

Feature Based Retrieval of Diatoms in an Image Database Using Decision Trees

Stefan Fischer, Michael Binkert and Horst Bunke
Institute of Computer Science and Applied Mathematics
University of Bern, Switzerland
{fischer,bunke}@iam.unibe.ch

November 24, 2000

Abstract

A feature based retrieval scheme for microscopic images of diatoms in an image database is presented in this report. Diatoms are unicellular algae found in water and other places wherever there is humidity and enough light for photo synthesis. Several methods for feature extraction are described and experimental results on real diatom images are presented. The proposed feature based retrieval scheme is based on symmetry measures, geometric properties, moment invariants, Fourier descriptors and simple textural features. Based on this features the image database is divided into classes using a decision tree based classification approach. We have evaluated the discriminant power of the features and show experimental results on a diatom image database.

CR Categories and Subject Descriptors: I.2.1 [Artificial Intelligence]: Applications and Expert Systems; I.5.4 [Pattern Recognition]: Applications.

General Terms: Algorithms.

Additional Key Words: Content based image retrieval, Diatom image database, Feature selection, Decision tree induction

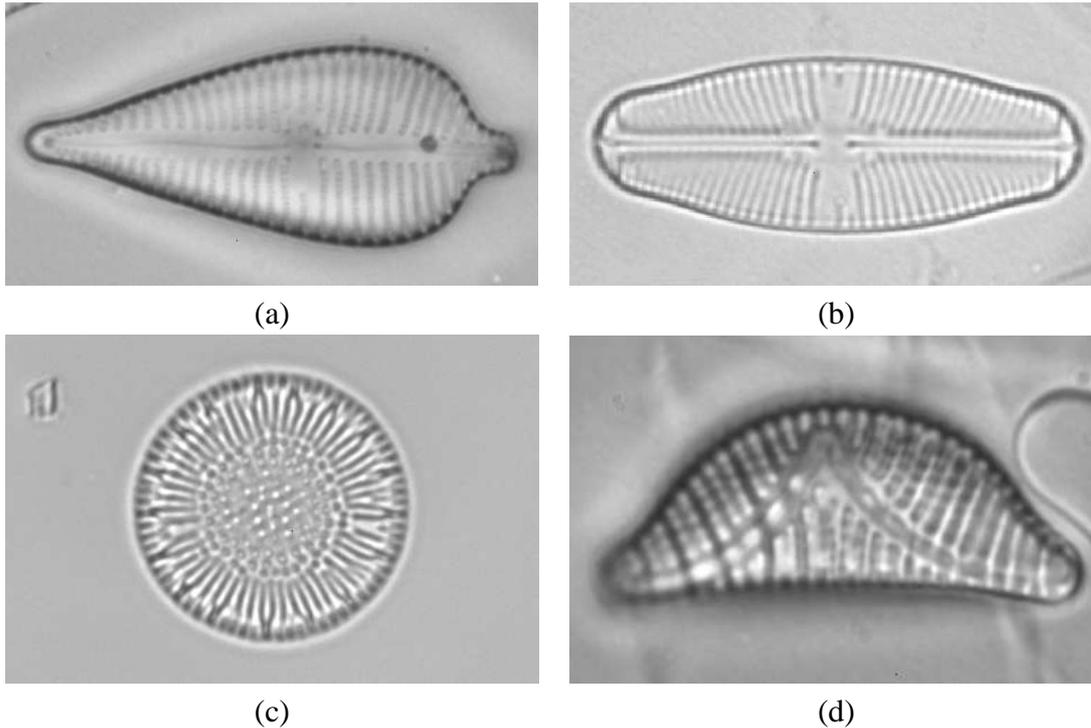


Figure 1: Example images of diatoms: (a) *Gomphonema augur*, (b) *Sellaphora pupula*, (c) *Cyclotella bodanica*, (d) *Epithemia sorex*

1 Introduction

Content based image retrieval (CBIR) is an emerging field in computer science. The aim of an image database system is to assist a human user to retrieve images. In systems which use query by image content, the query itself is an image. The system computes the similarity between the query and the stored images and returns those images which are “close” to the query. This implies that the system’s measure of image similarity corresponds to the user’s notion of similarity. The idea behind most of such systems is that similarity of image content can be characterized by combinations of low-level features such as color, texture or shape measures [8], [19].

In this report we present a feature based retrieval scheme for microscopic images in an image database. The work has been done in the framework of the ADIAC project which aims at the automatic identification and classification of diatoms [1]. Diatoms are unicellular algae found in water and other places wherever there is humidity and enough light for photo synthesis. Diatom identification and classification has a number of applications in areas such as environmental monitoring, climate research and forensic medicine [20]. One of the great challenges in diatom classification is the large number of classes involved. Experts estimate the number of diatom species between 15000 and 20000, or even higher. Example images of diatoms are shown in Figure 1. As can be seen, diatoms have various shapes and internal structures.

When biologists are going to identify these special kinds of algae, they follow a hierarchical classification procedure, where symmetry is one of the key features [3]. Subsequently other

features, such as shape properties and textural features of the internal structure are used. Because microscopic images are mostly available as gray level images or color images with minor color variations, the use of color features is not as useful in our application as, for example, in flower retrieval systems [4]. Thus we have limited the feature set to shape based features like moments, and global features of the internal structure of objects.

