual image representation and have used
285 the Winnow learning procedure.

286 Recall that there are 160,000 rectangle features as-
287 sociated with each image sub-window, a number far
288 larger than the number of pixels. Even though each
289 feature can be computed very efficiently, computing
290 the complete set is prohibitively expensive. Our hy-
291 pothesis, which is borne out by experiment, is that a
292 very small number of these features can be combined
293 to form an effective classifier. The main challenge is to
294 find these features.

295 In our system a variant of AdaBoost is used both
296 to select the features and to train the classifier (Freund
297 and Schapire, 1995). In its original form, the AdaBoost
298 learning algorithm is used to boost the classification
299 performance of a simple learning algorithm (e.g., it
300 might be used to boost the performance of a simple per-
301 ceptron). It does this by combining a collection of weak
302 classification functions to form a stronger classifier. In
303 the language of boosting the simple learning algorithm
304 is called a weak learner. So, for example the percep-
305 tron learning algorithm searches over the set of possible
306 perceptrons and returns the perceptron with the lowest
307 classification error. The learner is called weak because
308 we do not expect even the best classification function to
309 classify the training data well (i.e. for a given problem
310 the best perceptron may only classify the training data
311 correctly 51% of the time). In order for the weak learner
312 to be boosted, it is called upon to solve a sequence of
313 learning problems. After the first round of learning, the
314 examples are re-weighted in order to emphasize those
315 which were incorrectly classified by the previous weak
316 classifier. The final strong classifier takes the form of a
317 perceptron, a weighted combination of weak classifiers
318 followed by a threshold.[6]

319 The formal guarantees provided by the AdaBoost
320 learning procedure are quite strong. Freund and
321 Schapire proved that the training error of the strong
322 classifier approaches zero exponentially in the number
323 of rounds. More importantly a number of results
324 were later proved about generalization performance
325 (Schapire et al., 1997). The key insight is that gen-
326 eralization performance is related to the margin of the
327 examples, and that AdaBoost achieves large margins
328 rapidly.

329 The conventional AdaBoost procedure can be eas-
330 ily interpreted as a greedy feature selection process.
331 Consider the general problem of boosting, in which a
332 large set of classification functions are combined using
333 a weighted majority vote. The challenge is to associate
334 a large weight with each good classification function
335 and a smaller weight with poor functions. AdaBoost is
336 an aggressive mechanism for selecting a small set of
337 good classification functions which nevertheless have
338 significant variety. Drawing an analogy between weak
339 classifiers and features, AdaBoost is an effective pro-
340 cedure for searching out a small number of good "fea-
341 tures" which nevertheless have significant variety.

342 One practical method for completing this analogy is
343 to restrict the weak learner to the set of classification
344 functions each of which depend on a single feature.
345 In support of this goal, the weak learning algorithm is
346 designed to select the single rectangle feature which
347 best separates the positive and negative examples (this
348 is similar to the approach of Tieu and Viola (2000) in
349 the domain of image database retrieval). For each fea-
350 ture, the weak learner determines the optimal threshold
351 classification function, such that the minimum num-
352 ber of examples are misclassified. A weak classifier
353 ($h(x, f, p, \theta)$) thus consists of a feature ($f$), a thresh-
354 old ($\theta$) and a polarity ($p$) indicating the direction of the
355 inequality:

$$h(x, f, p, \theta) = \begin{cases} 1 & \text{if } pf(x) < p\theta \\ 0 & \text{otherwise} \end{cases}$$

356 Here $x$ is a $24 \times 24$ pixel sub-window of an image.

357 In practice no single feature can perform the classifi-
358 cation task with low error. Features which are selected
359 early in the process yield error rates between 0.1 and
360 0.3. Features selected in later rounds, as the task be-
361 comes more difficult, yield error rates between 0.4 and
362 0.5. Table 1 shows the learning algorithm.

363 The weak classifiers that we use (thresholded single
364 features) can be viewed as single node decision trees.

142     *Viola and Jones*

*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.

- Given example images $(x_1, y_1), \ldots, (x_n, y_n)$ 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 $t = 1, \ldots, T$:

  1. Normalize the weights, $w_{t,i} \leftarrow \frac{w_{t,i}}{\sum_{j=1}^{n} w_{t,j}}$
  2. Select the best weak classifier with respect to the weighted error

     $$\epsilon_t = \min_{f,p,\theta} \sum_i w_i \, |h(x_i, f, p, \theta) - y_i| \, .$$

     See Section 3.1 for a discussion of an efficient implementation.
  3. Define $h_t(x) = h(x, f_t, p_t, \theta_t)$ where $f_t$, $p_t$, and $\theta_t$ are the minimizers of $\epsilon_t$.
  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 $\beta_t = \frac{\epsilon_t}{1-\epsilon_t}$.
- The final strong classifier is:

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

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

Such structures have been called decision stumps in the machine learning literature. The original work of Freund and Schapire (1995) also experimented with boosting decision stumps.

### 3.1. Learning Discussion

The algorithm described in Table 1 is used to select key weak classifiers from the set of possible weak classifiers. While the AdaBoost process is quite efficient, the set of weak classifier is extraodinarily large. Since there is one weak classifier for each distinct feature/threshold combination, there are effectively $KN$ weak classifiers, where $K$ is the number of features and $N$ is the number of examples. In order to appreciate the dependency on $N$, suppose that the examples are sorted by a given feature value. With respect to the training process any two thresholds that lie between the same pair of sorted examples is equivalent. Therefore the total number of distinct thresholds is $N$. Given a task with $N = 20000$ and $K = 160000$ there are 3.2 billion distinct binary weak classifiers.

The wrapper method can also be used to learn a perceptron which utilizes $M$ weak classifiers (John et al., 1994) The wrapper method also proceeds incrementally by adding one weak classifier to the perceptron in each round. The weak classifier added is the one which when added to the current set yields a perceptron with lowest error. Each round takes at least $O(NKN)$ (or 60 Trillion operations); the time to enumerate all binary features and evaluate each example using that feature. This neglects the time to learn the perceptron weights. Even so, the final work to learn a 200 feature classifier would be something like $O(MNKN)$ which is $10^{16}$ operations.

The key advantage of AdaBoost as a feature selection mechanism, over competitors such as the wrapper method, is the speed of learning. Using AdaBoost a 200 feature classifier can be learned in $O(MNK)$ or about $10^{11}$ operations. One key advantage is that in each round the entire dependence on previously selected features is efficiently and compactly encoded using the example weights. These weights can then be used to evaluate a given weak classifier in constant time.

The weak classifier selection algorithm proceeds as follows. For each feature, the examples are sorted based on feature value. The AdaBoost optimal threshold for that feature can then be computed in a single pass over this sorted list. For each element in the sorted list, four sums are maintained and evaluated: the total sum of positive example weights $T^+$, the total sum of negative example weights $T^-$, the sum of positive weights below the current example $S^+$ and the sum of negative weights below the current example $S^-$. The error for a threshold which splits the range between the current and previous example in the sorted list is:

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

or the minimum of the error of labeling all examples below the current example negative and labeling the examples above positive versus the error of the converse. These sums are easily updated as the search proceeds.

