MORPHOLOGICAL COLOR SIZE DISTRIBUTIONS FOR IMAGE CLASSIFICATION AND RETRIEVAL

Jesús Angulo and Jean Serra

{angulo,serra}@cmm.ensmp.fr ENSMP, Centre de Morphologie Mathématique, 35 rue Saint-Honoré, Fontainebleau, FRANCE

ABSTRACT

Current content-based image retrieval techniques can typically perform efficient and effective searches on heterogeneous image databases. This contribution deals with an approach based on the integration of color and texture description which is applied to a very homogeneous database: a blood image bank. The content of images is very similar and therefore it becomes imperative to use very precise descriptors: the color is described by classical color distributions (histograms) and for the texture, we introduce the morphological color size distributions. The similarity is measured by computing distance metrics between histograms. In order to increase the accuracy of retrieval, the results of color-based and texture-based retrieval are integrated by combining the associated dissimilarity values. The effects of different integration methods on classification performance are explained by means of experimental tests in a database of 123 cell images (leukocyte color images). After learning processing, where different feature selection and classifier definition alternatives are tested, a definitive integrated approach is proposed (precision 94.44%).

1. INTRODUCTION

Due to the large amount of images that are generated by innumerable applications, image databases are being created and used. An efficient and automatic procedure is required for indexing and retrieving images from these databases. The content-based image retrieval involves describing an image by a digital signature and then matching the query image to the most similar images within the database according to the resemblance of their signatures. Traditionally, the content description is done according to the notions of color, shape and texture: QBIC [6], Photobook [13], CANDID [10], Pic-ToSeek [7], etc. Image databases can generally be classified in either two categories. The first one concerns homogeneous databases, they contain images of the same object class (faces, fingerprints, biomedical images, etc.). The second category includes databases of heterogeneous images (photographic stock, WWW, etc.).

In this contribution, we present a method for classifying and querying images based on the integration of color and texture description which is applied to a very homogeneous database: a blood image bank [5]. In this kind of databases, the content of images is very similar and therefore it becomes imperative to use very precise descriptors. Specifically, our goal is to develop a technique for classifying and indexing color leukocyte images into five categories [1]. The color is described by classical color distributions (histograms) and for the texture, we introduce the morphological color size distributions. In both cases, the dissimilarity is computed using distance metrics. The effects of different integration methods on classification performance are explained by means of experimental tests in a database of 123 cell images. Similar image-based approaches for automated leukocyte classification using morphological tools have been previously proposed [19].

The rest of the paper is structured as follows. First, the similarity between statistical pattern recognition and retrieval is presented in Section 2. In Section 3 a reminder of histogram definition and distance metrics is included. We continue in Section 4 with a presentation of the morphological color size distribution and in Section 5 with a brief discussion about the integration of color and texture distances. Then, the results of application to leukocyte image indexing are discussed. Finally, conclusions are given in Section 7.

2. BACKGROUND

2.1. Image attributes

In order to classify images according to image templates or in order to retrieve images, we must be able to efficiently compare two images to determine their resemblance. Let X be an image, and \mathbf{x} a d-dimensional feature vector of attributes of X. Let f represent a mapping from the image space onto the d-dimensional feature space $f : X \to \mathbf{x} = (x_1, x_2, \dots, x_d)$ where d is the number of features used to represent the image (each feature may have different dimensions). f is called the *feature extraction*. An efficient matching scheme depends on the structure (such as scalar, graphs, histograms, pdf's, etc.) and on the discriminatory information contained in the extracted features by f.