In the framework of automatic diatom identification, a CBIR system can be used to roughly identify objects. A possible scenario could be as follows: A user observes an unknown object during the analysis of samples. With a digital camera attached to the microscope, an image of the object is taken and transferred to the identification system. The system analyses the image and returns a list of possible diatom species together with similar sample images and textual descriptions as result. Based on this information, the user can eventually decide about the species of the unknown object.

The report is organized as follows. In Section 2 two methods to determine the symmetry of a given object are described. One method is based on the object contour, while the other uses the internal gray level distribution. In Section 3 methods for distinguishing between different boundary shapes are introduced. Global texture measures for the internal structure are described in Section 4. A classification approach based on decision tree induction is given in Section 5. Finally, experimental results obtained for decision tree induction using different combinations of features are given in Section 6, and conclusions are drawn in Section 7.

2 Symmetry

Diatoms can be grouped by their symmetry characteristics. There are species without symmetry axis while other species have one, two or even more axes. The latter occurs, for example, for circularly shaped diatoms.

Our terminology will be as follows. A straight line through the centroid of a two-dimensional figure is called a *symmetry axis* if the figure remains identical after a reflection at this straight line. A figure S is called *reflectionally symmetrical* with the degree m if it has m symmetry axes. S is called *rotationally symmetrical* with the degree $m > 1$ if there are m different angles, so that S remains identical after a rotation around $\alpha = k \cdot (360^\circ/m)$, $k = 1, 2, \dots$

Let us assume that the outline and the center of gravity of the object under consideration are known. An approach to finding the outlines of diatoms in microscopic images has been developed as part of the ADIAC project and is described in [5].

First, we describe a simple but very efficient way to determine reflectional symmetries of an object based on a 1D representation of its boundary. Next, an approach to symmetry detection based on gray level distribution is presented.

2.1 Symmetry Detection Using Distance List

The 2D shape of an object can be represented by a 1D function using the distances between the center of gravity and selected boundary points [14]. The boundary points are determined by sampling the contour in a clockwise manner with constant angle $\Delta\alpha$ as shown in Figure 2. For $\Delta\alpha = 1^\circ$ degree, for example, the distance list $\mathbf{d} = (d_0, \dots, d_{359})$ has an entry for each angle $0^\circ, \dots, 359^\circ$. For each angle the distance to the boundary is measured. For a non-convex

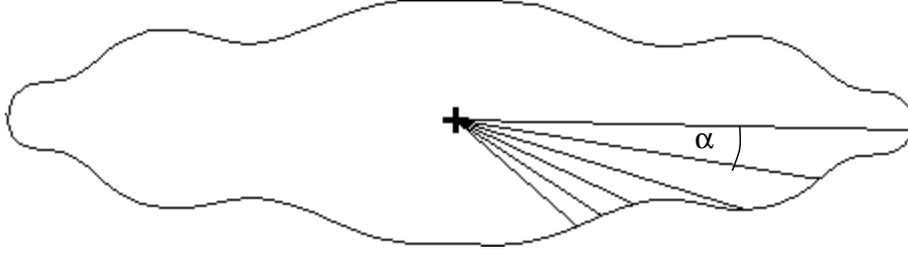


Figure 2: Contour sampling with constant angle $\Delta\alpha$

object, there can be more than one distance for the same angle. In this case the distance to the nearest boundary point is taken. To gain independence of the size of the object the entries in the distance list are normalized according to the largest distance. Using such a distance list, reflectional and rotational symmetry can be detected easily as it is shown in the following.

2.1.1 Reflectional Symmetry

The distance list allows the exact localization of all symmetry axes of a shape. Each entry i in the distance list defines a straight line with an angle $i \cdot \Delta\alpha$ through the center of the shape. If the straight line is a symmetry axis, then opposite distances must be equal. Thus the asymmetry of a shape with respect to a straight line i can be expressed as:

$$a(i) = \sum_{j=0}^{(n/2)-1} \left| d_{(i+j) \bmod n} - d_{(i-j) \bmod n} \right|$$

with n being the number of entries in the distance list.

To decide whether a straight line is a symmetry axis or not, a threshold T is used, and for each symmetry axis the condition $a(i) < T$ has to be fulfilled. Since $a(i)$ is very small for a straight line with a small offset to the real symmetry axis, a minimum angle between two different symmetry axes is defined.

