

Face Detection and Its Applications in Intelligent and Focused Image Retrieval

Zhongfei Zhang

Dept. of Computer Science
SUNY Binghamton
Binghamton, NY 13902
zzhang@binghamton.edu

Rohini K. Srihari Aibing Rao

Dept. of Comp. Sci. and Eng.
SUNY Buffalo
Buffalo, NY 14260
{rohini,arao}@cse.buffalo.edu

Abstract

This paper presents a face detection technique and its applications in image retrieval. Even though this face detection method has relatively high false positives and low detection rate (as compared with the dedicated face detection systems in the literature of image understanding), due to its simple and fast nature, it has been shown that this system may be well applied in image retrieval in certain focused application domains. Two application examples are given: one combining face detection with indexed collateral text for image retrieval regarding human beings, and the other combining face detection with conventional similarity matching techniques for image retrieval with similar background. Experimental results show that our proposed approaches have significantly improved the image retrieval precisions than the existing search engines in these focused application domains.

Key Words: Face Detection, Image Retrieval, Collateral Text, Similarity Matching, Focused Domain.

1 Introduction

Image retrieval plays an important role in the research of multimedia information indexing and retrieval, and the techniques for efficient and effective image retrieval have significant applications in data mining, internet/intranet search, as well as information retrieval in many unstructured, multimedia databases.

The efficiency of image retrieval concerns the response time it takes after a query is entered. Typically the efficiency issue relates to the size of the database, and more importantly, also relates to the indexing scheme. The effectiveness of image retrieval, on the other hand, refers to the relevance in content

of the retrieved documents to the query. This issue is related to the retrieval precision in the area of information retrieval. Since an effective image retrieval requires certain level of image understanding, and since image understanding typically is computationally expensive, in general, there are very few image search engines commercially available in real applications that actually conduct indexing at the image level. Instead, almost all the image search engines employ an indexing scheme through indexing the collateral text of the images. Since text indexing is much cheaper than image indexing, these image search engines are efficient in image retrieval, but typically not effective.

In order to address the effectiveness issue, certain level of image indexing must be performed. However, in the communities of computer vision and image understanding, it is well believed that there is no single algorithm that would be able to serve the general purpose of understanding the content of any images. A realistic approach to effectively indexing images in an image library to facilitate image retrieval is to first manually classify all the images into different “categories” or “domains”, and to design search engines for each of the domains. This is referred to *focused* image retrieval, i.e. image retrieval in certain domains.

In this paper, we focus our domain on images with human beings. Research shows [14] that this is probably one of the most popular image domains in many business or daily life applications, such as consumer photos, news photos, etc. In particular, this research focuses on a face detection system as well as two application areas: image retrieval of querying human beings and image retrieval of querying similar background. The application in the first area may be further classified into two scenarios: (1) querying identification of human beings in general, e.g. *find images of people in beach* (2) querying identification of specific individuals, e.g. *find images of Hillary Clinton in speech*.

In first application area, we assume that there is certain collateral textual information available accompanying the images, such as caption. This is a reasonable assumption as almost all news photos come with caption, and many consumer photos have annotation text associated with the pictures. Thus, the individual identification is approached through verification at the text level using text indexing as opposed to that at the image level using image indexing. Therefore, no face identification or recognition at the image level is intended; only face detection is conducted to verify the existence of human beings in the image. Take the example of the query *find images of Hillary Clinton in speech*. Image indexing is only performed for face detection, i.e. verification of the existence of human beings in images; the verification of the existence of the specific individual — Hillary Clinton, i.e. face recognition, is performed through verification at the text level — looking for the evidence of the presence of the name *Hillary Clinton* or related equivalents such as *First Lady*, *Mrs. Clinton*, etc. in the collateral text.

In order to facilitate an efficient image retrieval in this focused domain, it is critical to have an efficient face detection algorithm. In this paper, we first present an efficient face detection system. Then based on this face detection system, we combine it with a conventional text indexing system to conduct image retrieval for the first type of application. Next we combine this face detection system with a conventional color histogram based similarity matching system to solve for the second application of image retrieval. These methods are evaluated separately, and finally, the paper is concluded.