2.2. Classification and retrieval

In statistical pattern recognition (SPR), a pattern is represented by a set of d features or attributes, viewed as a ddimensional feature vector. The decision making process in SPR can be summarised as follows: a given pattern is assigned to one of c categories $\omega_1, \omega_2, \cdots, \omega_c$ based on the vector of d feature values $\mathbf{x} = (x_1, x_2, \cdots, x_d)$ [4]. The pattern vector x belonging to class ω_i can be viewed as an observation drawn randomly from the class-conditional probability function $p(\mathbf{x}|\omega_i)$. The "optimal" Bayes decision rules for minimising the risk can be stated as follows: assign input pattern x to class ω_i for which the conditional risk $R(\omega_i | \mathbf{x}) =$ $\sum_{j=1}^{c} L(\omega_i, \omega_j) P(\omega_j | \mathbf{x})$ is minimum, where $L(\omega_i, \omega_j)$ is the loss incurred in deciding ω_i when the true class is ω_j and $P(\omega_i | \mathbf{x})$ is the posterior probability. In the case of the 0/1loss function, the conditional risk becomes the conditional probability of misclassification $L(\omega_i, \omega_j) = 0$ if i = j or $L(\omega_i, \omega_j) = 1$ if $i \neq j$. For this choice of loss function, the Bayes decision rule can be simplified as follows (also called the maximum a posteriori (MAP) rule): assign input pattern **x** to the class ω_i if

$$P(\omega_i | \mathbf{x}) > P(\omega_j | \mathbf{x}) \quad for \ all \quad j \neq i.$$
(1)

The class-conditional densities are usually not known in practice and must be learnt from the available training patterns. Template matching is a natural approach to pattern classification: assign patterns to the most similar template. Template matching can easily be expressed mathematically: let **x** be the feature vector for the unknown input, and let $\mathbf{m}_1, \mathbf{m}_2, \cdots, \mathbf{m}_c$ be templates for the *c* classes. A minimum-error classifier computes $||\mathbf{x} - \mathbf{m}_k||$ for k = 1 to *c* and chooses the class for which this error is minimum, we call this a minimumdistance classifier [9].

In most of search and retrieval algorithms, a distance measure between the image attributes is used to rank the database images in ascending order of their distances to the query image, which is assumed to correspond to a descending order of similarity. The system returns the top-ranked images that are most similar to the query image. Using a SPR paradigm, retrieval operation can be interpreted as a two-class pattern classification approach. In doing so, we define two classes, the relevance class ω_A and the irrelevance class ω_B , in order to classify image pairs as similar or dissimilar. Let d be the distance between the query image and the image database image. The image pair is assigned to the relevance class if $P(\omega_A|d) > P(\omega_B|d)$ [2].

In short, we can say that in applications involving classification, the system assigns each image to the closest semantic class in the database. In image retrieval applications, the system retrieves the most similar images to the query. The same feature descriptors and distances can be used to classify or to retrieve.

3. COLOR DISTRIBUTIONS

Color is an important attribute for image retrieval: color is an intuitive feature for which it is possible to use an effective and compact representation. Color spaces provide the method to manipulate colors. In the field of image processing and computer graphics, a lot of color models have been proposed [14]. In this study we have worked with three wellknown color spaces: *RGB*, *HSV* and *Lab*. Regardless of the color space, color information in an image can be represented by a single 3-D histogram or three separate 1-D histograms. In a simplified way, we worked with the histogram of each color component, i.e. the *color histogram* is done by three histograms. These color representations are invariant under rotation and translation of the image. A suitable normalisation also provides scale/size invariance.

3.1. Histogram definition

Given an image f, of size M by N pixels, characterised by the color c at location (i, j), i.e. c = f(i, j), the first-order color distribution or histogram of the color set C is given by

$$h_f(\mathbf{c}) = \frac{1}{MN} \sum_{i=0}^{M-1} \sum_{j=0}^{N-1} \delta(f(i,j) - \mathbf{c}), \quad \forall \mathbf{c} \in \mathcal{C}.$$
(2)

In the equation above $\delta()$ is the unitary impulse function. We notice that the $h_f(c)$ values are normalised in order to sum to one [20]. This normalisation allows to compare regions or images of varying size. The value of each bin is thus the number of image pixels having the color c, or, after normalisation by MN, the probability that the color c appears in the image. If the cardinal of C is n (n is the number of bins), the histograms can be represented as feature vectors in a n-dimensional space. We can define this R^n space as the histogram space \mathcal{H}^n .