Many general feature selection procedures have been proposed (see chapter 8 of Webb (1999) for a review). Our final application demanded a very aggressive process which would discard the vast majority of features. For a similar recognition problem Papageorgiou et al. (1998) proposed a scheme for feature selection based

429 on feature variance. They demonstrated good results se-
430 lecting 37 features out of a total 1734 features. While
431 this is a significant reduction, the number of features
432 evaluated for every image sub-window is still reason-
433 ably large.
434 Roth et al. (2000) propose a feature selection process
435 based on the Winnow exponential perceptron learning
436 rule. These authors use a very large and unusual feature
437 set, where *each pixel* is mapped into a binary vector of $d$
438 dimensions (when a particular pixel takes on the value
439 $x$, in the range $[0, d - 1]$, the $x$-th dimension is set to
440 1 and the other dimensions to 0). The binary vectors
441 for each pixel are concatenated to form a single binary
442 vector with $nd$ dimensions ($n$ is the number of pixels).
443 The classification rule is a perceptron, which assigns
444 one weight to each dimension of the input vector. The
445 Winnow learning process converges to a solution where
446 many of these weights are zero. Nevertheless a very
447 large number of features are retained (perhaps a few
448 hundred or thousand).

449 *3.2. Learning Results*

450 While details on the training and performance of the
451 final system are presented in Section 5, several sim-
ple results merit discussion. Initial experiments demon-

452 strated that a classifier constructed from 200 features
453 would yield reasonable results (see Fig. 4). Given a
454 detection rate of 95% the classifier yielded a false pos-
455 itive rate of 1 in 14084 on a testing dataset. This is
456 promising, but for a face detector to be practical for
457 real applications, the false positive rate must be closer
458 to 1 in 1,000,000.
459 For the task of face detection, the initial rectangle
460 features selected by AdaBoost are meaningful and eas-
461 ily interpreted. The first feature selected seems to focus
462 on the property that the region of the eyes is often darker
463 than the region of the nose and cheeks (see Fig. 5). This
464 feature is relatively large in comparison with the detec-
465 tion sub-window, and should be somewhat insensitive
466 to size and location of the face. The second feature se-
467 lected relies on the property that the eyes are darker
468 than the bridge of the nose.
469 In summary the 200-feature classifier provides ini-
470 tial evidence that a boosted classifier constructed from
471 rectangle features is an effective technique for face de-
472 tection. In terms of detection, these results are com-
473 pelling but not sufficient for many real-world tasks. In
474 terms of computation, this classifier is very fast, re-
475 quiring 0.7 seconds to scan an 384 by 288 pixel im-
476 age. Unfortunately, the most straightforward tech-
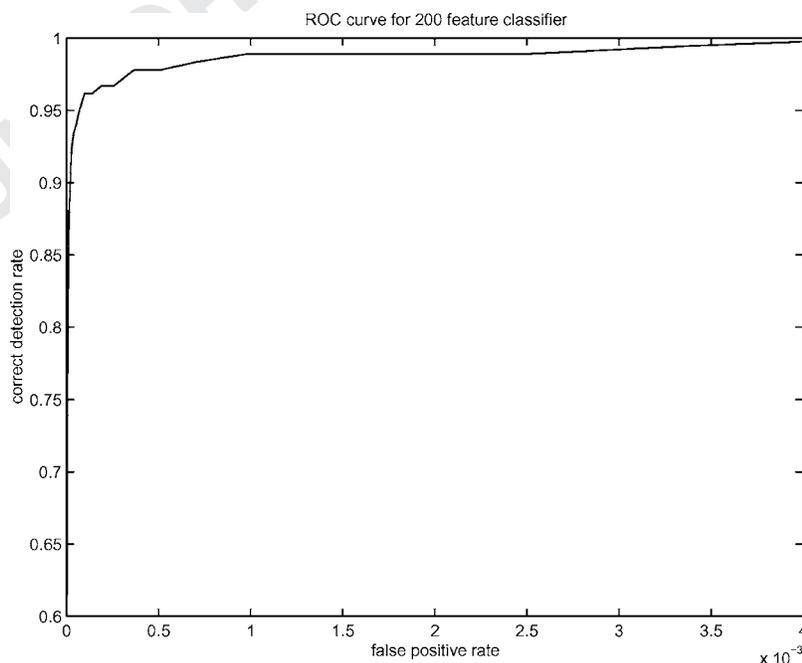nique for improving detection performance, adding



*Figure 4.* Receiver operating characteristic (ROC) curve for the 200 feature classifier.
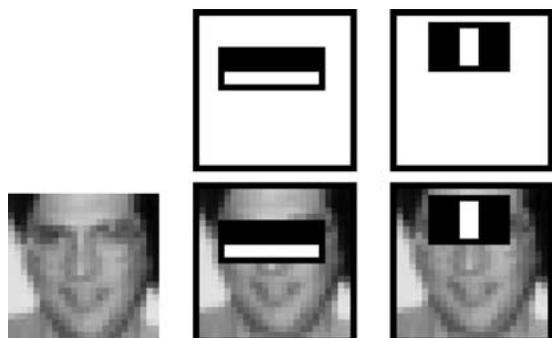
144    *Viola and Jones*



*Figure 5.* The first and second features selected by AdaBoost. The two features are shown in the top row and then overlayed on a typical training face in the bottom row. The first feature measures the difference in intensity between the region of the eyes and a region across the upper cheeks. The feature capitalizes on the observation that the eye region is often darker than the cheeks. The second feature compares the intensities in the eye regions to the intensity across the bridge of the nose.

477 features to the classifier, directly increases computation
478 time.

### 4. The Attentional Cascade

480 This section describes an algorithm for constructing a
481 cascade of classifiers which achieves increased detec-
482 tion performance while radically reducing computation
483 time. The key insight is that smaller, and therefore more
484 efficient, boosted classifiers can be constructed which
485 reject many of the negative sub-windows while detect-
486 ing almost all positive instances. Simpler classifiers are
487 used to reject the majority of sub-windows before more
488 complex classifiers are called upon to achieve low false
489 positive rates.

490 Stages in the cascade are constructed by training
491 classifiers using AdaBoost. Starting with a two-feature
492 strong classifier, an effective face filter can be obtained
493 by adjusting the strong classifier threshold to mini-
494 mize false negatives. The initial AdaBoost threshold,
495 $\frac{1}{2}\sum_{t=1}^{T}\alpha_t$, is designed to yield a low error rate on the
496 training data. A lower threshold yields higher detec-
497 tion rates and higher false positive rates. Based on per-
498 formance measured using a validation training set, the
499 two-feature classifier can be adjusted to detect 100% of
500 the faces with a false positive rate of 50%. See Fig. 5 for
501 a description of the two features used in this classifier.