2.1.2 Rotational Symmetry

For a rotational symmetry of degree 2, and an even number n of entries in the distance list, distances d_i and $d_{(i+n/2) \bmod n}$ are compared, and asymmetry is measured as follows:

$$a = \sum_{i=0}^{n-1} \left| d_i - d_{(i+n/2) \bmod n} \right|.$$

Similarly to reflectional symmetry, a shape is rotationally symmetric if the condition $a < T$ is fulfilled. If the degree of rotational symmetry m is higher than 2, a section-wise comparison is necessary. In this case the distance list is divided into m sections A_0, \dots, A_{m-1} with

$$A_i = (d_{i \cdot n/m}, \dots, d_{(i+1) \cdot n/m - 1}).$$

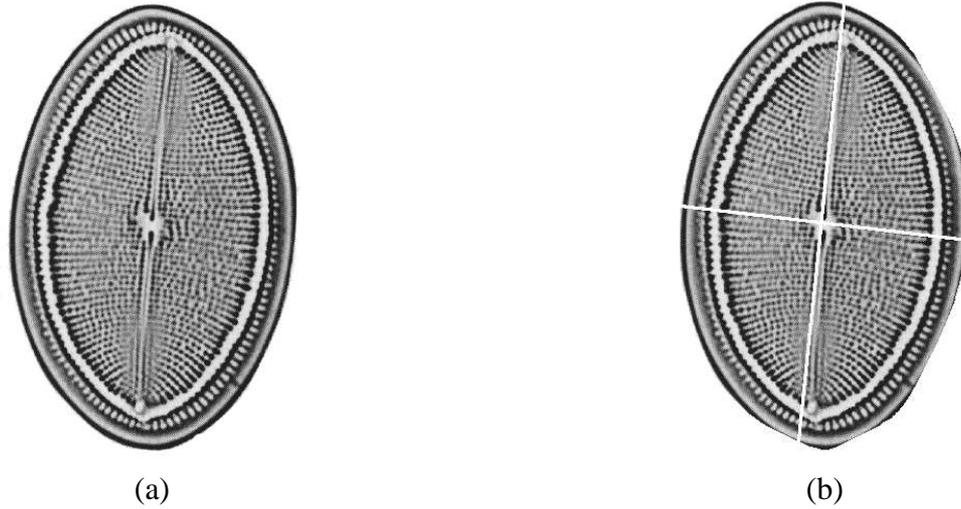


Figure 3: Example image [12] for dissimilar symmetries (a), overlaid with axes of internal symmetry (b)

Each distance d_i is contained in exactly one section, and in all other sections it has exactly one corresponding distance. Thus the asymmetry between two sections A_i and A_j is calculated by

$$a(A_i, A_j) = \sum_{k=0}^{(n/m)-1} \left| d_{i*n/m+k} - d_{j*n/m+k} \right|$$

and the asymmetry of the entire shape by

$$a = \min_{i \in [0, \dots, m-1]} \sum_{j=0}^{m-1} a(A_i, A_j).$$

The minimum is taken in order to increase the robustness under distortions.

2.2 Gray level Gradient Based Symmetry

Sometimes the symmetry of the contour of an object does not exactly coincide with the symmetry of its internal structure. This may happen due to noise or occlusion of the boundary. But there are some species of diatoms where exactly this phenomenon naturally occurs. A microscopic image of such a diatom is shown in Figure 3(a). As can be seen there is a slight offset between the location of the symmetry axes of the shape and those of the internal structure.

The direction of reflectional symmetry axes can be computed using the gradient orientation distribution of a gray level image [21]. The gradient orientation distribution is build based on the orientation of the gray level gradient vector and stored in a direction histogram. In Figure 4(a) the gradient vector direction histogram for the gray level image of Figure 3(a) is shown.

Apparently for a symmetric object the computed direction histogram is also symmetric. To determine the location of the symmetry axes the convolution c of the direction histogram h is

computed

$$c(x) = \sum_{m=0}^{n-1} h(m)h(x-m) \quad (1)$$

for $x = 0, \dots, n$. The quantities x and m denote positions in the histogram and n is the total number of entries. In the convoluted direction histogram, maxima occur at the positions of the symmetry axes. As can be seen in Figure 4(b) maxima are found at angle positions $7^\circ, 96^\circ, 187^\circ$ and 276° degree. In Figure 3(b) the gray level image is overlaid with the two symmetry axes corresponding to these maxima.

It is important to emphasize that the histogram property is a necessary but not sufficient condition for symmetry. This means that in particular non-symmetrical objects may exhibit reflectional symmetry in their orientation histogram. Thus if the orientation histogram shows reflectional symmetry, a further step is carried out to check whether the object is rotationally symmetrical or not. In our approach, this is done using the edge map. Based on the position and orientation of each potential symmetry axis, the percentage of edge elements with a counterpart one the other side of the symmetry axis is measured. If the percentage is below a certain threshold the symmetry axis is rejected. Edge elements are found by means of a standard edge detector.

3 Shape features

Besides simple shape-features such as length, width, size, and ratio, which represent geometric properties of objects, there are additional region and outline based shape measures [15]. In our approach we use moment invariants and Fourier descriptors as region and outline based features, respectively. These methods will be described in the following sections.

3.1 Geometric properties

Features that can directly be retrieved from the boundary of an object are length, width, size, and length-width ratio. Using such features in addition to symmetry information, the search space in the diatom identification process can significantly be restricted. Most diatoms have a known range of length, width, and length-width ratio. These values are listed, together with descriptions of other important features, for many diatom species in atlantes, e.g. [7], [13].