3.2. Histogram distance metrics

The histograms are feature vectors which are used as image indices. A distance measure is used in the histogram space to measure the similarity of two images. Formally, the histogram space \mathcal{H}^n is considered as a metric space and the histograms are points of the space if for every pair of histograms h_f and h_g a corresponding number $d(h_f, h_g)$ can be found, called the distance metric between the two points, which satisfies the following properties: non-negativity $(d(h_f, h_g) \ge$ 0), identity $(d(h_f, h_g) = 0 \Leftrightarrow h_f = h_g)$, commutativity $(d(h_f, h_g) = d(h_g, h_f))$ et triangular inequality $(d(h_f, h_g) \le$ $d(h_f, h_k) + d(h_k, h_g))$.

The similarity, or dissimilarity via the distance, is computed between a *pattern or query histogram* $h_Q \in \mathcal{H}^n$ and a *template or image database histogram* $h_I \in \mathcal{H}^n$. There are several mappings $d : \mathcal{H}^n \times \mathcal{H}^n \to [0, \infty)$ satisfying the metric properties which were used in the literature [16] in order to measure the relevance of the histogram dissimilarity. We focused on four distance metrics:

• Normalised histogram intersection, d_{\cap} . Normalised histogram intersection is equivalent to the use of the sum of absolute differences or city-block metric. The normalised histogram intersection distance is defined by

$$d_{\cap}(h_Q, h_I) = 1 - \frac{\sum_{i=0}^{n-1} \min(h_Q(i), h_I(i))}{\min(\sum_{i=0}^{n-1} h_Q(i), \sum_{i=0}^{n-1} h_Q(i))}.$$
(3)

Note that d_{\cap} is robust to substantial object occlusion and cluttering [18][7].

• *Histogram Euclidean distance*, d_E . The classical histogram Euclidean distance is defined as

$$d_E(h_Q, h_I) = \sqrt{\sum_{i=0}^{n-1} (h_Q(i) - h_I(i))^2}.$$
 (4)

• *Histogram* χ^2 *-metric distance,* d_{χ^2} . The histogram χ^2 -metric distance is done by

$$d_{\chi^2}(h_Q, h_I) = \sum_{i=0}^{n-1} \frac{(h_Q(i) - h_I(i))^2}{(h_Q(i) + h_I(i))}.$$
 (5)

This metric based on the χ^2 statistic has been proposed as a metric for image similarity in [3].

• *Histogram Mahalanobis distance*, d_M . The Mahalanobis distance is a special case of the quadratic-form distance metric in which the transform matrix is given by the covariance matrix obtained from a training set of histograms [6]. The Mahalanobis distance between two histograms is given by

$$d_M(h_Q, h_{I^k}) = (h_Q - h_{I^k})^T \mathbf{C}_k^{-1} (h_Q - h_{I^k}).$$
(6)

We remind that h_Q is the query histogram and h_{I^k} is the template histogram or mean histogram of a training set of histograms of class k. They both are *n*dimensional vectors and the covariance matrix C_k is a $n \times n$ matrix. In the special case when h_{I^k} (for all the classes) are statistically independent, but have unequal variances $\sigma_{I^k}^2$, C_k is a diagonal matrix. In this case, the Mahalanobis distance reduces to

$$d_M(h_Q, h_{I^k}) = \sum_{i=0}^{n-1} \frac{(h_Q(i) - h_{I^k}(i))^2}{\sigma_{I^k}^2(i)}.$$
 (7)

In our work, we have used three 1-D histograms, i.e. the color image Q is represented by three histograms $h_Q^{X_1}$, $h_Q^{X_2}$ and $h_Q^{X_3}$ such that $(X_1X_2X_3) = (RGB)$ or $(X_1X_2X_3) = (HSV)$ or $(X_1X_2X_3) = (Lab)$. Let $d(h_Q^{X_i}, h_I^{X_i})$ be the histogram distance of component X_i , the color distance in the color space $(X_1X_2X_3)$ is the sum of each component histogram distance,

$$d(h_Q, h_I) = d(h_Q^{X_1}, h_I^{X_1}) + d(h_Q^{X_2}, h_I^{X_2}) + d(h_Q^{X_3}, h_I^{X_3}).$$
(8)

4. MORPHOLOGICAL COLOR SIZE DISTRIBUTIONS

Texture has been one of the most important characteristics which have been used to classify and recognise objects and for image retrieval [2]. Granulometries are interesting for texture analysis because size distributions provide information about the *shape* and *size* of the patterns found in ordered textures as well as the degree of *granularity* of disordered textures [17]. Granulometries computed on binary or grey tone images provide signatures that can be used for pattern recognition or image retrieval. Texture classification and feature extraction by granulometries have many applications in the biomedical and industrial sector.