502 The detection performance of the two-feature clas-
503 sifier is far from acceptable as a face detection system.
504 Nevertheless the classifier can significantly reduce the

505 number of sub-windows that need further processing
506 with very few operations:

507 1. Evaluate the rectangle features (requires between 6
508     and 9 array references per feature).
509 2. Compute the weak classifier for each feature (re-
510     quires one threshold operation per feature).
511 3. Combine the weak classifiers (requires one multiply
512     per feature, an addition, and finally a threshold).

513 A two feature classifier amounts to about 60 mi-
514 croprocessor instructions. It seems hard to imagine
515 that any simpler filter could achieve higher rejection
516 rates. By comparison, scanning a simple image tem-
517 plate would require at least 20 times as many operations
518 per sub-window.

519 The overall form of the detection process is that of
520 a degenerate decision tree, what we call a "cascade"
521 (Quinlan, 1986) (see Fig. 6). A positive result from
522 the first classifier triggers the evaluation of a second
523 classifier which has also been adjusted to achieve very
524 high detection rates. A positive result from the second
525 classifier triggers a third classifier, and so on. A negative
526 outcome at any point leads to the immediate rejection
527 of the sub-window.

528 The structure of the cascade reflects the fact that
529 within any single image an overwhelming majority of
530 sub-windows are negative. As such, the cascade at-
531 tempts to reject as many negatives as possible at the
532 earliest stage possible. While a positive instance will



*Figure 6.* 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.

**533** trigger the evaluation of every classifier in the cascade,
**534** this is an exceedingly rare event.
**535**    Much like a decision tree, subsequent classifiers are
**536** trained using those examples which pass through all
**537** the previous stages. As a result, the second classifier
**538** faces a more difficult task than the first. The examples
**539** which make it through the first stage are "harder" than
**540** typical examples. The more difficult examples faced
**541** by deeper classifiers push the entire receiver operat-
**542** ing characteristic (ROC) curve downward. At a given
**543** detection rate, deeper classifiers have correspondingly
**544** higher false positive rates.

**545** *4.1.    Training a Cascade of Classifiers*

**546** The cascade design process is driven from a set of de-
**547** tection and performance goals. For the face detection
**548** task, past systems have achieved good detection rates
**549** (between 85 and 95 percent) and extremely low false
**550** positive rates (on the order of $10^{-5}$ or $10^{-6}$). The num-
**551** ber of cascade stages and the size of each stage must
**552** be sufficient to achieve similar detection performance
**553** while minimizing computation.
**554**    Given a trained cascade of classifiers, the false pos-
**555** itive rate of the cascade is

$$F = \prod_{i=1}^{K} f_i,$$

**556** where $F$ is the false positive rate of the cascaded clas-
**557** sifier, $K$ is the number of classifiers, and $f_i$ is the false
**558** positive rate of the $i$th classifier on the examples that
**559** get through to it. The detection rate is

$$D = \prod_{i=1}^{K} d_i,$$

**560** where $D$ is the detection rate of the cascaded classifier,
**561** $K$ is the number of classifiers, and $d_i$ is the detection
**562** rate of the $i$th classifier on the examples that get through
**563** to it.
**564**    Given concrete goals for overall false positive and
**565** detection rates, target rates can be determined for each
**566** stage in the cascade process. For example a detection
**567** rate of 0.9 can be achieved by a 10 stage classifier if
**568** each stage has a detection rate of 0.99 (since $0.9 \approx$
**569** $0.99^{10}$). While achieving this detection rate may sound
**570** like a daunting task, it is made significantly easier by the
**571** fact that each stage need only achieve a false positive
**572** rate of about 30% ($0.30^{10} \approx 6 \times 10^{-6}$).

**573** The number of features evaluated when scanning
**574** real images is necessarily a probabilistic process. Any
**575** given sub-window will progress down through the cas-
**576** cade, one classifier at a time, until it is decided that
**577** the window is negative or, in rare circumstances, the
**578** window succeeds in each test and is labelled positive.
**579** The expected behavior of this process is determined
**580** by the distribution of image windows in a typical test
**581** set. The key measure of each classifier is its "positive
**582** rate", the proportion of windows which are labelled as
**583** potentially containing a face. The expected number of
**584** features which are evaluated is:

$$N = n_0 + \sum_{i=1}^{K} \left( n_i \prod_{j<i} p_j \right)$$

**585** where $N$ is the expected number of features evaluated,
**586** $K$ is the number of classifiers, $p_i$ is the positive rate of
**587** the $i$th classifier, and $n_i$ are the number of features in the
**588** $i$th classifier. Interestingly, since faces are extremely
**589** rare, the "positive rate" is effectively equal to the false
**590** positive rate.
**591** The process by which each element of the cascade
**592** is trained requires some care. The AdaBoost learning
**593** procedure presented in Section 3 attempts only to min-
**594** imize errors, and is not specifically designed to achieve
**595** high detection rates at the expense of large false positive
**596** rates. One simple, and very conventional, scheme for
**597** trading off these errors is to adjust the threshold of the
**598** perceptron produced by AdaBoost. Higher thresholds
**599** yield classifiers with fewer false positives and a lower
**600** detection rate. Lower thresholds yield classifiers with
**601** more false positives and a higher detection rate. It is
**602** not clear, at this point, whether adjusting the threshold
**603** in this way preserves the training and generalization
**604** guarantees provided by AdaBoost.
**605** The overall training process involves two types of
**606** tradeoffs. In most cases classifiers with more features
**607** will achieve higher detection rates and lower false pos-
**608** itive rates. At the same time classifiers with more fea-
**609** tures require more time to compute. In principle one
**610** could define an optimization framework in which

- the number of classifier stages, **611**
- the number of features, $n_i$, of each stage, **612**
- the threshold of each stage **613**

**614** are traded off in order to minimize the expected num-
**615** ber of features $N$ given a target for $F$ and $D$. Unfortu-
**616** nately finding this optimum is a tremendously difficult
**617** problem.

146     *Viola and Jones*

*Table 2.*   The training algorithm for building a cascaded detector.

- User selects values for $f$, the maximum acceptable false positive rate per layer and $d$, the minimum acceptable detection rate per layer.
- User selects target overall false positive rate, $F_{target}$.
- $P$ = set of positive examples
- $N$ = set of negative examples
- $F_0 = 1.0; D_0 = 1.0$
- $i = 0$
- while $F_i > F_{target}$
  - $i \leftarrow i + 1$
  - $n_i = 0; F_i = F_{i-1}$
  - while $F_i > f \times F_{i-1}$
    - $* \ n_i \leftarrow n_i + 1$
    - $*$ Use $P$ and $N$ to train a classifier with $n_i$ features using AdaBoost
    - $*$ Evaluate current cascaded classifier on validation set to determine $F_i$ and $D_i$.
    - $*$ Decrease threshold for the $i$th classifier until the current cascaded classifier has a detection rate of at least $d \times D_{i-1}$ (this also affects $F_i$)
  - $N \leftarrow \emptyset$
  - If $F_i > F_{target}$ then evaluate the current cascaded detector on the set of non-face images and put any false detections into the set N

In practice a very simple framework is used to produce an effective classifier which is highly efficient. The user selects the maximum acceptable rate for $f_i$ and the minimum acceptable rate for $d_i$. Each layer of the cascade is trained by AdaBoost (as described in Table 1) with the number of features used being increased until the target detection and false positive rates are met for this level. The rates are determined by testing the current detector on a validation set. If the overall target false positive rate is not yet met then another layer is added to the cascade. The negative set for training subsequent layers is obtained by collecting all false detections found by running the current detector on a set of images which do not contain any instances of faces. This algorithm is given more precisely in Table 2.