To calculate these geometric values the major axis of an object is computed. Based on this axis the length of the object is taken as the largest distance between intersections of the boundary and the axis. The width is calculated in the same way for an axis perpendicular to the major axis. If the contour of the object is polygonal approximated with vertices p_0, \dots, p_n and $p_{n+1} = p_0$ then the area enclosed by the contour can be derived from Green's theorem as:

$$area = \frac{1}{2} \left| \sum_{k=0}^{n-1} p_k \times p_{k+1} \right|. \quad (2)$$

If the resolution of the original image is known the values can be converted into μm and compared with those found in atlantes.

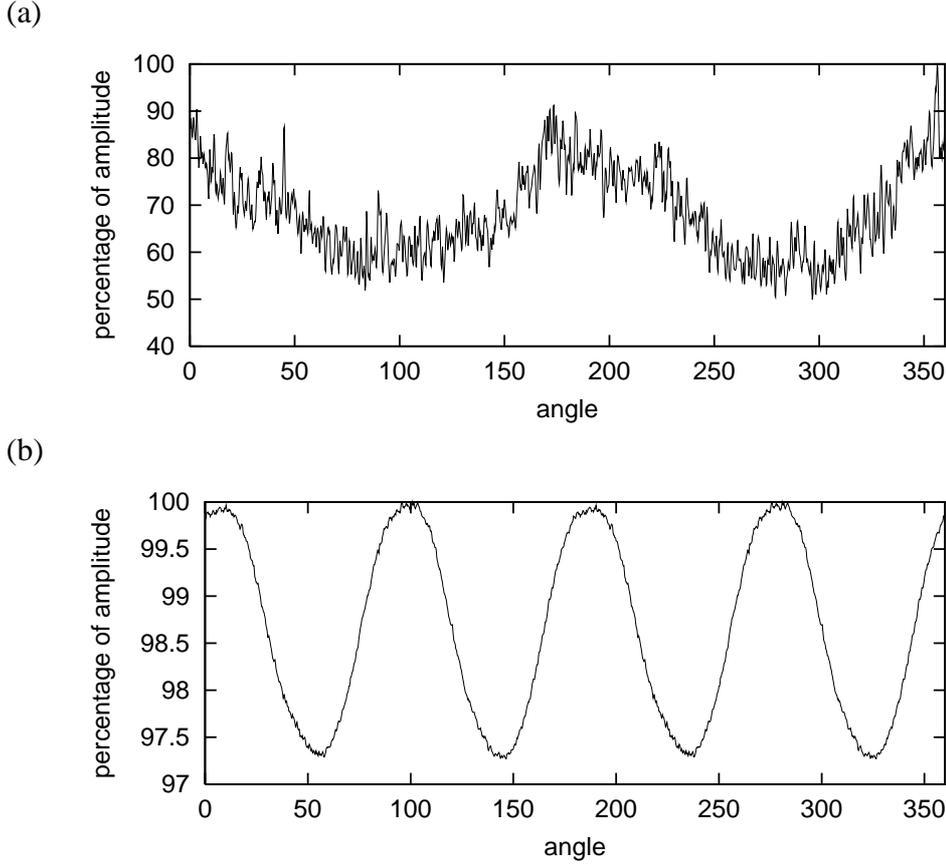


Figure 4: (a) Direction histogram of the image shown in Figure 3(a), (b) convoluted direction histogram

3.2 Moment Invariants

Moment invariants first introduced by Hu [11] are widely used in pattern recognition and have shown good results in various image recognition tasks, e.g. [6], [10]. As shape measure they have the property of being invariant under translation, scale change and rotation. For our image retrieval application we have used the seven moment invariants reported in [11] which are computed over all pixels including the boundary and its associated interior part.

Given the intensity function of an image $f(x, y)$, regular moments m_{pq} of the order $p + q$ ($p, q = 0, 1, 2, \dots$) are defined as:

$$m_{pq} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x^p y^q f(x, y) dx dy. \quad (3)$$

Based on the regular moments, central moments μ_{pq} are defined as:

$$\mu_{pq} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - \bar{x})^p (y - \bar{y})^q f(x, y) dx dy \quad (4)$$

where $\bar{x} = m_{10}/m_{00}$ and $\bar{y} = m_{01}/m_{00}$ denote the centroid of the image. The central moments are invariant under translation of the image. To get invariance under rotation the following

moment invariants are defined:

$$\begin{aligned}
I_1 &= \mu_{20} + \mu_{02}, \\
I_2 &= (\mu_{20} - \mu_{02})^2 + 4\mu_{11}^2, \\
I_3 &= (\mu_{30} - 3\mu_{12})^2 + (3\mu_{21} - \mu_{03})^2, \\
I_4 &= (\mu_{30} + \mu_{12})^2 + (\mu_{21} + \mu_{03})^2, \\
I_5 &= (\mu_{30} - 3\mu_{12})(\mu_{30} + \mu_{12}) \\
&\quad \cdot [(\mu_{30} + \mu_{12})^2 - 3(\mu_{21} + \mu_{03})^2] \\
&\quad + (3\mu_{21} - \mu_{03})(\mu_{21} + \mu_{03}) \\
&\quad \cdot [3(\mu_{30} + \mu_{12})^2 - (\mu_{21} + \mu_{03})^2], \\
I_6 &= (\mu_{20} - \mu_{02})[(\mu_{30} + \mu_{12})^2 - (\mu_{21} + \mu_{03})^2] \\
&\quad + 4\mu_{11}(\mu_{30} + \mu_{12})(\mu_{21} + \mu_{03}), \\
I_7 &= (3\mu_{21} - \mu_{03})(\mu_{30} + \mu_{12}) \\
&\quad \cdot [(\mu_{30} + \mu_{12})^2 - 3(\mu_{21} + \mu_{03})^2] \\
&\quad - (\mu_{30} - 3\mu_{12})(\mu_{21} + \mu_{03}) \\
&\quad \cdot [3(\mu_{30} + \mu_{12})^2 - (\mu_{21} + \mu_{03})^2].
\end{aligned} \tag{5}$$

These values present much of the information given by the boundary of the object in a compressed form and are therefore good candidates for an indexing scheme.

3.3 Fourier Descriptors

Fourier descriptors are used for the representation of the boundary of two-dimensional shapes. The basic idea is to represent a closed curve by a periodic function of a continuous parameter, or alternatively, by a set of Fourier coefficients of this function [2].

Starting with an arbitrary point, coordinate pairs $(x_0, y_0), (x_1, y_1), \dots, (x_n, y_n)$ are recorded while traversing the boundary in clockwise direction. The boundary can be represented as a sequence of complex numbers $s(k) = x_k + i \cdot y_k, k = 0, 1, \dots, n$. Applying the discrete Fourier transform to the sequence $s(k)$ leads to complex coefficients

$$a(u) = \frac{1}{n} \sum_{k=0}^n s(k) e^{-\frac{i2\pi uk}{n}}, u = 0, \dots, n \tag{6}$$

the so-called Fourier descriptors.

Fourier descriptors are information preserving and allow the original boundary to be restored. If only low-frequency components are used for the reconstruction, sharp features such as corners are lost, but the global shape of the object is still captured. In our approach a set of 30 Fourier descriptors is used from a total of 128 coefficients. This is enough to distinguish between the different boundary shapes. To gain independence under translation and rotation the center of gravity and the major axis of the object are computed and the boundary point on the major axis with the largest distance from the center of gravity is selected as starting point.

4 Texture

An important characteristic of diatoms is their internal texture. The identification of texture has been extensively studied in the computer vision community [17]. There are statistical methods that measure variance, entropy or energy, and perceptual techniques identifying the underlying direction, orientation and regularity. The simplest texture descriptors are based on the intensity histogram of an image. To achieve invariance under intensity changes that might arise due to the image capturing environment and its configuration, we calculate textural features on an edge direction histogram. Such a method takes the internal structure of objects into account and is largely invariant under changes of lighting conditions. The direction histogram has entries for each angle between 0° and 359° degrees and counts the number of occurrences of each specific orientation of the gradient vector computed on the gray level image (see also Section 2.2).

Texture features considered in our approach include the first four central moments corresponding to mean, variance, skewness and kurtosis [17].

The mean value μ of the normalized direction histogram h with entries h_0, \dots, h_{n-1} is given by

$$\mu = \sum_{i=0}^{n-1} i h_i \quad (7)$$

and the variance σ^2 by

$$\sigma^2 = \sum_{i=0}^{n-1} (i - \mu)^2 h_i. \quad (8)$$

The skewness μ_3 of the direction histogram is defined as

$$\mu_3 = \frac{1}{\sigma^3} \sum_{i=0}^{n-1} (i - \mu)^3 h_i \quad (9)$$

and the kurtosis as

$$\mu_4 = \frac{1}{\sigma^4} \sum_{i=0}^{n-1} (i - \mu)^4 h_i - 3 \quad (10)$$

where n is the number of angles used in the histogram.

The mean value μ gives an estimate of the average orientation of the edges within the object and the variance σ^2 is an indication of the dispersion. The histogram skewness is a measure of the histogram's symmetry. Kurtosis is a measure of the peakedness of the histogram. The main advantage of all these texture measures is their computational simplicity compared with other methods based on gray value co-occurrence matrix or spectral features.

5 Classification

There are many different classification methods from the areas of neural networks and statistical decision theory [16]. For the problem considered in this report we adopted a decision tree based approach. The reason is that decision trees resemble the way human experts identify a diatom. Moreover, because of the huge number of classes involved (see Section 1), a one level decision procedure, as performed by a neural network or a statistical classifier, seems infeasible. An

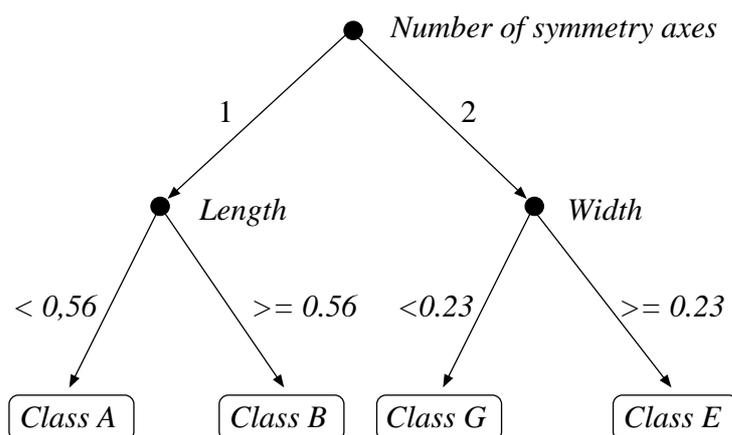


Figure 5: Example decision tree for diatom identification

additional problem with the application of a statistical classifier or a neural network is the lack of enough training data.