4.1. Basic definitions

Matheron [11] introduced the notion of granulometry, and the extension to grey tone functions was made by Serra [15]. A granulometry is the study of the size distributions of the objects of an image. Formally, a *granulometry* can be defined as a decreasing family of *openings* having a size parameter λ , $\Gamma = (\gamma_{\lambda})_{\lambda \geq 0}$, and satisfying the absorption law:

$$\forall \lambda \ge 0, \forall \mu \ge 0, \gamma_{\lambda} \gamma_{\mu} = \gamma_{\mu} \gamma_{\lambda} = \gamma_{\max(\lambda, \mu)}$$
(9)

Moreover, granulometries by *closings* (or *anti-granulometry*) can also be defined as families of closings.

Performing the granulometric analysis of an image I with Γ is equivalent to mapping each opening of size λ with a measure $m(\gamma_{\lambda}(I))$ of the opened image $\gamma_{\lambda}(I)$. This measure is typically the volume in the greyscale case. The granulometric curve, or *pattern spectrum*, of I with respect to γ , denoted $PS_{\gamma(I)}$ is defined as the following normalised mapping:

$$PS_{\gamma(I)}(\lambda) = PS_I(\lambda) = \frac{m(\gamma_\lambda(I)) - m(\gamma_{\lambda+1}(I))}{m(I)}, n \ge 0.$$
(10)

The pattern spectrum $PS_{\gamma(I)}$ maps each size λ to some measure of the bright image structures with this size. By duality,

the concept of pattern spectrum extends to anti-granulometry by closings $\varphi_{\lambda}(I)$,

$$PS_{\varphi(I)}(-\lambda) = PS_I(-\lambda) = \frac{m(\varphi_{\lambda}(I)) - m(\varphi_{\lambda-1}(I))}{m(I)}.$$
(11)

and is used to characterise the size distribution of dark image structures.

The normalised pattern spectrum is a probability density function [15], i.e. $PS_I(\lambda)$ is a size histogram of *I*. A large impulse in the pattern spectrum at a given scale indicates the presence of many image structures at that scale.

Both size distributions can be collated into a unique curve with closings versus size on the left (negative) side and openings versus size on the right (positive) side,

$$PS_I = \{ PS_I(-\lambda), PS_I(\lambda) \}.$$
(12)

If the size of the openings/closings is $1 \le \lambda \le n$, the size distributions can be represented as feature vectors in a 2n-dimensional space, called the granulometry space $\mathcal{G}^{2n} \subset \mathbb{R}^{2n}$. On this space, all the histogram distance metrics can be used in order to quantify the dissimilarity according to PS. For instance, let PS_Q be a query size distribution and PS_I be a reference size distribution, both belong to \mathcal{G}^{2n} , the size distribution Euclidean distance is defined by $d_E(PS_Q, PS_I) = \sqrt{\sum_{i=-n}^{n} (PS_Q(i) - PS_I(i))^2}$.