### 4.2.   Simple Experiment

In order to explore the feasibility of the cascade approach two simple detectors were trained: a monolithic 200-feature classifier and a cascade of ten 20-feature classifiers. The first stage classifier in the cascade was trained using 5000 faces and 10000 non-face sub-windows randomly chosen from non-face images. The second stage classifier was trained on the same 5000 faces plus 5000 false positives of the first classifier. This process continued so that subsequent stages were trained using the false positives of the previous stage.

The monolithic 200-feature classifier was trained on the union of all examples used to train all the stages of the cascaded classifier. Note that without reference to the cascaded classifier, it might be difficult to select a set of non-face training examples to train the monolithic classifier. We could of course use all possible sub-windows from all of our non-face images, but this would make the training time impractically long. The sequential way in which the cascaded classifier is trained effectively reduces the non-face training set by throwing out easy examples and focusing on the "hard" ones.

Figure 7 gives the ROC curves comparing the performance of the two classifiers. It shows that there is little difference between the two in terms of accuracy. However, there is a big difference in terms of speed. The cascaded classifier is nearly 10 times faster since its first stage throws out most non-faces so that they are never evaluated by subsequent stages.

### 4.3.   Detector Cascade Discussion

There is a hidden benefit of training a detector as a sequence of classifiers which is that the effective number of negative examples that the final detector sees can be very large. One can imagine training a single large classifier with many features and then trying to speed up its running time by looking at partial sums of features and stopping the computation early if a partial sum is below the appropriate threshold. One drawback of such an approach is that the training set of negative examples would have to be relatively small (on the order of 10,000 to maybe 100,000 examples) to make training feasible. With the cascaded detector, the final layers of the cascade may effectively look through hundreds of millions of negative examples in order to find a set of 10,000 negative examples that the earlier layers of the cascade fail on. So the negative training set is much larger and more focused on the hard examples for a cascaded detector.

A notion similar to the cascade appears in the face detection system described by Rowley et al. (1998). Rowley et al. trained two neural networks. One network was moderately complex, focused on a small region of the image, and detected faces with a low false positive rate. They also trained a second neural network which
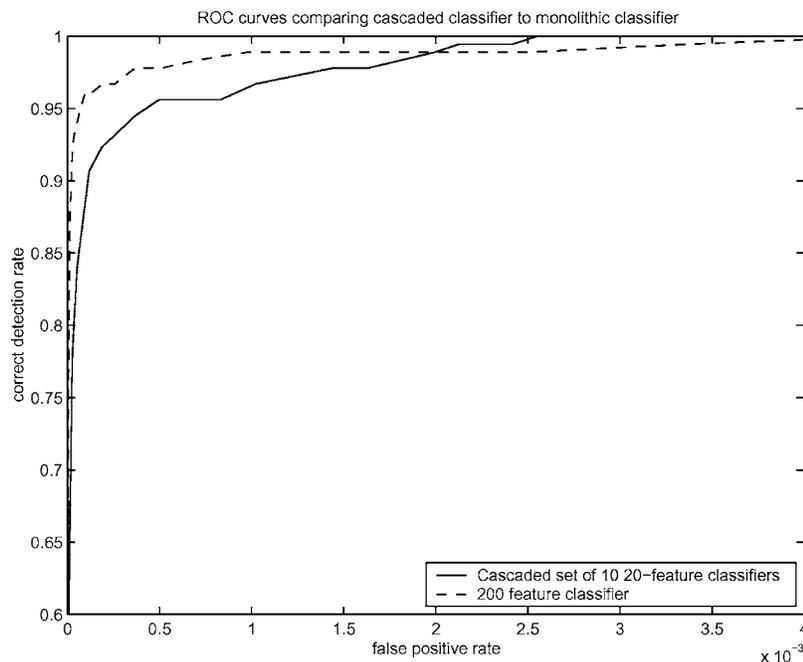
*Figure 7.* ROC curves comparing a 200-feature classifier with a cascaded classifier containing ten 20-feature classifiers. Accuracy is not significantly different, but the speed of the cascaded classifier is almost 10 times faster.

**689** was much faster, focused on a larger regions of the **690** image, and detected faces with a higher false positive **691** rate. Rowley et al. used the faster second network to **692** prescreen the image in order to find candidate regions **693** for the slower more accurate network. Though it is **694** difficult to determine exactly, it appears that Rowley **695** et al.'s two network face system is the fastest existing **696** face detector.[7] Our system uses a similar approach, but **697** it extends this two stage cascade to include 38 stages.

**698** The structure of the cascaded detection process is **699** essentially that of a degenerate decision tree, and as **700** such is related to the work of Amit and Geman (1999). **701** Unlike techniques which use a fixed detector, Amit and **702** Geman propose an alternative point of view where un- **703** usual co-occurrences of simple image features are used **704** to trigger the evaluation of a more complex detection **705** process. In this way the full detection process need not **706** be evaluated at many of the potential image locations **707** and scales. While this basic insight is very valuable, **708** in their implementation it is necessary to first evaluate **709** some feature detector at every location. These features **710** are then grouped to find unusual co-occurrences. In **711** practice, since the form of our detector and the fea- **712** tures that it uses are extremely efficient, the amortized **713** cost of evaluating our detector at *every scale and lo-*

**714** *cation* is much faster than finding and grouping edges **715** throughout the image.

**716** In recent work Fleuret and Geman (2001) have pre- **717** sented a face detection technique which relies on a **718** "chain" of tests in order to signify the presence of a **719** face at a particular scale and location. The image prop- **720** erties measured by Fleuret and Geman, disjunctions **721** of fine scale edges, are quite different than rectangle **722** features which are simple, exist at all scales, and are **723** somewhat interpretable. The two approaches also differ **724** radically in their learning philosophy. Because Fleuret **725** and Geman's learning process does not use negative **726** examples their approach is based more on density es- **727** timation, while our detector is purely discriminative. **728** Finally the false positive rate of Fleuret and Geman's **729** approach appears to be higher than that of previous ap- **730** proaches like Rowley et al. and this approach. In the **731** published paper the included example images each had **732** between 2 and 10 false positives. For many practical **733** tasks, it is important that the expected number of false **734** positives in any image be less than one (since in many **735** tasks the expected number of true positives is less than **736** one as well). Unfortunately the paper does not report **737** quantitative detection and false positive results on stan- **738** dard datasets.

148     *Viola and Jones*

**739**   **5.   Results**

**740**   This section describes the final face detection system.
**741**   The discussion includes details on the structure and
**742**   training of the cascaded detector as well as results on
**743**   a large real-world testing set.