A decision tree based classification method is a supervised learning technique that builds a decision tree from a set of training samples. The result of the learning procedure is a tree in which each leaf carries a class name, and each interior node specifies a test on a particular feature, with one branch corresponding to each possible value or range of that feature. In Figure 5 a simple tree using the proposed features is shown. In this example in the first step the number of symmetry axes is chosen to distinguish between different kinds of diatoms. Subsequently length and width are used until classes *A*, *B*, *E* and *G* are reached.

Such a kind of a tree can easily be translated into a set of rules, or if-then-else clauses, which are human readable. Hence, the classifier inferred from a set of training samples can be interpreted by a human expert. For example the rule for Class B in Figure 5 is

```

IF Number of symmetry axes = 1 AND
   Length >= 0.56
THEN class = Class B
  
```

This property, which distinguishes decision tree based classification from neural networks and statistical approaches, is very important for taxonomic identification tasks, such as the one considered in this report. It will allow a human expert to change or extend a tree obtained from the decision tree induction procedure. Such a modification may be necessary to obtain a robust system in case of many classes, particularly if training data are sparse.

In the decision tree induction process, a tree is grown top-down in the following way. Starting with the whole training set represented by the root node, a feature-based search is done to recursively construct the subsequent layers of the tree. The basic idea is to split the training set into subsets in such a way that finally each subset holds only samples belonging to exactly one class. In our system we have employed the well-known C4.5 algorithm [18]. In this algorithm subsets are build based on an information theoretic gain criterion. At each node in the decision tree, the training set is split into subsets choosing the feature that maximizes the information gain.

Class of symmetry (0 . . . 4)
Moment invariants (1 . . . 7)
Fourier descriptors (30)
Geometric (length, width, ratio, size)
Texture (mean, variance, skewness, kurtosis)

Table 1: Proposed features for the retrieval of diatoms

For samples of the training set T the probability that the samples belongs to class C_i can be estimated as follows:

$$P(C_i) = \frac{\text{number of samples in } T \text{ belonging to } C_i}{\text{number of samples in } T}. \quad (11)$$

The entropy I of the training set T is the sum over all classes C_1, \dots, C_c

$$I(T) = - \sum_{i=1}^c P(C_i) \log_2 P(C_i). \quad (12)$$

If a feature is selected and the training set is split into subsets according to all possible values of that feature, the entropy for the split can be expressed as the weighted sum over all subsets T_1, \dots, T_n as

$$B(T, A) = \sum_{i=1}^n \frac{|T_i|}{|T|} I(T_i). \quad (13)$$

In Equation (13), A indicates the selected feature and n is the number of possible values of feature A . The entropy is small if and only if the subsets T_i hold samples from only a single or a few classes. Based on this criterion, for each feature A the information gained by splitting T according to the possible outcomes is measured by

$$G(T, A) = I(T) - B(T, A). \quad (14)$$

Consequently, when dividing the training set that feature is selected at each node of the tree which maximizes the information gain according to Equation (14). The set of samples associated with the node is divided into subsets according to the different values of feature A , and each subset represents a child node. The test is recursively continued for all child nodes until the associated subsets hold only samples belonging to a single class.

In the next section experimental results for the use of decision tree induction to categorize our test image database based on the features described in Sections 2 to 4 are shown.

6 Experimental Results

The proposed methods were tested using a database of 468 gray level images of diatoms. In the near future a significantly extended version of this database will be available on the ADIAC web page [1].

Feature	Errors	Size
Test with single groups of features		
Class of symmetry	380	6
Moment invariants	16	535
Fourier descriptors	7	391
Geometric features	26	431
Textural features	24	627
Fourier descriptors and		
Class of symmetry	0	427
Moment invariants	1	373
Geometric features	0	383
Textural features	0	389
Fourier descriptors, moment invariants and		
Class of symmetry	1	408
Geometric features	1	355
Textural features	0	373
Fourier descriptors, moment invariants, geometric features and		
Class of symmetry	0	386
Textural features	0	359
All features	0	381

Table 2: Number of errors occurring during tree induction and size of the constructed tree using the C4.5 algorithm for different combinations of features

For each of the images in the database the class¹ of the diatom is known. At the moment the database holds images of 59 different genus which is further divided into 191 different species. In most cases there are one to three images available per species. As this is not sufficient for decision tree learning and testing, we decided to distinguish between 82 classes with an average of 5 samples per class. The classes considered in our experiments are partly on genus level and partly on species level.