4.2. Color granulometries and texture distance

In this work, we are interested in the characterisation of texture in color images and therefore we have to extend the notion of granulometry, a grey tone concept, to color images. The simplest way is to calculate the pattern spectrum (by openings and closings) for each color component, i.e, the color size distribution or color granulometry is really the set of three granulometries. Moreover, we used only the representation of the color image in the *RGB* color space (the application of morphological operators on other color space components has added difficulties). The color image *Q* is represented by three size distributions PS_Q^R , PS_Q^G and PS_Q^B and the *texture distance* is the sum of each component size distribution distance,

$$d(PS_Q, PS_I) = d(PS_Q^R, PS_I^R) + d(PS_Q^G, PS_I^G) + d(PS_Q^B, PS_I^B).$$
(13)

The *RGB* components are strongly correlated, however preliminary tests on our haematological database had shown that the approach using three granulometries performs generally better than the technique using only a granulometry, typically on the luminance component. There are other alternative approaches in order to define a more genuine color size distribution. They involve the definition of morphological openings/closings in color lattices. At present, we are working on this kind of techniques.

5. INTEGRATION OF COLOR AND TEXTURE ATTRIBUTES

The use of a single image attribute for image retrieval or classification may lack in sufficient discriminatory information. In order to increase the accuracy of the retrievals, the results obtained from the query based on individual features have to be integrated. In our approach the features of color (spectral distributions in the different color spaces) and texture (morphological color granulometries in *RGB*) are defined by histograms (see figure 1), and the similarity rank or index is done by means of histogram distances. The results of the color-based and the texture-based retrieval can be integrated by combining the associated dissimilarity values. Weighted linear combination of vectors is a possible method. For instance, we define an integrated distance between Q and I as

$$d(Q, I) = \frac{\omega_c d(h_Q, h_I) + \omega_t d(PS_Q, PS_I)}{\omega_c + \omega_t},$$
 (14)

where ω_c and ω_t are the weights assigned to the color-based and texture-based similarity, respectively [8]. There are other non-linear combination methods [12]: voting, decision trees, etc., i.e,

$$d(Q, I) = f(d(h_Q, h_I), d(PS_Q, PS_I))$$
(15)

where f represent a general combination mapping. In section 6, we discuss the effects of integration methods on classification performance.

6. APPLICATION: LEUKOCYTE CLASSIFICATION BY COLOR-TEXTURE ANALYSIS

Leukocytes, or white blood cells, may be subdivided into five different categories: (1) monocytes, (2) neutrophils, (3) basophils, (4) eosinophils and (5) lymphocytes. When a peripheral blood smear, stained with May-Grünwuald Giemsa, is examined under a microscope at $\times 100$ magnification, the five classes of leukocytes may be differentiated according to their morphological and spectral features (see figure 2). The images of work were acquired at controlled light intensities, the camera is calibrated with a black/white-field image to ensure correct color registration.

In this section, we first evaluate the capability, using the different distances, of the color histograms and the color size distribution separately for image classification. Then, effects of several integration methods on classification performance are shown, a final algorithm is proposed. For the experiments described we used 123 leukocyte color images of size 256×256 pixels¹.

One of the key questions about any classifier or retrieval system is how well it will perform, i.e., what is the *error rate*

¹The leukocyte image database may be obtained from the authors upon request.



Figure 1: Example of color-texture description: (a) leukocyte color image, (b) histograms of the image the RGB space, and (c) size distributions in the RGB.

or the *precision rate*. In order to quantify the performance of our classification approaches, we define its precision as, $Precision(\%) = \frac{No. of images rightly detected}{No. of images} \times 100$ (the proportion of images correctly classified).

6.1. Experimental results

Following the SPR paradigm, the experimental study is divided into two modes: training or learning step and classification or testing step. The images have been divided into two sets: 51 images in the training data set (9 images belong to each class, except for class 3, with 6 images), 74 images in the classification data set (mixture of 5 classes).

6.1.1. Training templates

The goal of training procedure was to create a valid estimate of templates for each one of the five classes of leukocyte images. We developed it in three stages. First, computing the color histograms and the color size distributions. Second, calculating the means and the covariance matrices of each



Figure 2: *Examples of leukocyte color images: (a) lymphocytes, (b) monocytes, (c) eosinophils, (d) neutrophils and (e) basophils.*



Figure 3: Examples of templates of class 1 and 5: (a) histograms of component H and (b) size distributions of component G (* is class 1 and \triangle is class 5).

leukocyte class k (k = 1, 2, ..., 5). Third, testing on the training data (one way to estimate the preliminary precision is to see how well the technique classifies the examples used to estimate the parameters, however, it is a very poor measure of generalisation ability) [4]. Examples of mean templates are shown in figure 3.

