

Regions-of-Interest and Spatial Layout for Content-Based Image Retrieval

Baback Moghaddam
Henning Biermann
Dimitris Margaritis

TR-2000-35 November 2000

Abstract

To date most “content-based image retrieval” (CBIR) techniques rely on *global* attributes such as color or texture histograms which tend to ignore the spatial composition of the image. In this paper, we present an alternative image retrieval system based on the principle that it is the *user* who is most qualified to specify the query “content” and *not* the computer. With our system, the user can select multiple “regions-of-interest” and can specify the relevance of their spatial layout in the retrieval process. We also derive similarity bounds on histogram distances for pruning the database search. This experimental system was found to be superior to global indexing techniques as measured by statistical sampling of multiple users’ “satisfaction” ratings.

This work may not be copied or reproduced in whole or in part for any commercial purpose. Permission to copy in whole or in part without payment of fee is granted for nonprofit educational and research purposes provided that all such whole or partial copies include the following: a notice that such copying is by permission of Mitsubishi Electric Information Technology Center America; an acknowledgment of the authors and individual contributions to the work; and all applicable portions of the copyright notice. Copying, reproduction, or republishing for any other purpose shall require a license with payment of fee to Mitsubishi Electric Information Technology Center America. All rights reserved.

Publication History:-

1. First printing, TR-2000-35, November 2000

Regions-of-Interest and Spatial Layout for Content-Based Image Retrieval

Baback Moghaddam¹, Henning Biermann² and Dimitris Margaritis³

¹ Mitsubishi Electric Research Laboratory, Cambridge MA, USA

² Department of Computer Science, New York University, USA

³ Department of Computer Science, Carnegie Mellon University, USA

Abstract. To date most “content-based image retrieval” (CBIR) techniques rely on *global* attributes such as color or texture histograms which tend to ignore the spatial composition of the image. In this paper, we present an alternative image retrieval system based on the principle that it is the *user* who is most qualified to specify the query “content” and *not* the computer. With our system, the user can select multiple “regions-of-interest” and can specify the relevance of their spatial layout in the retrieval process. We also derive similarity bounds on histogram distances for pruning the database search. This experimental system was found to be superior to global indexing techniques as measured by statistical sampling of multiple users’ “satisfaction” ratings.

1 Introduction

Current content-based image retrieval systems often rely only on *global* image characteristics such as color and texture histograms. While these simple global descriptors are fast and often do succeed in partially capturing the essence of the user’s query, they more often fail due to the lack of higher-level knowledge about what exactly was of interest to the user in the query image - *ie.*, the user-defined content. Figure 1 illustrates the multiplicity of user-defined “content” in a single image. The goal of this research was to develop and test a new technique for image retrieval using *local* image representations in a bottom-up fashion. Our localized representations can be easily grouped into multiple user-specified “regions-of-interest” and constrained to preserve their relative spatial configuration during retrieval. We posit that this leads to a more *user-centric* and thus a more powerful search engine.

The fact that spatial information is a critical component of image description and subsequent matching has only recently been addressed by researchers in the field. Recently, the community has witnessed a gradual shift towards spatially-encoded image representations (see Smith [13]). Current techniques range widely from fixed image partitioning as in the “ImageRover” system of Sclaroff *et al.* [12], to highly local characterizations like the “color correlograms” of Huang *et al.* [6]. Somewhere in between these two extremes, one can find various techniques which deal with “regions” or “blobs”. For example, the “configural templates” of Lipson *et al.* [8] specify a class of images (*e.g.*, snow-capped mountain scenes) by means of photometric and geometric constraints on pre-defined image regions. Other techniques use automatic blob segmentation and description, as in Howe’s “percentile blob” technique [5] or the more sophisticated “Blobworld” segmentation system of Carson *et al.* [1].