2 Face Detection

Identifying people in general imagery, although with its long history in literature [10], is still a difficult problem, since it requires solutions to both face detection and face recognition, which have been proven to be difficult [3]. Even though we have assumed that the face recognition part is solved by indexing in the collateral text of the image [16], the face detection part still needs to be done in general appearance. By general appearance, it is meant that the faces appearing in the images may exhibit arbitrary variations in their pose (from frontal views to different angles of side and tilt views), orientation (head inclined to the left, right, etc.), scale (from very small to very large), expression (smiling, crying, sober, etc.), contrast (from very dark to very bright), as well as background (may be very complex). Given this generality, it poses great

challenge to research. Many reported face detection systems in the literature either have restrictions on the general appearance, e.g. only frontal faces are assumed, or are expensive in computation. In this work, due to the nature of image retrieval, we cannot afford a solution with expensive computation. On the other hand, we have to handle the general appearance requirement in image retrieval. The face detection developed here satisfies the both requirements.

In the literature of face detection, since faces are 3D deformable objects, almost all the techniques are based on deformable models. Typical deformable models include algebraic models (e.g. Yullie et al [20]), neural networks (e.g. Rowley et al [15]), segmentation and grouping based models (e.g. Govindaraju [8]), color classification (e.g. Wang and Chang [19]), learning in feature space (e.g. Sung and Poggio [18] and Osuna et al [13]), and information theory (e.g. Moghadam and Pentland [11], Colmenarez and Huang [4]), as well as some combination of different “basic” methods (e.g. Oliver et al [12] used 2D blob features based on color feature space classification to detect human faces, and Burl and Perona [2] combined local feature analysis with probabilistic reasoning in the image spatial domain to locate frontal human faces in cluttered scenes). Algebraic models usually give rise to accurate segmentation of faces in an image. However, since these models are essentially based on non-linear search in a multidimensional space, they always need initial guess, which needs *a priori* knowledge of the solutions. Neural network approaches are usually robust, but they need large scale training. Segmentation and grouping based techniques normally involve many parameters, and the parameters need tuning. Color classification techniques are also very robust, but they may typically produce large false positives, and thus, need to be combined with other techniques in order to give rise to accurate detection. Actually color classification is a special case of machine learning in a color feature space. As a general machine learning approach to face detection, depending on different learning strategies, data training is always a key problem that could lead to success or failure of the system. Information theory based approaches are similar to neural network approaches and machine learning approaches, and thus suffer from the same drawbacks. The advantages, though, are potential improvement in efficiency [4].

Since color imagery is very popular in use and is very easy to obtain in these days, in this paper, we focus our research on color imagery. The face detection is approached as pattern classification in a color feature space. The detection process is accomplished in

two major steps: *Feature Classification* and *Candidate Generation*.

2.1 Feature Classification

In order to conduct pattern classification in a feature space, the first problem is to define the proper features. In this work, we use two color features: *hue* and *chrominance*. Hue refers to the basic color as it appears to the observer in the sense that certain value of hue is roughly proportional to the wavelength of its particular color [1]. The chrominance, on the other hand, is the measure of the lack of whiteness in a color [1]. Chrominance is used as a filter to reject those pixels that are saturated. The hue feature is then used to classify between face and non-face based on a Bayesian rule [6]. Intensity is not used so that the proposed method can be independent of lighting condition to a certain degree [1]. Given a particular triple of RGB values of a color image, there are many ways to transform them into hue and chrominance values. Since the method proposed by Gong and Sakauchi [7] appears to have the least singularity, it is employed in this work.

Once a transformation method is determined, the hue and chrominance ranges typical for human face (actually for human skin) were computed based on training over 100 face images of people of different races, manually collected from different WWW sites. Fig. 1 shows the hue histograms of typical White, Black, and Asian faces. It is obvious that all the histograms exhibit very similar clustered ranges of hue feature for human faces of these different races. Similar consistency is also obtained for histograms of chrominance over different races. Based on all the histograms of the training data, *a priori* probabilities for face and non-face classes may be obtained which are used in our face detection system. Then the problem of face detection is reduced to the problem of classification between face and non-face classes based on a standard Bayesian rule [6]. This classification result is a binary image as shown in Fig. 2(b), where Fig. 2(a) is the original input color image.

2.2 Candidate Generation

The binary image obtained from Bayesian classification is examined for potential candidates of faces. Presumably, those condensed, “big” clusters indicate the presence of candidates. In order to “clean out” the noise of the binary image before any clustering technique is used to obtain the real face candidates, a morphological operation [9] is applied. Then, a con-

nected component search [9] is conducted to collect all the “big” clusters.