6.1.2. Performance of distances and attributes

Then, for each leukocyte pattern image (from the classification data set) the dissimilarity to the templates was computed. The leukocyte image was assigned to the most similar template, i.e, we chose the class for which the distance is minimum. The goal of these tests was to select the best color histogram representation and the best distance metric when the different features are considered separately. Results of the precision are given in table 1.

None of the four attributes alone are enough for the classification, regardless of the chosen distance. The best distance and the corresponding accuracy were: for RGB histograms, 20.83% using d_M ; for HSV histograms, 30.55% using d_{\cap} , for Lab histograms, 34.72% using d_{\cap} and for color size distributions, 62.50% using d_E (this is the best-case classification accuracy). The worst-case classification accuracy was 2.77% using RGB histograms according to the Euclidean distance. We believe that it may be due to the fact that small variations in the luminance component of the color image

	Distance			
Feature	d_{\cap}	d_E	d_{χ^2}	d_M
	(%)	(%)	(%)	(%)
$(h^R h^G h^B)$	18.05	2.77	18.05	20.83
$(h^H h^S h^V)$	30.55	23.71	29.16	16.66
$(h^L h^a h^b)$	34.72	18.05	31.94	23.61
$(PS^R PS^G PS^B)$	19.44	62.50	61.11	25.00

Table 1: Performance of the different distances and attributes.

give rise to strong influence in the color distribution. The general low accuracy for the color distances is due to the variation of quality of the staining procedure of blood images (the color quantified is just the result of the staining). The control of the staining conditions and the control of the image acquisition are fundamental as these conditions will determine the validity and the robustness of the classification. The color granulometries presented a much higher performance in terms of accuracy (more robust with regard to the staining result).

6.1.3. Weighted linear combinations

In order to combine the color distance and the texture distance, we used the equation 14. More specifically, the combination of $d_E(PS_Q, PS_I)$ with $d_M(h_Q^{\vec{R}GB}, h_I^{RGB})$ or $d_{\cap}(h_Q^{HSV}, h_I^{HSV})$ or $d_{\cap}(h_Q^{Lab}, h_I^{Lab})$ were tested. The choice of ω_c and ω_t is determined more or less empirically. There are several possibilities, for instance, the simplest is ω_c = $\omega_t = 1$. Another possibility is to choose the weights on the basis of the accuracy of the individual feature-based classifications [8]. Since the classification on the basis of color provides between 20.83% and 34.72% accuracy in the best case and the texture-based classification provides an accuracy of 62.50% in the best case, we have chosen $\omega_c = 3$ and $\omega_t = 7$. Table 2 presents the results of classification on the basis of weighted linear integration. We observe that the results are unsatisfactory, the best performance is just 31.72%. Therefore a linear combination approach does not guarantee a better accuracy.

6.1.4. Weighted voting strategy

The weighted voting strategy operates as follows: in response to a matching operation, each feature gives five grades of similarity to the templates by increasing order distance. The ranks of the dissimilarity within each feature can then be combined by a weighted sum to output the global distance to each template according to the combination of features. The implementation of this approach is as follows. After sorting the Euclidean distance of color granulometries, 40 votes {12,10,8,6,4}, are distributed between the 5 leukocyte class

Features	Weighted Linear Combination(%)		
	$\omega_c = 1, \omega_t = 1$	$\omega_c = 3, \omega_t = 7$	
$d_M(h_Q^{RGB}, h_I^{RGB})$			
and $d_E(PS_Q, PS_I)$	22.55	25.10	
$d_{\cap}(h_Q^{HSV}, h_I^{HSV})$			
and $d_E(PS_Q, PS_I)$	26.30	30.12	
$d_{\cap}(h_Q^{Lab}, h_I^{Lab})$			
and $d_E(PS_Q, PS_I)$	24.62	31.72	
Features	Weighted Voting Strategy(%)		
	20 votes color, 40 votes texture		
The four	30.55		

Table 2: Performance of weighted linear combinations and weighted voting strategy.

distances. In the same way, according to the distance of the best three color distance histograms, 20 votes $\{6,5,4,3,2\}$, are distributed for each color space. The votes of each class are added up: the most voted is the assigned class. Again the result, 30.55%, is disappointing and this integration technique was immediately rejected.