**744**   *5.1.   Training Dataset*

**745**   The face training set consisted of 4916 hand labeled
**746**   faces scaled and aligned to a base resolution of 24 by
**747**   24 pixels. The faces were extracted from images down-
**748**   loaded during a random crawl of the world wide web.
**749**   Some typical face examples are shown in Fig. 8. The
**750**   training faces are only roughly aligned. This was done
**751**   by having a person place a bounding box around each
**752**   face just above the eyebrows and about half-way be-
**753**   tween the mouth and the chin. This bounding box was
**754**   then enlarged by 50% and then cropped and scaled to
**755**   24 by 24 pixels. No further alignment was done (i.e.
**756**   the eyes are not aligned). Notice that these examples
**757**   contain more of the head than the examples used by

**758**   Rowley et al. (1998) or Sung and Poggio (1998). Ini-
**759**   tial experiments also used 16 by 16 pixel training im-
**760**   ages in which the faces were more tightly cropped,
**761**   but got slightly worse results. Presumably the 24 by
**762**   24 examples include extra visual information such as
**763**   the contours of the chin and cheeks and the hair line
**764**   which help to improve accuracy. Because of the nature
**765**   of the features used, the larger sized sub-windows do
**766**   not slow performance. In fact, the additional informa-
**767**   tion contained in the larger sub-windows can be used
**768**   to reject non-faces earlier in the detection cascade.

**769**   *5.2.   Structure of the Detector Cascade*

**770**   The final detector is a 38 layer cascade of classifiers
**771**   which included a total of 6060 features.
**772**   The first classifier in the cascade is constructed us-
**773**   ing two features and rejects about 50% of non-faces
**774**   while correctly detecting close to 100% of faces. The
**775**   next classifier has ten features and rejects 80% of non-
**776**   faces while detecting almost 100% of faces. The next
**777**   two layers are 25-feature classifiers followed by three
**778**   50-feature classifiers followed by classifiers with a



*Figure 8.*   Example of frontal upright face images used for training.

**778** variety of different numbers of features chosen accord-
**779** ing to the algorithm in Table 2. The particular choices
**780** of number of features per layer was driven through
**781** a trial and error process in which the number of fea-
**782** tures were increased until a significant reduction in the
**783** false positive rate could be achieved. More levels were
**784** added until the false positive rate on the validation set
**785** was nearly zero while still maintaining a high correct
**786** detection rate. The final number of layers, and the size
**787** of each layer, are not critical to the final system perfor-
**788** mance. The procedure we used to choose the number
**789** of features per layer was guided by human intervention
**790** (for the first 7 layers) in order to reduce the training time
**791** for the detector. The algorithm described in Table 2 was
**792** modified slightly to ease the computational burden by
**793** specifying a minimum number of features per layer by
**794** hand and by adding more than 1 feature at a time. In
**795** later layers, 25 features were added at a time before
**796** testing on the validation set. This avoided having to
**797** test the detector on the validation set for every single
**798** feature added to a classifier.

**799** The non-face sub-windows used to train the first
**800** level of the cascade were collected by selecting ran-
**801** dom sub-windows from a set of 9500 images which
**802** did not contain faces. The non-face examples used to
**803** train subsequent layers were obtained by scanning the
**804** partial cascade across large non-face images and col-
**805** lecting false positives. A maximum of 6000 such non-
**806** face sub-windows were collected for each layer. There
**807** are approximately 350 million non-face sub-windows
**808** contained in the 9500 non-face images.

**809** Training time for the entire 38 layer detector was on
**810** the order of weeks on a single 466 MHz AlphaStation
**811** XP900. We have since parallelized the algorithm to
**812** make it possible to train a complete cascade in about a
**813** day.

**814** *5.3.    Speed of the Final Detector*

**815** The speed of the cascaded detector is directly related
**816** to the number of features evaluated per scanned sub-
**817** window. As discussed in Section 4.1, the number of fea-
**818** tures evaluated depends on the images being scanned.
**819** Since a large majority of the sub-windows are dis-
**820** carded by the first two stages of the cascade, an av-
**821** erage of 8 features out of a total of 6060 are eval-
**822** uated per sub-window (as evaluated on the MIT +
**823** CMU (Rowley et al., 1998). On a 700 Mhz Pentium
**824** III processor, the face detector can process a 384 by
**825** 288 pixel image in about .067 seconds (using a starting

**826** scale of 1.25 and a step size of 1.5 described below).
**827** This is roughly 15 times faster than the Rowley-Baluja-
**828** Kanade detector (Rowley et al., 1998) and about 600
**829** times faster than the Schneiderman-Kanade detector
**830** (Schneiderman and Kanade, 2000).

**831** *5.4.    Image Processing*

**832** All example sub-windows used for training were vari-
**833** ance normalized to minimize the effect of different
**834** lighting conditions. Normalization is therefore neces-
**835** sary during detection as well. The variance of an image
**836** sub-window can be computed quickly using a pair of
**837** integral images. Recall that $\sigma^2 = m^2 - \frac{1}{N}\sum x^2$, where
**838** $\sigma$ is the standard deviation, $m$ is the mean, and $x$ is
**839** the pixel value within the sub-window. The mean of a
**840** sub-window can be computed using the integral image.
**841** The sum of squared pixels is computed using an integral
**842** image of the image squared (i.e. two integral images
**843** are used in the scanning process). During scanning the
**844** effect of image normalization can be achieved by post
**845** multiplying the feature values rather than operating on
**846** the pixels.

**847** *5.5.    Scanning the Detector*

**848** The final detector is scanned across the image at multi-
**849** ple scales and locations. Scaling is achieved by scaling
**850** the detector itself, rather than scaling the image. This
**851** process makes sense because the features can be eval-
**852** uated at any scale with the same cost. Good detection
**853** results were obtained using scales which are a factor of
**854** 1.25 apart.

**855** The detector is also scanned across location. Sub-
**856** sequent locations are obtained by shifting the window
**857** some number of pixels $\Delta$. This shifting process is af-
**858** fected by the scale of the detector: if the current scale is
**859** $s$ the window is shifted by $[s\Delta]$, where [] is the round-
**860** ing operation.

**861** The choice of $\Delta$ affects both the speed of the de-
**862** tector as well as accuracy. Since the training images
**863** have some translational variability the learned detector
**864** achieves good detection performance in spite of small
**865** shifts in the image. As a result the detector sub-window
**866** can be shifted more than one pixel at a time. However,
**867** a step size of more than one pixel tends to decrease the
**868** detection rate slightly while also decreasing the number
**869** of false positives. We present results for two different
**870** step sizes.

150    *Viola and Jones*

**871**    *5.6.    Integration of Multiple Detections*

**872** Since the final detector is insensitive to small changes
**873** in translation and scale, multiple detections will usually
**874** occur around each face in a scanned image. The same
**875** is often true of some types of false positives. In practice
**876** it often makes sense to return one final detection per
**877** face. Toward this end it is useful to postprocess the
**878** detected sub-windows in order to combine overlapping
**879** detections into a single detection.