Like all other face detection systems reported in the literature, the granularity is always a problem, i.e. how “big” a cluster is qualified for being a candidate? In the current system implemented, a threshold of 400 pixels (corresponding to a 20×20 square bounding box in area) is used to reject any smaller clusters. This means that we are not interested in detecting those very small faces in an image. It is worth noting that this threshold is a reasonable one as it is used by some other successful face detection systems reported in the literature (e.g. [15]).

Since this method essentially returns all human skin parts as face candidates rather than just faces themselves, plus there are other objects exhibiting their hue and chrominance features within the human skin ranges, the false positives are typically high. Fig. 2(c) shows the results at this step.

In order to reduce the false positives, a simple heuristic based on human face *golden ratio* is used to reject most of the false positives¹. Fig. 2(d) is the result after applying this heuristic.

3 Combining Face Detection with Collateral Text Indexing

Since image processing/understanding is in general expensive as compared with text processing, there are few commercial search engines that actually employ image level indexing techniques. Instead, images are typically indexed at the collateral text level. While this is efficient, it is not an effective approach. To give an example, a query of *Bill Clinton* was entered into *Lycos* image search engine, and Table 1 recorded the precisions returned in this experiment.

Table 1
Lycos Image Retrieval Example

Number of Top Retrieved Images	Precision
5	0%
10	20%
20	40%
30	50%

Manual examination of those irrelevant documents that were retrieved shows that they have strong relationship to the query *Bill Clinton* in the text even though in the images there is no Bill Clinton at all.

¹It is believed that for a typical human face, the ratio of the width to the height of the face is always around the magic value of $2/(1 + \sqrt{5})$, which is called the golden ratio [5].

An example is a document consisting of a picture of White House and in the collateral text the name of Bill Clinton is mentioned several times. That is why that image was ranked so high in the retrieved document list even though the image has no Bill Clinton in it. This shows that only using text indexing for image retrieval may mislead the actual retrieval. In fact, Bill Clinton is probably one of the most popular celebrities. However, even with this query, the precision is so low. Thus, it is no surprise that often how low the precision would be for a typical query like this. This shows that the text based image indexing is inappropriate and ineffective. Hence, it is necessary to employ a face detection system to help verify the existence of human faces. Note that even though the detection rate of our face detection system is not 100%, and even though it has large false positives, since the computation is very fast, it is still worth employing this system. The rationale is that even if the face detection rate drops to 50%, at least we would have half of the retrieved images that are relevant, which brings up to 50% precision. That is considered as a very high precision in real applications.

Here we apply face detection as a filter after the search based on the text indexing [16]. This way we only need to apply face detection to the retrieved documents from the first round of search of the whole database based on text indexing. Consequently, the additional computation for face detection is significantly reduced. On the other hand, due to this filtering, the precision is significantly improved. Fig. 4 shows the algorithm of this approach.

4 Background Similarity Matching

It has been shown that the conventional similarity based matching techniques do not work well for images that contain human beings [17]. On the other hand, in certain application areas such as consumer photos and news photos, it is not unusual that images contain human beings. An effective solution to this problem is to apply face detection to the image first to crop out the areas with human bodies. Then the rest of the image (i.e. the background) may be ported to the conventional similarity based matching for image retrieval based on the background similarity. Fig. 5 shows an example of retrieved images with background similarity.

5 Experimental Evaluations

5.1 Face Detection Evaluation

An evaluation for the face detection system has been conducted for 460 images with 1169 faces from different sources. These are all very general color images with significantly varied appearances of human faces. Fig. 3 shows one of the examples of the output of our face detection system. Note the significant variations in these images of lighting, race, face pose and orientation, face scale, and background. There are also a lot of faces in these images that are partially occluded and/or overlapped². Based on all these 460 images with 1169 faces collected so far for detection evaluation, the detection rate is 60.7% and the false positive rate is 56.4%. It is worth noting that this detection rate is lower than those of some dedicated systems reported in the literature of image understanding and computer vision (e.g. [15]). The reason is that there is *no* constraint imposed in this system, as opposed to many of the other systems (e.g. frontal face constraint is typically used in the systems reported in the literature). Thus, the lower detection rate is traded off by the generality and the complexity of the face appearances. The higher false positive rate is expected since this approach is based on human skin classification. When combined with other sources of information, the high false positives may typically be rejected (e.g. high false positives may typically be rejected by low text indexing confidence).