6.1.5. Conditional supported combination

Before defining the new integrated technique, we developed a color normalisation phase in order to control any variability of staining procedure. Analysing the leukocyte image database, we observed that it is possible to normalise using as reference the background color, i.e. the color of the plasma. This color must be constant and moreover very uniform in all the blood images. The spectral uniformity involves that the corresponding mode of the histogram is very narrow and we have also confirmed that this mode is the maximum of the histogram. By means of detection of this maximum, the background normalisation can be implemented on each component. The histogram is centred according to the maximum (shifting all the values). Afterwards, the distances are computed using the new centred histograms. The notable performance of this normalisation on the HSV histograms, the accuracy after normalisation is 55.55% (30.55% without normalisation), suggested us to include it in the definitive approach.

After that some general alternatives of combination of color and texture have been tested, we introduced a specific method. The approach is based on a priority use of the color size distribution distances. The innovation is that the color histograms distances, using both spaces HSV (normalised) and Lab, supported the texture distances conditionally. The conditions of classification were the consequence of a deep study of the preliminary results. The details are as follows. The distances to the five templates are calculated. The three distances vectors are sorted in descending order $\{d_E^{PS}\}^k$,

 $\{d_{\cap}^{HSV}\}^k$ and $\{d_{\cap}^{Lab}\}^k$. Then, the following conditions (clas-

sification tree) are checked: *Ist*.- If $\{d_E^{PS}\}^1 = 2$ or $\{d_E^{PS}\}^2 = 2 \Rightarrow \text{Class } 2$. *2nd*.- If $\{d_E^{PS}\}^1 = 1$ or $\{d_E^{PS}\}^2 = 1$ or $\{d_E^{PS}\}^3 = 1$ and $\{d_{\cap}^{HSV}\}^1 = 1 \Rightarrow \text{Class } 1$. *3rd*.- If $\{d_E^{PS}\}^1 = 4$ or $\{d_E^{PS}\}^2 = 4$ or $\{d_E^{PS}\}^3 = 4$ and $\{d_{\cap}^{HSV}\}^1 = 4$ and $\{d_{\cap}^{Lab}\}^1 = 4 \Rightarrow \text{Class } 4$. *4th*.- If $\{d_E^{PS}\}^1 = 3$ and $\{d_{\cap}^{HSV}\}^1 = 3 \Rightarrow \text{Class } 3$. *5th*.- If $\{d_E^{PS}\}^1 = 5$ or $\{d_E^{PS}\}^2 = 5$ or $\{d_E^{PS}\}^3 = 5 \Rightarrow \text{Class } 5$. *6th*.- If it is still not classified $\Rightarrow \text{Class } = \{d_E^{PS}\}^1$.

The accuracy of the conditional supported combination approach is equal to 94.44%. Note that this performance is very optimal.

7. CONCLUSIONS

We have presented the use of morphological color size distributions as a low-level textural feature for classifying and retrieval images. In order to improve the performance, they can be combined with other features, typically the color distributions. The difficulty lied in the fact that all images are very similar are therefore need very precise description. We conclude based on the experiments presented that in a homogeneous database, the main factors in limiting color-texture classification performance are not only the inaccurate description of the feature distributions, but the right integration of feature similarities. There, the learning procedure in a representative image database is a must in order to define the classification strategy. An efficient content-based retrieval scheme must have the following features: accuracy and stability, but speed as well. In order to improve the speed, fast granulometry algorithms are being included [21]. We are currently investigating the definition of other morphological color size distributions, involving openings/closing in color lattices and exploring the application of the presented techniques to other homogeneous databases. This approach is on the basis of an important module of our haematological cytology image analysis and semantic indexing system.

