Image retrieval by examples

ITC-irst Tech. Rep. 9910-11 To appear on IEEE Trans. on Multimedia

R. Brunelli and O. Mich ITC-irst I-38050 Povo, Trento, ITALY

Abstract

A currently relevant research field in information sciences is the management of non-traditional distributed multimedia databases. Two related key issues are achieving an efficient content-based query by example retrieval and a fast response time. This paper presents the architecture of a distributed image retrieval system which provides novel solutions to these key issues. In particular, a way to quantify the effectiveness of low level visual descriptors in database query tasks is presented. The results are also used to improve the system response time, an important issue when querying very large databases. A new mechanism to adapt system query strategies to user behavior is also introduced to improve the effectiveness of relevance feedback and overall system response time. Finally, the issue of browsing multiple distributed databases is considered and a solution proposed using multidimensional scaling techniques.

Keywords: image retrieval, relevance feedback, distributed databases, multidimensional scaling, clustering, image database browsing.

1 Introduction

The current ever growing amount of multimedia data requires a big integrated effort in the research fields of Computer Vision, Information Retrieval and Database Management for its effective management. In particular, retrieving information from multimedia repositories requires the development of techniques to supplement traditional methods based on textual descriptions and searches. The reason for this necessity is twofold: associating textual descriptions to multimedia data can be very expensive, and, what is even more important, textual descriptions may not characterize data adequately for subsequent retrieval. The latter issue is of particular relevance for multimedia material, whose searching criteria and features are highly dependent on user goals.

An attempt to overcome this limitation is through query by example where non textual queries are formulated by the user using multimedia items related to the material he/she is looking for (e.g. images or video clips for searching footage). Recently many multimedia retrieval systems have investigated the query by example framework. Some of the most relevant are $OBIC^{(TM)}$, which is IBM's Query By Image Content system (see [1]), Excalibur Visual RetrievalWare^(r), a comprehensive application development software to provide content-based, highperformance retrieval for multiple types of digital visual media, Visual Information Retrieval (VIR) Image Engine by Virage, a set of libraries for analyzing and comparing the visual content of images, MARS (Multimedia Analysis and Retrieval System), an application developed by the Beckman Institute and Department of Computer Science at the University of Illinois, whose aim is to integrate various techniques in the fields of Image Processing and Information Retrieval into an Image Data Base Management System that is accessible from the web.

In the *query by example* framework, the user formulates a query by providing examples of objects similar to the one he/she wishes to retrieve. The system converts them into an internal representation used for assessing their similarity to the items stored in the database to be searched. The main advantage of *query by example* is that the user is not required to provide an explicit description of the items which is instead computed by the system. In order for this paradigm to be effective, good content descriptions must be computed automatically by the system and ways to compare them obtaining results in accordance with human judgments should also be available.

This paper discusses the use of pattern analysis techniques, such as density estimation, clustering, and multidimensional scaling, for the development of a computer assisted image search system: COMPASS.

The architecture of the system is described in Section 2. The issue of small yet effective image descriptors is considered in Section 3, while different query strategies with relevance feedback are described in Section 4. Algorithms for optimizing search strategies are presented in Section 5. Finally, some browsing issues are considered in Section 6 and the concluding remarks are reported in Section 7.

2 System architecture

The overall structure of COMPASS, an image retrieval system to support the query by example paradigm for multiple distributed databases, is presented in Figure 1. The system is configured as a client-server architecture in which a client application can submit a user query to multiple image servers. The answers from multiple image servers are then merged and proposed to the user as a single result.

Following the query by example paradigm, users rely on the images themselves to formulate queries. A generic image \mathcal{I} is characterized as a triple (I, F, M) whose elements represent a complete description of the image pixels I, possibly indirectly by pointing to the corresponding memory storage, a derived feature description $F = \{F_i\}$, automatically computed by the system, and associated meta data Mproviding information on image contents. Derived image descriptions can be computed directly by the client application while meta information is not usually provided automatically.

A query by example Q is defined by giving a set E of images and, possibly, by selecting a subset f of F and a comparison strategy S to be used by the image servers when comparing the query images to those stored in the database:

$$\boldsymbol{Q} = (\boldsymbol{E}, \boldsymbol{f}, \boldsymbol{S}) \tag{1}$$

The query images can be

- local or remote images accessible from the client application: the user may provide appropriate meta-data which can be used to supplement the visual similarity search with a more traditional textual search;
- images from a previous query considered as relevant by the user;
- images from selected image servers, relying on the browsing functionalities of the client.

In order to answer a query, the image server compares the images in the query set E to the stored ones using strategy S, obtaining a *dissimilarity* score for each of them. The dissimilarity of images could be computed using both the derived descriptors f and image meta-data M. Derived descriptors are often represented as numerical vectors while meta data are usually in textual form. The analysis presented in this paper will be limited to the use of derived de-

scriptors represented as numerical vectors, leaving out any available meta data in the computation of image similarity.

As the set of query images must be compared to other images, a function to compute the dissimilarity of an image from an image set must be introduced. If we restrict to metric spaces for the derived descriptors, the distance between an image I and an image set E can be computed using the following formula:

$$D(\mathcal{I}, \boldsymbol{E}) = \min_{\mathcal{I}' \in \boldsymbol{E}} d(\mathcal{I}, \mathcal{I}')$$
(2)

where d represent the distance defined in the metric space.