Our system differs from above in one key aspect: there are no pre-segmented regions. Instead, the user defines “blobs” or “regions-of-interest” (ROIs) directly on a query image (and implicitly their relative spatial configuration) in order to better communicate to the search engine the intended “content” (which could possibly represent only a subset or partial aspect of the query image selected).⁴ The disadvantage of this scheme, however, is that region-matching and subsequent database indexing must be done in an online fashion and moreover in “interactive-time” to be tolerated by the user. Aggressive search pruning and database re-organization does, however, alleviate this problem to some extent. The advantage, on the other hand, is that the user is not limited to working with the available set of pre-defined blobs, as in “Blobworld” [1].

2 Representation and Similarity

An image retrieval system is completely defined by two basic specifications: representation (of features) and a corresponding similarity metric (for comparison of features or their distributions). Familiar examples include global color histograms for color composition and various “texture” measures (*e.g.*, typically the

⁴ Malki *et al.* [9] have recently proposed a similar technique, which avoids pre-segmentation by means of multi-resolution quad-trees.



Fig. 1. Multiple “contents” in a single image. Photo courtesy of Philip Greenspun.

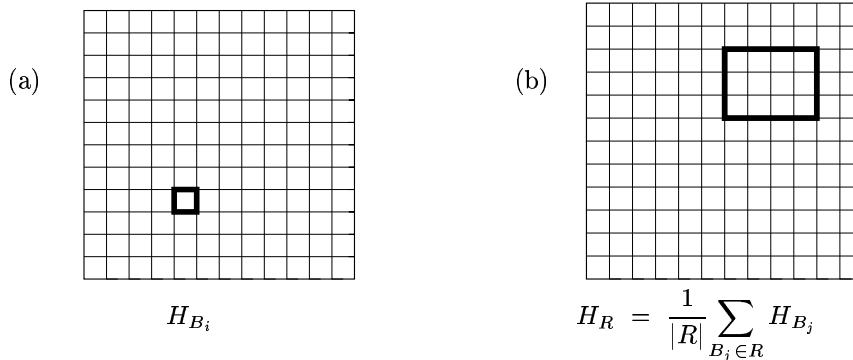


Fig. 2. (a) An image block B_i with corresponding histogram H_{B_i} (b) A region R composed of individual blocks B_j and its “pooled” histogram H_R

output of a set of linear filters at multiple scales). Further refinements could include color coherency and higher-level textural attributes such as granularity, periodicity or directionality. A key component, as important (if not more) than the choice of representation, is the similarity metric used (typically a Euclidean norm on a feature vector or histogram).

While global similarity matching has certain desirable properties (*eg.*, invariance to translation or rotation) it fails to capture the spatial layout and structure of the image. Moreover, what the user typically thinks of as the “content” is seldom captured by the whole image or its global properties. Therefore, it is better to let the user identify the regions in the image which he/she is interested in (the “content”), with the possibility of specifying the spatial layout as a search constraint. This demands a *local* representation at the finest resolution possible, which can be readily grouped into larger regions and perhaps even integrated into a single global description. Before we present the details of our system design, we should point out that our choice of feature representation (and similarity metrics) was made deliberately generic, so as to focus attention on the key proposition of this paper which is the incorporation of spatial similarity in matching multiple ROIs in user queries.

2.1 Local Feature Distributions

Our system divides the image into an array of 16-by-16 pixel blocks wherein each pixel yields a LUV color coordinate and three texture measurements; edge strength: $\log(G_x^2 + G_y^2)$, Laplacian: $G_{xx} + G_{yy}$ and edge orientation: $\arg(G_x, G_y)$, where G_x and G_{xx} are the 1st and 2nd derivatives of a Gaussian filter with specified scale σ . In our experiments, two separate scale parameters were used: $\sigma = 1$ and $\sigma = 2$, yielding two sets of (“independent” or at least uncorrelated) texture measurements. This particular texture representation scheme was selected based on the robust image matching results obtained by Schiele & Crowley [11], but other texture features could also have been used.