**880**    In these experiments detections are combined in a
**881** very simple fashion. The set of detections are first par-
**882** titioned into disjoint subsets. Two detections are in the
**883** same subset if their bounding regions overlap. Each
**884** partition yields a single final detection. The corners of
**885** the final bounding region are the average of the corners
**886** of all detections in the set.

**887**    In some cases this postprocessing decreases the num-
**888** ber of false positives since an overlapping subset of
**889** false positives is reduced to a single detection.

**890**    *5.7.    Experiments on a Real-World Test Set*

**891** We tested our system on the MIT + CMU frontal face
**892** test set (Rowley et al., 1998). This set consists of 130

**893** images with 507 labeled frontal faces. A ROC curve
**894** showing the performance of our detector on this test
**895** set is shown in Fig. 9. To create the ROC curve the
**896** threshold of the perceptron on the final layer classifier
**897** is adjusted from $+\infty$ to $-\infty$. Adjusting the threshold to
**898** $+\infty$ will yield a detection rate of 0.0 and a false positive
**899** rate of 0.0. Adjusting the threshold to $-\infty$, however,
**900** increases both the detection rate and false positive rate,
**901** but only to a certain point. Neither rate can be higher
**902** than the rate of the detection cascade minus the final
**903** layer. In effect, a threshold of $-\infty$ is equivalent to re-
**904** moving that layer. Further increasing the detection and
**905** false positive rates requires decreasing the threshold
**906** of the next classifier in the cascade. Thus, in order to
**907** construct a complete ROC curve, classifier layers are
**908** removed. We use the *number* of false positives as op-
**909** posed to the *rate* of false positives for the *x*-axis of
**910** the ROC curve to facilitate comparison with other sys-
**911** tems. To compute the false positive rate, simply divide
**912** by the total number of sub-windows scanned. For the
**913** case of $\Delta = 1.0$ and starting scale = 1.0, the number
**914** of sub-windows scanned is 75,081,800. For $\Delta = 1.5$
**915** and starting scale = 1.25, the number of sub-windows
**916** scanned is 18,901,947.

**917**    Unfortunately, most previous published results on
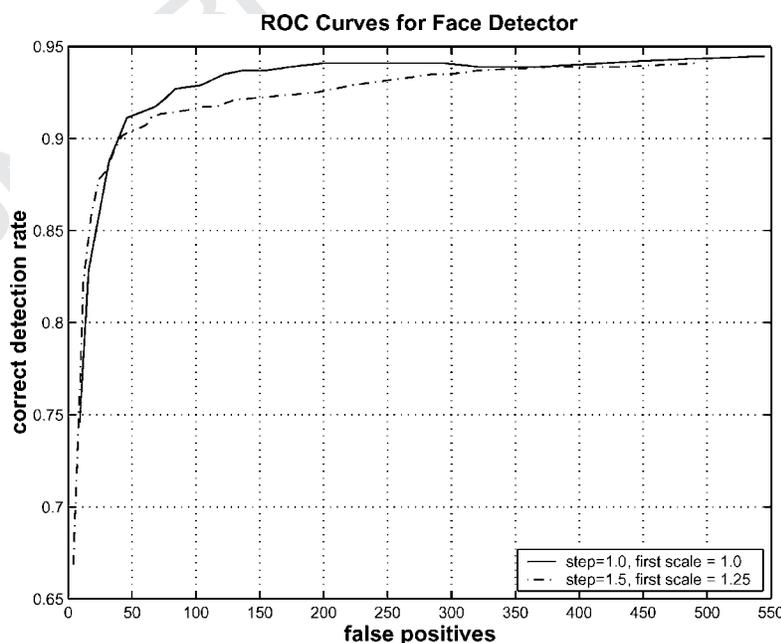face detection have only included a single operating



*Figure 9.*    ROC curves for our face detector on the MIT + CMU test set. The detector was run once using a step size of 1.0 and starting scale of 1.0 (75,081,800 sub-windows scanned) and then again using a step size of 1.5 and starting scale of 1.25 (18,901,947 sub-windows scanned). In both cases a scale factor of 1.25 was used.

*Table 3.*   Detection rates for various numbers of false positives on the MIT + CMU test set containing 130 images and 507 faces.

| Detector | False detections | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | 10 | 31 | 50 | 65 | 78 | 95 | 167 | 422 |
| Viola-Jones | 76.1% | 88.4% | 91.4% | 92.0% | 92.1% | 92.9% | 93.9% | 94.1% |
| Viola-Jones (voting) | 81.1% | 89.7% | 92.1% | 93.1% | 93.1% | 93.2% | 93.7% | – |
| Rowley-Baluja-Kanade | 83.2% | 86.0% | – | – | – | 89.2% | 90.1% | 89.9% |
| Schneiderman-Kanade | – | – | – | 94.4% | – | – | – | – |
| Roth-Yang-Ahuja | – | – | – | – | (94.8%) | – | – | – |

regime (i.e. single point on the ROC curve). To make comparison with our detector easier we have listed our detection rate for the same false positive rate reported by the other systems. Table 3 lists the detection rate for various numbers of false detections for our system as well as other published systems. For the Rowley-Baluja-Kanade results (Rowley et al., 1998), a number of different versions of their detector were tested yielding a number of different results. While these various results are not actually points on a ROC curve for a particular detector, they do indicate a number of different performance points that can be achieved with their approach. They did publish ROC curves for two of their detectors, but these ROC curves did not represent their best results. For the Roth-Yang-Ahuja detector (Roth et al., 2000), they reported their result on the MIT + CMU test set minus 5 images containing line drawn faces removed. So their results are for a subset of the MIT + CMU test set containing 125 images with 483 faces. Presumably their detection rate would be lower if the full test set was used. The parentheses around their detection rate indicates this slightly different test set. The Sung and Poggio face detector (Sung and Poggio, 1998) was tested on the MIT subset of the MIT + CMU test set since the CMU portion did not exist yet. The MIT test set contains 23 images with 149 faces. They achieved a detection rate of 79.9% with 5 false positives. Our detection rate with 5 false positives is 77.8% on the MIT test set.

Figure 10 shows the output of our face detector on some test images from the MIT + CMU test set.

### 5.7.1. A Simple Voting Scheme Further Improves Results.
The best results were obtained through the combination of three detectors trained using different initial negative examples, slightly different weighting on negative versus positive errors, and slightly different criteria for trading off false positives for classifier size. These three systems performed similarly on the final task, but in some cases errors were different. The detection results from these three detectors were combined by retaining only those detections where at least 2 out of 3 detectors agree. This improves the final detection rate as well as eliminating more false positives. Since detector errors are not uncorrelated, the combination results in a measurable, but modest, improvement over the best single detector.

### 5.7.2. Failure Modes.
By observing the performance of our face detector on a number of test images we have noticed a few different failure modes.

The face detector was trained on frontal, upright faces. The faces were only very roughly aligned so there is some variation in rotation both in plane and out of plane. Informal observation suggests that the face detector can detect faces that are tilted up to about $\pm 15$ degrees in plane and about $\pm 45$ degrees out of plane (toward a profile view). The detector becomes unreliable with more rotation than this.