Acknowledgements

The authors would like to thank Prof. Georges Flandrin of Hôpital Necker-Enfants Malades, Paris, France, for providing the image base and for interesting discussions about morphological and spectral characteristics of leukocytes. Thanks are also addressed to anonymous reviewers for their detailed comments and suggestions.

8. REFERENCES

[1] J. Angulo, J. Serra and G. Flandrin, "Leukocyte classification by color-texture analysis. Exploring some alternatives," *CMM-Ecole des Mines de Paris*, Internal Note N-46/01/MM, June 2001.

- [2] S. Aksoy and R. M. Haralick, "Using texture in image similarity and retrieval," in *Texture Analysis in Machine Vision, eds. M. Pietikainen and H. Bunke*, World Scientific, 2000.
- [3] J. D. Bonet, P. Viola and J. F. III, "Flexible Histograms: A Multiresolution Target Discrimination Model," *Proceedings of SPIE*, Vol. 3370, 1998.
- [4] R. O. Duda and P. E. Hart, "Pattern Classification and Scene Analysis," *Wiley-Interscience*, New York, 1973.
- [5] G. Flandrin, "Image Bank, diagnostic codification and telediagnosis in hematology," *Leukemia and Lymphoma*, Vol. 25, pp. 97–109, 1997.
- [6] M. Flicker et al., "The QBIC project: Querying images by content using color, texture and shape," *SPIE Storage and Retrieval of Image and Video Databases*, pp. 173–181, 1993.
- [7] T. Gevers and A. W. M. Smeulders, "PicToSeek: Combining Color and Shape Invariant Features for image Retrieval," *IEEE Trans. on Image Processing*, Vol. 9, No.1, 2000.
- [8] A. K. Jain and A. Vailaya, "Image Retrieval using Color and Shape," *Pattern Recognition*, Vol. 29, No. 8, 1996.
- [9] A. K. Jain, R. P. W. Duin and J. Mao, "Statistical Pattern Recognition: A review," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol. 22, No. 1, pp. 4–37, January 2000.
- [10] P. M. Kelly and T. M. Cannon, "CANDID: Comparison algorithm for navigating digital image databases," *Proc. of 7th International Working Conference on Scientific and Statistical Database Management*, pp. 252– 258, 1994.
- [11] G. Matheron, "Eléments pour une théorie des milieux poreux," *Masson*, Paris, 1967.
- [12] C. Nastar et al., "Efficient Query Refinement for Image Retrieval," *Proc. of CVPR'98*, Santa Barbara, California, 1998.
- [13] A. Pentland et al., "Photobook: Content-based manipulation of image databases," *SPIE Storage and Retrieval of Image and Video Databases II*, pp. 34–47, 1994.
- [14] G. Sharma and H. J. Trussell, "Digital Color imaging," *IEEE Trans. on Image Processing*, Vol. 6, No. 7, pp. 901–932, 1997.

- [15] J. Serra, "Image Analysis and Mathematical Morphology. Vol I and Vol II: Theoretical Advances," *Academic Press*, London, 1982 and 1988.
- [16] J. R. Smith, "Integrated Spatial and Feature Image Systems: Retrieval, Analysis and Compression," *Ph. D. Thesis*, Columbia University, USA, 1997.
- [17] P. Soille, "Morphological image analysis," *Springer-Verlag*, Berlin-Heidelberg, 1999.
- [18] M. J. Swain and D. H. Ballard, "Color indexing," *International Journal of Computer Vision*, Vol. 7, pp. 11–32, 1991.
- [19] N. Theera-Umpon and P. D. Gader, "Counting white blood cells using morphological granulometries," *Journal of Electronic Imaging*, Vol. 9, No. 2, pp. 170–177, 2000.
- [20] C. Vertan and Nozha Boujemaa, "Upgrading Color Distributions for Image Retrieval: can we do better?," *Proc. of International Conference on Visual Information Systems*, Lyon, France, 2000.
- [21] L. Vincent, "Fast grayscale granulometry algorithms," *Proc. of the EURASIP Workshop ISMM*'94, pp. 265– 272, Fontainebleau, France, 1994.