Estimates of the *joint distribution* of the features for color and texture were obtained non-parametrically by means of a joint 3D histogram in LUV color space (implemented with 5-by-5-by-5 bins) and a joint 3D histogram of edge magnitude, Laplacian and orientation (implemented with 4-by-7-by-4 bins), computed

at two octave scales. The edge strength was quantized (classified) into only 4 values: {no edge, weak edge, average edge, strong edge}. Similarly, the edge orientation was classified into 4 values corresponding to {horizontal, vertical, diagonal left, diagonal right}. Note that in both histograms, the total number of bins is about 120 and given the 256 pixels in a 16-by-16 block, we average out to roughly 2 observations per bin. To aid the estimation process, we also used Bayesian *m-estimates* [4] in counting hits, using database-derived prior distributions in order to balance the tradeoff between prior belief and the observed data.

2.2 Histogram Similarity Measures

We implemented and tested 3 different histogram similarity measures for our data representation: Histogram Intersection [14], Chi-squared statistic [10] and Bhattacharyya distance [3], each of which has a well-defined probabilistic interpretation (in contrast to Euclidean distance norms on histograms, which in our opinion are hard to justify, despite their prevalent use). We validated and compared the performance of these 3 measures on the **Vis-Tex** database [7] with a 58-class texture classification task and found that simple nearest-neighbor classification (using the above similarity measures) yielded acceptable performance (88-90% accuracy). We found no statistically significant difference between the 3 measures to justify selecting one over the other and all 3 were made available to the user in the browser interface.

3 Region Matching

These non-parametric densities represent *local* color and texture and due to the additive property of histograms, can be easily combined (summed) to form densities for larger image blocks, including the entire image at which point they become identical to global histograms. When the user specifies a region of interest, its underlying block histograms are "pooled" to represent a "meta-block" histogram as illustrated in Figure 2. A region is then used to index into the database, where an online search for the best matching region (of the same size) is conducted using the aforementioned similarity metrics. Multiple region queries are processed in parallel and the best region match scores are then combined (usually by summation) to determine the final visual similarity ranking.

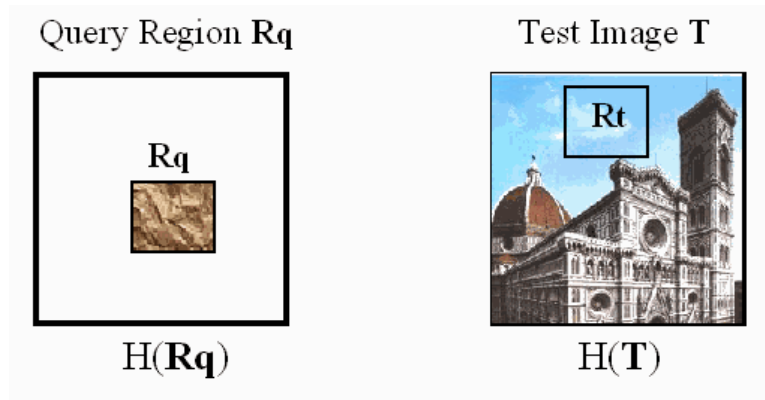


Fig. 3. A query R_q and its candidate best matching region R_t in a database image T .

To speed up the online search, the entire database is first pruned to find a small subset (typically 5-10%) of "compatible" images with fast *global* histogram indexing using a *branch and bound* algorithm. The key observation is that we can find lower bounds for the query distances using pooled-ROI and global histogram comparisons. Consider a query region histogram H_R^q and a database test image with candidate regions H_R^t and a global histogram H^t , as shown in Figure 3. It can be easily shown that with a distance metric based on histogram intersection (distance denoted here by \perp) one can compute an exact lower bound on the region-to-region distance

$$H_R^q \perp H_R^t \geq H_R^q \perp H^t \quad \dots \quad \text{for any } H_R^t \quad (1)$$

Since similarity (inverse distance) is simply a count of the number of pixels in common, the similarity between the query and the global test image can never be smaller than the similarity between the query

region and a corresponding subregion in the test image. For other metrics such as Chi-squared and Bhattacharyya, exact lower bounds are difficult to compute. Nevertheless, approximate lower bounds similar to Equation 1 can still be of practical use. Ultimately, no matter what comparison metric is used, there is no reason to search a test image whose global histogram has little in common with the query region. For example, consider the futility of searching for a query region of black and yellow zebra stripes in a (mostly blue) underwater image of a school of sharks. Once the database is sufficiently pruned, we search for all combinations of regions in the target image in order to find the best matching regions (this is potentially very slow but the user is expected to specify only few ROIs).