We have also noticed that harsh backlighting in which the faces are very dark while the background is relatively light sometimes causes failures. It is interesting to note that using a nonlinear variance normalization based on robust statistics to remove outliers improves the detection rate in this situation. The problem with such a normalization is the greatly increased computational cost within our integral image framework.

Finally, our face detector fails on significantly occluded faces. If the eyes are occluded for example, the detector will usually fail. The mouth is not as important and so a face with a covered mouth will usually still be detected.
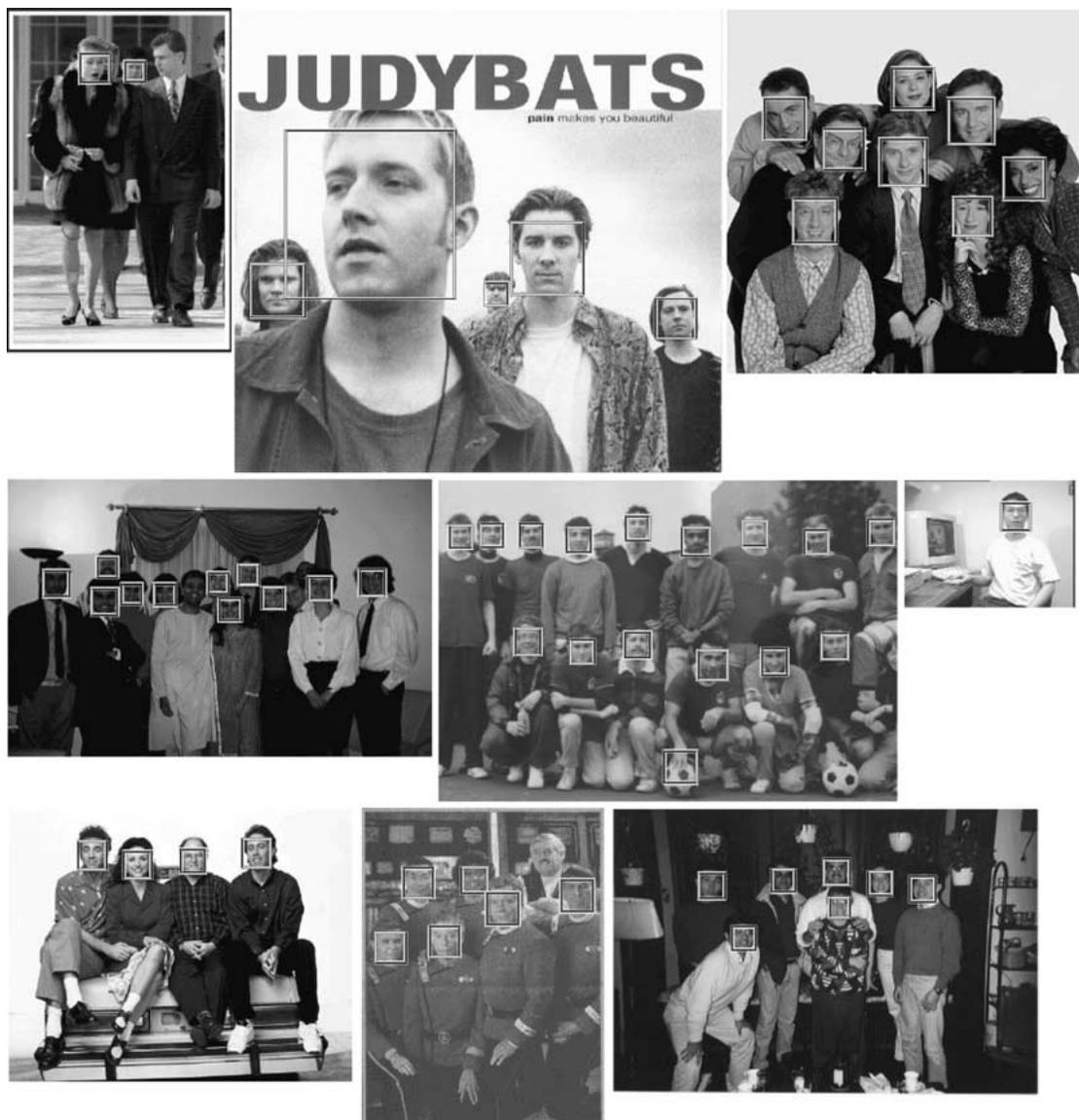
152       *Viola and Jones*



*Figure 10.*    Output of our face detector on a number of test images from the MIT + CMU test set.

## 6.    Conclusions

We have presented an approach for face detection which minimizes computation time while achieving high detection accuracy. The approach was used to construct a face detection system which is approximately 15 times faster than any previous approach. Preliminary experiments, which will be described elsewhere, show that highly efficient detectors for other objects, such as pedestrians or automobiles, can also be constructed in this way.

This paper brings together new algorithms, representations, and insights which are quite generic and may well have broader application in computer vision and image processing.

The first contribution is a new a technique for computing a rich set of image features using the integral image. In order to achieve true scale invariance, almost all face detection systems must operate on multiple image scales. The integral image, by eliminating the need to compute a multi-scale image pyramid, reduces the initial image processing required for face detection

1011 significantly. Using the integral image, face detection
1012 is completed in almost the same time as it takes for an
1013 image pyramid to be computed.
1014    While the integral image should also have immedi-
1015 ate use for other systems which have used Haar-like
1016 features such as Papageorgiou et al. (1998), it can fore-
1017 seeably have impact on any task where Haar-like fea-
1018 tures may be of value. Initial experiments have shown
1019 that a similar feature set is also effective for the task
1020 of parameter estimation, where the expression of a
1021 face, the position of a head, or the pose of an object is
1022 determined.
1023    The second contribution of this paper is a simple
1024 and efficient classifier built from computationally ef-
1025 ficient features using AdaBoost for feature selection.
1026 This classifier is clearly an effective one for face detec-
1027 tion and we are confident that it will also be effective in
1028 other domains such as automobile or pedestrian detec-
1029 tion. Furthermore, the idea of an aggressive and effec-
1030 tive technique for feature selection should have impact
1031 on a wide variety of learning tasks. Given an effective
1032 tool for feature selection, the system designer is free to
1033 define a very large and very complex set of features as
1034 input for the learning process. The resulting classifier
1035 is nevertheless computationally efficient, since only a
1036 small number of features need to be evaluated during
1037 run time. Frequently the resulting classifier is also quite
1038 simple; within a large set of complex features it is more
1039 likely that a few critical features can be found which
1040 capture the structure of the classification problem in a
1041 straightforward fashion.
1042    The third contribution of this paper is a technique for
1043 constructing a cascade of classifiers which radically
1044 reduces computation time while improving detection
1045 accuracy. Early stages of the cascade are designed to
1046 reject a majority of the image in order to focus subse-
1047 quent processing on promising regions. One key point
1048 is that the cascade presented is quite simple and ho-
1049 mogeneous in structure. Previous approaches for at-
1050 tentive filtering, such as Itti et al. (1998) propose a
1051 more complex and heterogeneous mechanism for fil-
1052 tering. Similarly Amit and Geman (1999) propose a
1053 hierarchical structure for detection in which the stages
1054 are quite different in structure and processing. A ho-
1055 mogeneous system, besides being easy to implement
1056 and understand, has the advantage that simple tradeoffs
1057 can be made between processing time and detection
1058 performance.
1059    Finally this paper presents a set of detailed exper-
1060 iments on a difficult face detection dataset which has

1061 been widely studied. This dataset includes faces under
1062 a very wide range of conditions including: illumina-
1063 tion, scale, pose, and camera variation. Experiments on
1064 such a large and complex dataset are difficult and time
1065 consuming. Nevertheless systems which work under
1066 these conditions are unlikely to be brittle or limited to a
1067 single set of conditions. More importantly conclusions
1068 drawn from this dataset are unlikely to be experimental
1069 artifacts.