This face detection system is very fast. Since the actual time taken in detecting faces for a specific image varies with the image size, in order to show the speed information, we ran an image of 267×203 on a Sparc 10. The conversion time of RGB to the hue and chrominance features was 1.458 second; the remaining face candidate generation took 0.494 second. This shows that almost three fourth of the time is spent on color format conversion, which may be performed offline. Thus, the “real” face detection work only takes about half second for this sized image. Note that there is no attempt made for optimization of the code yet. It is clear that it is possible to achieve real-time performance using this face detection system.

²In this evaluation, we count all partially occluded or overlapped faces as ground truth faces. Since they tend to be missed or incorrectly grouped with other faces, this contributes to the lower detection rate in this evaluation.

5.2 Collateral Text Based Image Retrieval Evaluation

In this application, queries are first processed using text indexing methods [16]. This produces a ranking $P_{x_1} \dots P_{x_n}$ as indicated in Figure 4. We then eliminate those pictures p_{x_i} which do not contain faces; this information is obtained by the face detection system.

Table 2 presents the results of an experiment conducted on 198 images that were downloaded from various news web sites. The original dataset consisted of 941 images; of these, 117 were empty (white space) and 277 were discarded as being graphics. From these, we chose a subset of 198 images for this experiment. There were 10 queries, each involving the search for pictures of named individual(s); some specified contexts also, such as *find pictures of Hillary Clinton at the Democratic Convention*. Due to the demands of truthing, we report the results for one query; more comprehensive evaluation is currently underway. The 3 columns in Table 2 indicate precision rates using various criteria for content verification: (i) using text indexing alone [16], (ii) using text indexing and manual visual inspection, and (iii) using text indexing and automatic face detection. As the table indicates, using text alone can be misleading — when inspected, many of these pictures do not contain the specified face. By applying face detection to the result of text indexing, we discard photographs that do not have a high likelihood of containing faces. The last column indicates that this strategy is effective in increasing precision rates. The number in parentheses indicates the number of images in the given quantile that were discarded due to failure to detect faces.

Table 2
Retrieval Results

	Text Only	Text + Manual Insp	Text + Face Det
At 5 docs	1.0	1.0	1.0 (3)
At 10 docs	1.0	.70	1.0 (3)
At 15 docs	.80	.75	1.0 (2)
At 30 docs	.77	.67	NA

Sample output is shown in figures 6 and 7. Figure 6 illustrates the output based on text indexing alone. The last picture illustrates that text alone can be misleading! Figure 7 illustrates the re-ranked output based on results from face detection. This has the desired result that the top n images all are relevant. However, a careful examination reveals that due to the face detector's occasional failure to detect faces in images, relevant images are inadvertently being discarded. Thus, this technique increases precision, but lowers recall. However, if the source of images is the WWW, this may not be of concern.

5.3 Background Similarity Matching Evaluation

In order to evaluate the proposed approach to background similarity matching, we first set up an evaluation database consisting of 58 pure scenery images (which are typically good for conventional similarity based matching) and another 66 images from a consumer database which contain human faces in each of the images (and thus are *not* good for the conventional similarity based matching). Then we ran our algorithm in this combined database, and the averaged precision-recall is shown in Fig. 8. Note that the performance of this proposed approach is about 20% improvement over that of a conventional approach on this database.

6 Conclusion

We have proposed an efficient and effective method for face detection from color images. We have shown that even though this method is not a perfect solution to face detection in general (i.e. with relatively high false positives and low detection rate as compared with several dedicated face detection systems in the image understanding literature), this face detection system is good enough to be applied in image retrieval in certain popular domains such as consumer photos or news photos, due to its appealing efficiency. We also have shown two application examples of this face detection system in querying and retrieving images, one being image retrieval with indexed collateral text, and the other being image retrieval with background similarity. Evaluation results show that our proposed methods have significantly improved the precisions of the retrieval.

References

- [1] D.H. Ballard and C.M. Brown. *Computer Vision*. Prentice-Hall, Inc., 1982.
- [2] M.C. Burl and P. Perona. Recognition of planar object classes. In *Proc. CVPR*, 1996.
- [3] R. Chellappa, C.L. Wilson, and S. Sirohey. Human and machine recognition of faces: A survey. *Proc. of the IEEE*, 83(5), 1995.
- [4] A. Colmenarez and T. Huang. Face detection with information-based maximum discrimination. In *Proc. CVPR*, 1997.