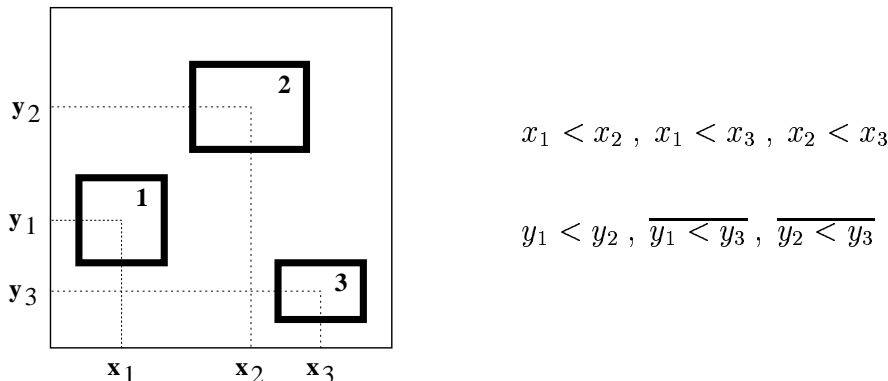


Fig. 4. Three regions and the complete set of binary relationships corresponding to their spatial configuration.

4 Spatial Constraints

In addition to querying by visual similarity, the user also has the option of specifying whether the selected regions should maintain their respective spatial configuration in the retrieved matches. We considered and investigated various techniques for spatial representation and matching, including elastic spring models and graph matching. But in the end we opted for a much simpler formulation based on the consistency of binary relations on the centroid coordinates of the regions, as illustrated in Figure 4. Given the user-defined query Q , consisting of n regions, its spatial configuration similarity to a candidate configuration T (with corresponding best matching regions), $S(Q, T)$ is

$$S(Q, T) = \sum_{i=1}^{n-1} \sum_{j=i+1}^n f(x_i^t - x_j^t) \text{sign}(x_i^q - x_j^q) + f(y_i^t - y_j^t) \text{sign}(y_i^q - y_j^q) \quad (2)$$

where x_i^q and x_i^t are the region centroid coordinates of the query Q and candidate T , respectively. The function f is a bipolar sigmoid (hyperbolic tangent) and its product with the sign function will essentially result in a “fuzzy” or “soft” count of the total number of satisfied constraints (in the set of binary relations) between Q and T . The scale parameter of the sigmoid function can be adjusted to specify how strictly a binary constraint is imposed (in the limit f can be made into a sign function as well).⁵ We note that this formulation is an approximate similarity measure as it assumes that the x and y coordinates of a region can be treated independently in determining the correct spatial relationship of two regions (it is also *not* rotation-invariant). Nevertheless, we have found it to be quick and easy to compute and quite adequate in measuring similarity of spatial configurations. Finally we note that the spatial similarity score is combined (typically by weighted summing) with the visual similarity score of all the regions to obtain a single final score by which the candidate entries in the database are ranked.