First we evaluated the discriminatory power of the different features. In Table 1 all features are listed. The classes of symmetry correspond to the number of symmetry axes, where number 4 stands for all cases with 4 or more symmetry axes. Some of the features listed in Table 1 are expected very important for classification, while others may be less important. To evaluate the usefulness of each feature we started to build the decision tree with just a single group of features, and subsequently combined the most powerful features. For this test the number of errors obtained during tree induction was considered as quality measure.

In Table 2 the experimental results are given. In the first column the employed features are listed while the second column gives the number of errors. The number of errors corresponds to the number of samples in the training set which were assigned to a wrong class. An error occurs

¹in terms of biologist diatoms are classified in *genus*, *species*, *subspecies* and so forth, but here we use the term *class* in the pattern recognition sense

if no combination of features can be found to distinguish between two diatoms that belong to two different classes. In column 3 of Table 2, the size of the tree which reflects the complexity of the decision procedure dependent on the chosen features is shown. All numbers reported in Table 2 are based on the full training set of 468 images. It can be observed that no single group of features is strong enough to allow an error free tree induction. The minimum number of errors resulted from the Fourier descriptors. Combining Fourier descriptors with any of the other features except for moment invariants gives an error free decision tree already. In general there are various combinations of features that lead to an error free decision tree. The size of all possible error free decision trees varies from 359 to 427.

The results shown in Table 2 indicate that it is easy to construct a decision tree that can classify the given training set without any error. However, for practical applications it is more interesting to know how well the induced tree is able to generalize the training set, i.e., how well it classifies unknown samples from a test set that were not used for decision tree construction. In order to test this ability, the error rate was measured under the *leave-one-out* method. Using the combination of all features, all but one image from the training set were used to build the tree, and the single image not involved in the decision tree induction process was used as test sample. This procedure was repeated using each sample in the training set exactly once for testing. This leads to a recognition rate of approximately 45%. The low recognition rate is due to classes with only a small number of samples included in the training set. If only samples of classes with at least n samples are considered, the recognition rate continuously increases with n .

This behaviour can be seen in Figure 6(a.). The maximum recognition rate of 69% was achieved if only classes with at least 10 samples were considered. It has to be taken into account that the number of images involved in the test decreases with the number of images per class as shown in Figure 6(b.). At present there are not more than 10 examples per class available. But in the near future the database will be significantly enlarged. Thus a further increase in recognition performance can be expected as more samples per class become available.

7 Conclusions

In this report, we first have proposed several features for use in a system for content based image retrieval. The proposed features include two different methods for symmetry detection, moment invariants, Fourier descriptors, and texture measures based on the central moments of the direction histogram. We have tested the discriminatory power of these features for the purpose of content based retrieval of microscopic images from an image database using decision trees. While no individual group of the low-level features has the strength to allow an error free tree induction, it can be concluded that the combination of features has this potential. The performance of the decision tree classifier was tested using the “leave-one-out” technique. It has been shown that the recognition rate continuously rises with the number of available samples per class. The methods described in this report are very useful for the retrieval of diatoms in the image database developed in the context of the ADIAC project. It has to be pointed out, however, that none of the methods described in this report includes any assumptions or knowledge specific to the domain of diatoms. Therefore these methods are potentially useful for other content based image retrieval tasks as well. Future work will include the identification of not only one candidate from the database, but the retrieval of the n most similar classes.