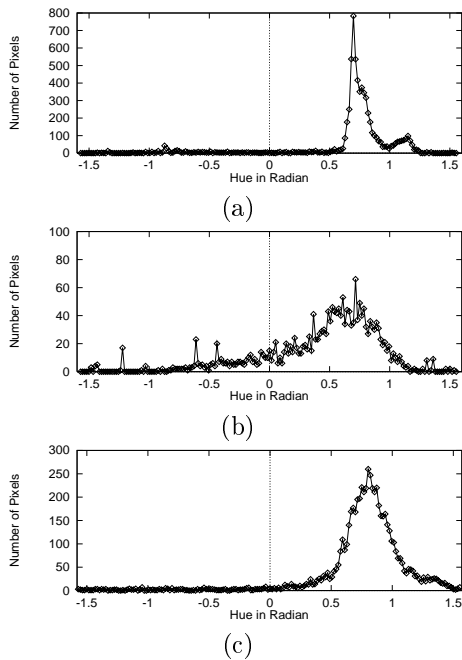


Figure 1: Typical hue histograms for different races (a) White (b) Asian (c) Black.

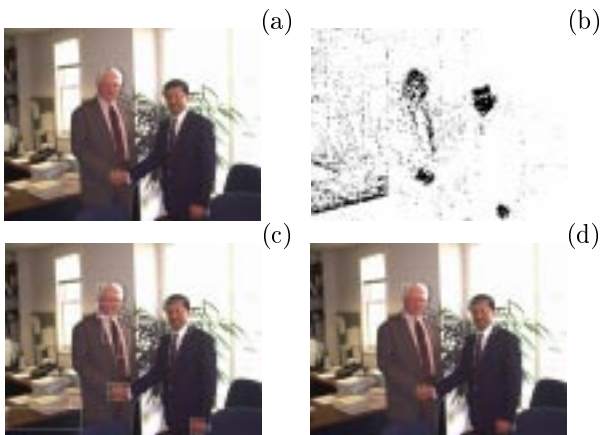


Figure 2: An example that shows the face detection results. (a) original image (b) after Bayesian classification (c) candidate generation before heuristic filtering (d) final candidates after heuristic filtering.



Figure 3: One example of the face detection results. Note the complexity of the appearances of the faces, and the variations of the scale, race, and lighting change of the faces in the image. Also note the false positives and negatives.

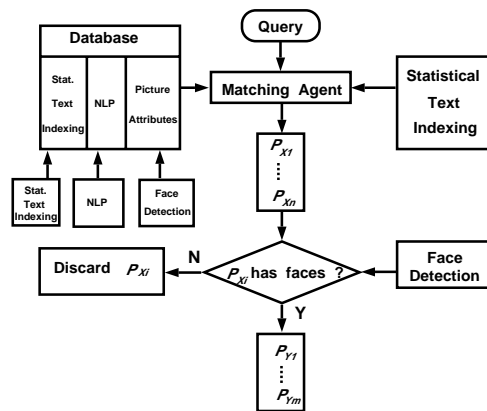


Figure 4: Single-ranking strategy for combining text and image information in retrieval.

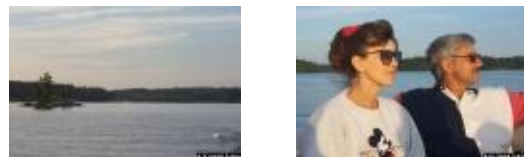


Figure 5: A pair of images that share a similar background. Due to the presence of the people (as foreground), the conventional similarity based matching techniques may fail to retrieve them as a pair of “similar” images.

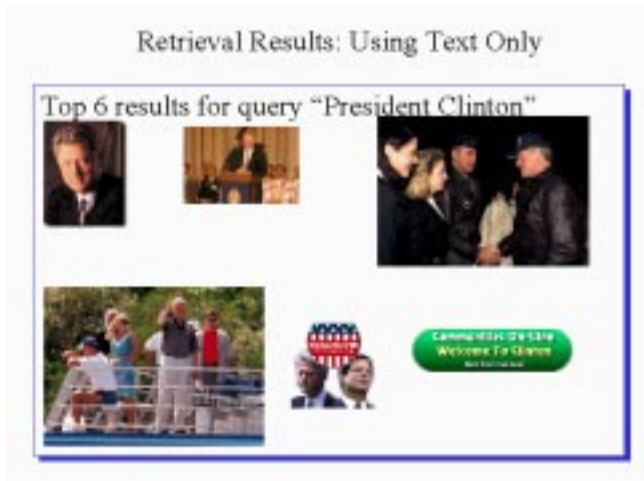


Figure 6: Top 6 images based on text indexing alone.