⁵ This formulation is related to the technique of “2D strings” proposed for iconic indexing by Chang et al. [2]

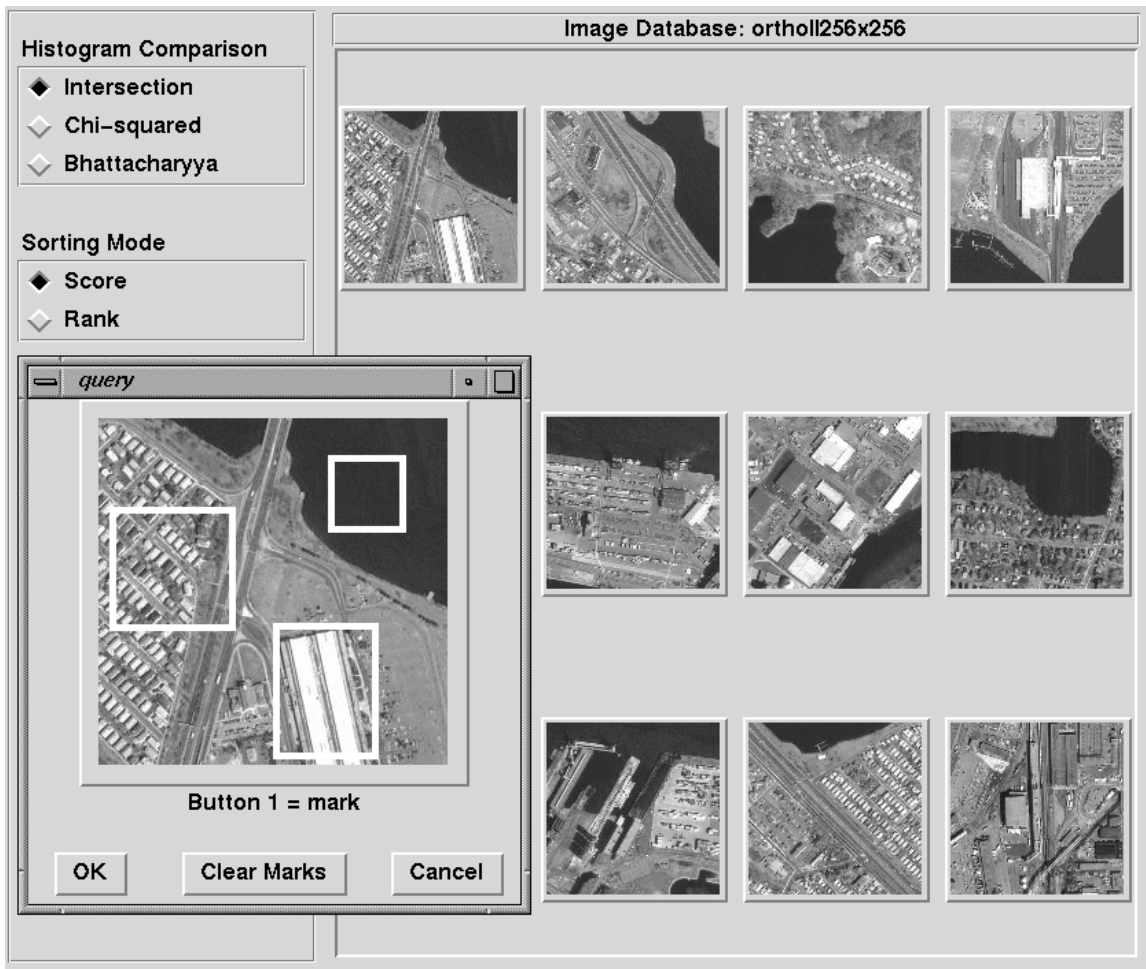


Fig. 5. An example of a multiple ROI query with a database of B&W aerial imagery.

5 Experimental Results

One of the unfortunate aspects of our user-defined multiple ROI query method is that no automatic image self-matching is possible in order to perform large classification and retrieval experiments to quantify performance. Our technique is inherently *interactive* and *user-based*, thus requiring a human in the testing loop. In other words, “content” is no longer defined by the unique and fixed global attributes of database images, but rather by a myriad of user-defined queries all of which can exist within a single image.