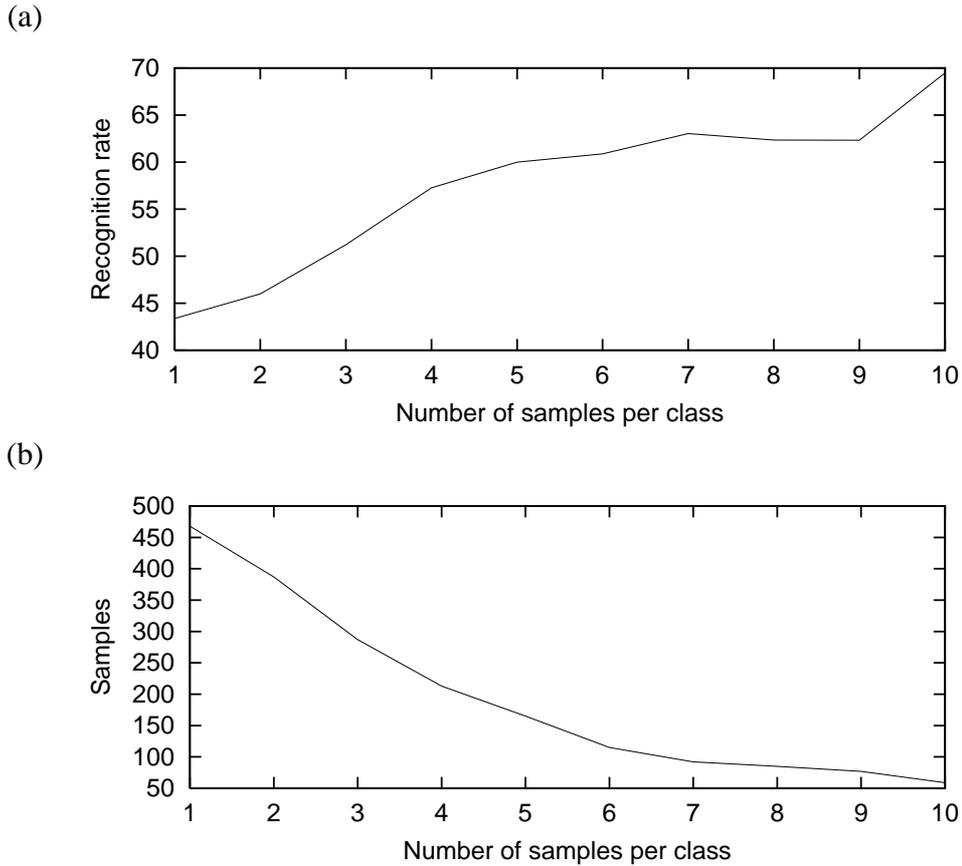


Figure 6: (a) Recognition rate (in percentage) for different numbers of samples per class. (b) Total number of samples for different numbers of samples per class.

8 Acknowledgement

The work has been done in the framework of the EU-sponsored MARine Science and Technology Programme (MAST-III), under contract no. MAS3-CT97-0122. We thank our project partners Micha Bayer and Stephen Droop for preparing the images in the ADIAC image database and for useful discussions and hints.

References

- [1] ADIAC. Automatic Diatom Identification And Classification. Project homepage: <http://www.ualg.pt/adiac/>.
- [2] K. Arbter, W. E. Snyder, H. Burkhardt, and G. Hirzinger. Application of affine-invariant fourier descriptors to the recognition of 3-d objects. *IEEE Trans. on PAMI*, 12(7):640–647, July 1990.

- [3] H. G. Barber and E. Y. Haworth. *A guide to the morphology of the DIATOM FRUSTULE*. Scientific publication No. 44,. The Freshwater Biological Association, The Ferry House, Far Sawrey, Ambleside, Cumbria, UK, 1981.
- [4] M. Das, R. Manmatha, and E. M. Riseman. Indexing flower patent images using domain knowledge. *IEEE Intelligent Systems*, pages 24–33, September 1999.
- [5] H. du Buf et al. Diatom identification: A double challenge called ADIAC. In *Proceedings of the 10th Int. Conf. on Image Analysis and Processing*, pages 734–739, Venice, Italy, Sept 27-29 1999.
- [6] S. Dudani. Aircraft identification by moment invariants. *IEEE Trans. on Computers*, 26:39–45, 1977.
- [7] F.E. Round, R.M. Crawford, and D.G. Mann. *The Diatoms: Biology & Morphology of the Genera*. Cambridge University Press, Cambridge, 1990.
- [8] M. Flickner, H. Sawhney, W. Niblack, J. Ashley, Q. Huang, B. Dom, M. Gorkani, J. Hafner, D. Lee, D. Petkovic, D. Steele, and P. Yanker. Query by image and video content: The QBIC system. *IEEE Computer Magazine*, 28(9):23–32, September 1995.
- [9] V. Ganti, J. Gehrke, and R. Ramakrishnan. Mining very large databases. *IEEE Computer Magazine*, August 1999.
- [10] L. Gupta. Invariant planar shape recognition. *Pattern Recognition*, 21:235–239, 1988.
- [11] M. Hu. Visual pattern recognition by moment invariants. *IEEE Trans. Information Theory*, 8(2):179–187, February 1962.
- [12] K. Krammer. *Kieselalgen: Biologie, Baupläne der Zellwand, Untersuchungsmethoden (in German)*. Franckh, Stuttgart, 1986.
- [13] K. Krammer and H. Lange-Bertalot. Bacillariophyceae. In H. Ettl, J. Gerloff, H. Heynig, and D. Mollenhauer, editors, *Süßwasserflora von Mitteleuropa (in German)*. Gustav Fischer Verlag, Stuttgart, 1986.
- [14] S. Loncaric. A survey of shape analysis techniques. *Pattern Recognition*, 31(8):983–1001, 1998.
- [15] B. M. Mehtre, M. S. Kankanhalli, and W. F. Lee. Shape measures for content based image retrieval: A comparison. *Information Processing & Management*, 33(3):319–337, 1997.
- [16] D. Michie, D. Spiegelhalter, and C. Taylor, editors. *Machine Learning, Neural and Statistical Classification*. Ellis Horwood, 1994.
- [17] I. Pitas. *Digital image processing algorithms*. Prentice Hall, London, 1993.
- [18] R. Quinlan. *C4.5: Programs for Machine Learning*. Morgan Kaufmann, San Mateo, USA, 1993.

- [19] J. R. Smith and S. F. Chang. Visualeek: A fully automated content-based image query system. In *ACM Multimedia Conference, Boston*, pages 87–98, November 1996.
- [20] E. F. Stoermer and J. P. Smol, editors. *The Diatoms: Applications for the Environmental and Earth Science*. Cambridge University Press, 1999.
- [21] C. Sun. Symmetry detection using gradient information. *Pattern Recognition Letters*, 16(14):987–996, 1995.