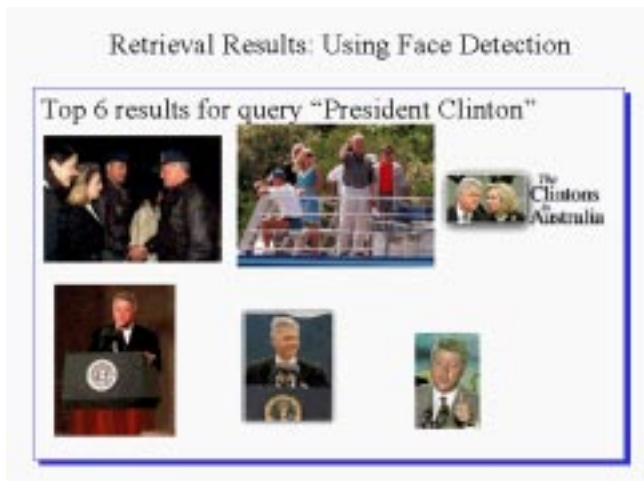


Figure 7: Top 6 images based on combining text indexing and face detection.

Precision-Recall of Different Matching Schemes on TOTAL Data Set

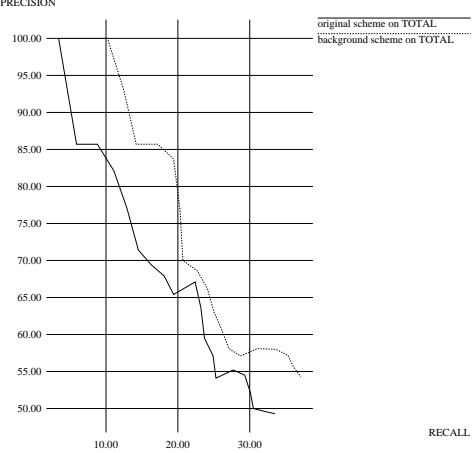


Figure 8: Average Precision-Recall for different matching schemes on a combined database consisting of purely scenery images and consumer images. The background matching scheme greatly improves the performance on this database. Number of retrieved images for each query is 20.

- [5] L.G. Farkas and I.R. Munro. *Anthropometric Facial Proportions in Medicine*. Charles C. Thomas, 1987.
- [6] K. Fukunaga. *Introduction to Statistical Pattern Recognition, 2nd Ed.* Academic Press, 1990.
- [7] Y. Gong and M. Sakauchi. Detection of regions matching specified chromatic features. *Computer Vision and Image Understanding*, 61(2), 1995.
- [8] V. Govindaraju. Locating human faces in photographs. *International Journal of Computer Vision*, 19(2), 1996.
- [9] R.M. Haralick and L.G. Shapiro. *Computer and Robot Vision*. Addison-Wesley, 1992.
- [10] T. Kanade. *Picture Processing System By Computer Complex and Recognition of Human Faces*. Ph.D. Thesis, Kyoto University, 1973.
- [11] B. Moghaddam and A. Pentland. Maximum likelihood detection of faces and hands. In *Proc. Workshop on Automatic Face and Gesture Recognition*, 1995.
- [12] N. Oliver, A. Pentland, and F. Berard. Lafter: Lips and face real time tracker. In *Proc. CVPR*, 1997.
- [13] E. Osuna, R. Freund, and F. Girosi. Training support vector machines: an application to face detection. In *Proc. CVPR*, 1997.

- [14] D.M. Romer. Research agenda for cultural heritage on information networks: Image and multimedia retrieval. *Kodak Internal Tech Report*, 1993.
- [15] H.A. Rowley, S. Baluja, and T. Kanade. Neural network-based face detection. *Trans. on PAMI*, 20(1), 1998.
- [16] R.K. Srihari and Z. Zhang. Exploiting multimodal context in image retrieval. *Library Trends*, July 1999.
- [17] R.K. Srihari and Z. Zhang. Image background search: Combining object detection techniques into content-based similarity image retrieval(cbsir) systems. In *Proc. of IEEE International Workshop on Content-Based Access of Image and Video Libraries*. IEEE Press, June 1999.
- [18] K. Sung and T. Poggio. Example-based learning for view-based human face detection. *Trans. on PAMI*, 20(1), 1998.
- [19] H. Wang and S-F Chang. A highly efficient system for automatic face region detection in mpeg video. *Trans. Circuits and systems for video technology*, 7(4), 1997.
- [20] A. Yullie, D. Cohen, and P. Hallinan. Feature extraction from faces using deformable templates. In *Proc. CVPR*, 1989.