Therefore, the only sensible performance measure is one that quantifies the user’s overall “satisfaction” with the retrieved matches. Our experimental design was simple: 31 naive users were instructed in the basic operations of the multiple ROI query interface and asked to perform a minimum of 20 different queries on various databases.⁶ Each user-defined region-based retrieval was presented with the resulting *global* search with the same query image. The user then had to select (forced choice) which set of retrievals (local or global) captured the “essence” of their intended query content. From this sample of more than 600 (mostly independent) selections the average percentage of acceptable local first-rank matches was found to be 88% ($F_{1,30} = 14.8$, $p \leq 0.05$). This indicates that local searches were in fact favored over global ones (50% would indicate no apparent preference for local *vs.* global).

Figure 5 shows an example query in our browser, running on a database of GIS Orthophoto Imagery of the state of Massachusetts (available at <http://ortho.mit.edu>). The smaller window in the lower left allows

⁶ Our collection of $O(10^4)$ images consists of separate databases of Corel stock photos, CD covers (from Amazon.com), GIS aerial imagery, 2D MRI medical images and a large and varied assortment of images gathered by crawling the web.

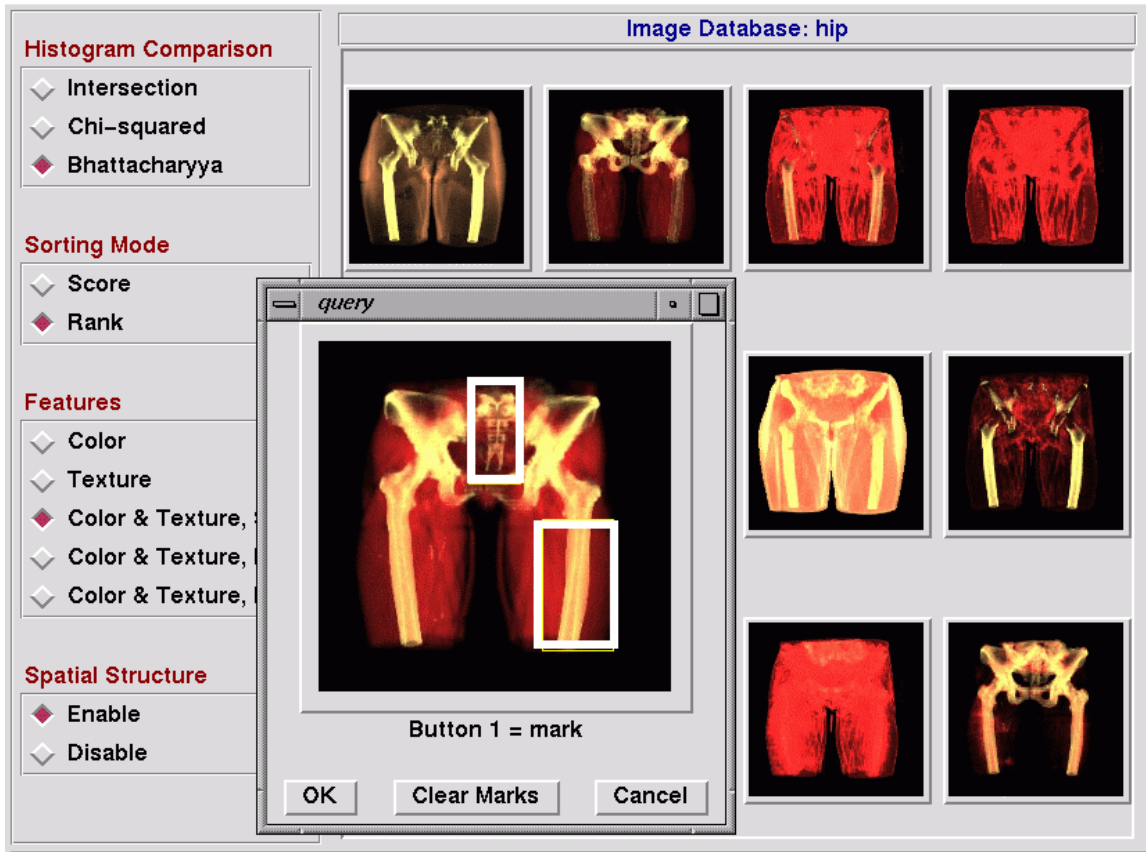


Fig. 6. An example of a multiple ROI query with a database of MRI imagery.