## Notes    1077

1. Supervised refers to the fact that the attentional operator is trained 1078 to detect examples of a particular class.
2. Henry Rowley very graciously supplied us with implementations 1079 of his detection system for direct comparison. Reported results 1080 are against his fastest system. It is difficult to determine from 1081 the published literature, but the Rowley-Baluja-Kanade detector 1082 is widely considered the fastest detection system and has been 1083 heavily tested on real-world problems.
3. A complete basis has no linear dependence between basis ele- 1084 ments and has the same number of elements as the image space, 1085 in this case 576. The full set of 160,000 features is many times 1086 over-complete.
4. There is a close relation to "summed area tables" as used in graph- 1087 ics (Crow, 1984). We choose a different name here in order to em- 1088 phasize its use for the analysis of images, rather than for texture 1089 mapping.
5. The availability of custom hardware and the appearance of spe- 1090 cial instruction sets like Intel MMX can change this analysis. 1091 It is nevertheless instructive to compare performance assuming 1092 conventional software algorithms.
6. In the case where the weak learner is a perceptron learning al- 1093 gorithm, the final boosted classifier is a two layer perceptron. A 1094 two layer perceptron is in principle much more powerful than any 1095 single layer perceptron.
7. Among other published face detection systems some are poten- 1096 tially faster. These have either neglected to discuss performance 1097 in detail, or have never published detection and false positive rates 1098 on a large and difficult training set.

## References    1099

Amit, Y. and Geman, D. 1999. A computational model for visual 1100 selection. *Neural Computation*, 11:1691–1715.    1101

154    *Viola and Jones*

**1102** Crow, F. 1984. Summed-area tables for texture mapping. In *Proceed-*
**1103**     *ings of SIGGRAPH*, 18(3):207–212.
**1104** Fleuret, F. and Geman, D. 2001. Coarse-to-fine face detection. *Int.*
Au: Pls. **1105**     *J. Computer Vision*.
provide **1106** Freeman, W.T. and Adelson, E.H. 1991. The design and use of steer-
vol & **1107**     able filters. *IEEE Transactions on Pattern Analysis and Machine*
page range **1108**     *Intelligence*, 13(9):891–906.
**1109** Freund, Y. and Schapire, R.E. 1995. A decision-theoretic generaliza-
**1110**     tion of on-line learning and an application to boosting. In *Com-*
**1111**     *putational Learning Theory: Eurocolt 95*, Springer-Verlag, pp.
**1112**     23–37.
**1113** Greenspan, H., Belongie, S., Gooodman, R., Perona, P., Rakshit, S.,
**1114**     and Anderson, C. 1994. Overcomplete steerable pyramid filters
**1115**     and rotation invariance. In *Proceedings of the IEEE Conference*
**1116**     *on Computer Vision and Pattern Recognition*.
**1117** Itti, L., Koch, C., and Niebur, E. 1998. A model of saliency-based
**1118**     visual attention for rapid scene analysis. *IEEE Patt. Anal. Mach.*
**1119**     *Intell.*, 20(11):1254–1259.
**1120** John, G., Kohavi, R., and Pfeger, K. 1994. Irrelevant features and
**1121**     the subset selection problem. In *Machine Learning Conference*
**1122**     *Proceedings*.
**1123** Osuna, E., Freund, R., and Girosi, F. 1997a. Training support
**1124**     vector machines: An application to face detection. In *Proceed-*
**1125**     *ings of the IEEE Conference on Computer Vision and Pattern*
**1126**     *Recognition*.
**1127** Osuna, E., Freund, R., and Girosi, F. 1997b. Training support vector
**1128**     machines: an application to face detection. In *Proceedings of the*
**1129**     *IEEE Conference on Computer Vision and Pattern Recognition*.
**1130** Papageorgiou, C., Oren, M., and Poggio, T. 1998. A general frame-
**1131**     work for object detection. In *International Conference on Com-*
**1132**     *puter Vision*.

**1133** Quinlan, J. 1986. Induction of decision trees. *Machine Learning*,
**1134**     1:81–106.
**1135** Roth, D., Yang, M., and Ahuja, N. 2000. A snowbased face detector.
**1136**     In *Neural Information Processing* 12.
**1137** Rowley, H., Baluja, S., and Kanade, T. 1998. Neural network-based
**1138**     face detection. *IEEE Patt. Anal. Mach. Intell.*, 20:22–38.
**1139** Schapire, R.E., Freund, Y., Bartlett, P., and Lee, W.S. 1997. Boost-
**1140**     ing the margin: A new explanation for the effectiveness of voting
**1141**     methods. In *Proceedings of the Fourteenth International Confer-*
**1142**     *ence on Machine Learning*.
**1143** Schapire, R.E., Freund, Y., Bartlett, P., and Lee, W.S. 1998. Boost-
**1144**     ing the margin: A new explanation for the effectiveness of voting
**1145**     methods. *Ann. Stat.*, 26(5):1651–1686.
**1146** Schneiderman, H. and Kanade, T. 2000. A statistical method for
**1147**     3D object detection applied to faces and cars. In *International*
**1148**     *Conference on Computer Vision*.
**1149** Simard, P.Y., Bottou, L., Haffner, P., and LeCun, Y. (1999). Boxlets:
**1150**     A fast convolution algorithm for signal processing and neural net-
**1151**     works. In M. Kearns, S. Solla, and D. Cohn (Eds.), *Advances*
**1152**     *in Neural Information Processing Systems*, vol. 11, pp. 571–
**1153**     577.
**1154** Sung, K. and Poggio, T. 1998. Example-based learning for view-
**1155**     based face detection. *IEEE Patt. Anal. Mach. Intell.*, 20:39–51.
**1156** Tieu, K. and Viola, P. 2000. Boosting image retrieval. In *Proceedings*
**1157**     *of the IEEE Conference on Computer Vision and Pattern Recog-*
**1158**     *nition*.
**1159** Tsotsos, J., Culhane, S., Wai, W., Lai, Y., Davis, N., and Nuflo,
**1160**     F. 1995. Modeling visual-attention via selective tuning. *Artificial*
**1161**     *Intelligence Journal*, 78(1/2):507–545.
**1162** Webb, A. 1999. Statistical Pattern Recognition. Oxford University
**1163**     Press: New York.