the user to graphically define and edit (in this case) three regions corresponding roughly to “dense urban row housing”, “water” and “factory” region types (note that these “classes” are entirely user-defined). The user can either retrieve images which respect the spatial configuration of the query, or alternatively, disable spatial scoring to simply retrieve images containing similar types of regions. Figure 6 shows an example of multi-region querying with a database of medical (MRI) images of the human hip.

6 Discussion

Currently the online search for individual regions is computationally intensive and more sophisticated pruning strategies should be implemented in order to avoid searching every region of every image in the database. Global histogram indexing is partially effective in pruning the database size down to a reasonably small candidate set. Furthermore, search schemes exploiting hierarchical database organization (based on global and/or local features) should significantly decrease the size of the candidate set and hence the search time.

Another speed-up possibility is to immediately reject candidate images based on partial spatial configurations (*e.g.*, if the best match for region 2 is already on the wrong side of region 1, reject the current image). The total computation time is the sum of the time used for the pruning (linear with a small constant) and subsequent matching in the candidate set (exponential in the number of ROIs and linear in the size of the candidate set). While it may not be possible to rival the speeds of retrieval engines with pre-segmentation — like “Blobworld” [1] — we believe our system offers the flexibility of online user-designed queries, thus leading to more accurate representations of “content”.

Finally, our system should be useful not only for general image retrieval, but also for domain-specific databases such as the GIS aerial imagery shown in Figure 5 where global descriptors are a rather poor representation of content. ROI querying can also benefit medical applications, where both appearance and spatial factors can play a significant diagnostic role (Figure 6).

References

1. C. Carson, S. Belongie, H. Greenspan, and J. Malik. Region-based image querying. In *Proc. IEEE Workshop on Content-Based Access of Image and Video Libraries*, June 1997.
2. S. Chang, Q. Shi, and S. Yan. Iconic indexing using 2-D strings. *IEEE Trans. on Pattern Analysis & Machine Intelligence*, 9(3):413–428, May 1987.
3. K. Fukunaga. *Introduction to Statistical Pattern Recognition*. Academic Press, 1971.
4. A. Gelman, J. B. Carlin, H. S. Stern, and D. B. Rubin. *Bayesian Data Analysis*. Chapman & Hall, 1995.
5. N. Howe. Percentile blobs for image similarity. In *Proc. IEEE Workshop on Content-Based Access of Image and Video Libraries*, June 1998.
6. J. Huang, S. K. Kumar, M. Mitra, W. Zhu, and R. Zabih. Image indexing using color correlograms. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 1997.
7. MIT Media Laboratory. Vistex vision texture database, 1995.
8. P. Lipson, E. Grimson, and P. Sinha. Configuration based scene classification and image indexing. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 1997.
9. J. Malki, N. Boujemaa, C. Nastar, and A. Winter. Region queries without segmentation for image retrieval by content. In *Third International Conference on Visual Information Systems (Visual'99)*, Amsterdam, June 1999.
10. A. Papoulis. *Probability, Random Variables, and Stochastic Processes*. McGraw Hill, 1991.
11. B. Schiele and J. L. Crowley. Object recognition using multidimensional receptive field histograms. In *European Conference on Computer Vision*, volume 1, pages 610–619. ECCV, April 1996.
12. S. Sclaroff, L. Taycher, and M. La Cascia. Imagerover: A content-based image browser for the world wide web. In *Proc. IEEE Workshop on Content-Based Access of Image and Video Libraries*, June 1997.
13. J. R. Smith. *Integrated Spatial and Feature Image Systems: Retrieval, Compression and Analysis*. PhD thesis, Columbia University, 1997.
14. M. Swain and D. Ballard. Color indexing. *International Journal of Computer Vision*, 7(1):11–32, 1991.