The effectiveness of feature comparison is improved by the use of *relevance feedback* which modifies the distance of the metric space using information derived from the interaction of the user with the system. The servers then sort database items by increasing dissimilarity. The set A of the top ranked ones is returned to the client together with their dissimilarity value and, if available and requested, associated meta-data. The client, upon receiving the answer from each server, sorts the resulting complete set by dissimilarity and offers to the user a single answer. The interaction of the user with the client is based on a graphical interface (see Figure 2), which, in close resemblance to the interfaces for querying traditional databases, provides:

- an area for the specification of the query: words are replaced by small image icons;
- the possibility of restricting searchable image content: document structure specifications are replaced by the different image content descriptors;
- an area where retrieved items are displayed;
- a way to limit the number of items retrieved (e.g. by imposing a threshold on the minimum required image similarity).

Besides database query by example, COMPASS also supports another very important activity: database browsing. This operation is important in the case of multimedia databases used as repositories of material to be creatively (re)assembled into new multimedia products such as cd-rom, composite images, footage, etc. Browsing support requirements differ from those of querying: database items should be organized and presented to the user in such a way that exploration of content is possible. The solution investigated in COMPASS is the organization of databases in clusters of similar images (see Section 6). Each cluster is represented by the client with a *key image* and cluster elements can be displayed on user demand providing a more detailed view of available images.

3 Image description

One of the key issues in querying image databases by similarity is the choice of appropriate image descriptors and



Figure 1. A general architecture for an image retrieval system based on the query by example paradigm. The shaded blocks are considered in detail by the current paper.



Figure 2. The COMPASS client GUI. The areas corresponding to the different functionalities of the client are outlined.

corresponding similarity measures. In a recent paper [2] the problem of quantifying the effectiveness of several low level visual descriptors was addressed. The proposed solution relies on the following definitions:

Definition 1 Given an n-dimensional histogram space \mathcal{H} and a dissimilarity measure¹ d on \mathcal{H} , the capacity curve C of \mathcal{H} is defined as the density distribution of the dissimilarity between the two elements of all possible histogram couples within \mathcal{H} .

Histogram capacity curves provide a basis on which the effectiveness, i.e. the discrimination ability, of different image descriptors can be compared. The shape of C(t) is an indicator of the distribution of histograms in \mathcal{H} with the topology induced by the selected comparison dissimilarity measure. If the average value of dissimilarity is low, histograms are not sparse enough in \mathcal{H} and histogram indexing is not effective. This can be formalized by the following definition:

Definition 2 The indexing effectiveness \mathcal{E} of an histogram space \mathcal{H} is given by the average dissimilarity value:

$$\mathcal{E} = \int y C(y) \, dy \tag{3}$$

The indexing effectiveness \mathcal{E} can be used to assess several descriptor-dissimilarity combinations for image retrieval applications. Several ways to compute image dissimilarity were considered in [2]: χ^2 , Kolmogorov-Smirnov, Kuiper, and L_p norms. The L_1 norm provided the best overall results in terms of indexing effectiveness and stability with respect to the number of histogram bins used. The main findings of [2] on the discrimination ability of some basic image descriptors are summarized in Table 1 and Figure 3. Reported data are based on two databases:

- VIDEO: a set of 40000 frames from nine different video clips. The video material was varied, ranging from comics, news, to documentaries and action movies.
- STILLS: a set of 3500 still images from a commercial collection, providing more colorful and high quality images than the *average* video material of the above database.

The effectiveness of the different descriptors can also be used to optimize the order in which they are compared. In image retrieval tasks, a threshold on the minimum acceptable similarity is usually imposed to limit the number of retrieved items. The computation of image dissimilarity can be stopped as soon as its monotonically increasing value exceeds the retrieval threshold. When multiple histograms

Descriptor	$\mathcal{E}_{ ext{VIDEO}}$	$\mathcal{E}_{ ext{STILLS}}$
Co-occurrence (hue)	57	68
Hue	55	70
Co-occurrence (lum)	52	50
Luminance	43	46
Edgeness	22	32

Table 1. The effectiveness \mathcal{E} of some lowlevel visual descriptors. Edgeness is defined as the magnitude of the luminance gradient while the co-occurrence of hue or luminance is a two dimensional histogram obtained by partitioning the image space into couples of pixels by means a a binary spatial relation: a pixel located at (x, y) is associated to a pixel at $(x + \Delta x, y + \Delta y)$ and the descriptor values at the two pixels are used as indices in a 2dimensional histogram.

are used to characterize an image, they can be concatenated in many different ways to obtain a single numerical vector describing the image. The order in which the histograms are concatenated impacts on the performance of the system. Comparing the descriptors sorted by decreasing effectiveness is expected to increase the computational savings associated to the use of a retrieval threshold. Experiments on the same data used in [2] are reported in Figure 4 and confirm this expectation.

4 Query by examples with Relevance Feedback

Relevance feedback is a fundamental mechanism by which system response can be improved by using information fed by the user [3, 4, 5, 6]. Whenever the system presents to the user a set of images considered to be similar to the provided examples, the user can pick among them the images he/she considers most *relevant* to the submitted query and add them to the original query. The resulting extended set E_R can be used to improve system response in a variety of ways [6]. A common approach to the implementation of relevance feedback for a system using image descriptors in numerical form is that of feature weighting and is based on the vector model used for textual documents.

Image derived descriptors $F = \{F_i\}$ are obtained by binning, with the same number of bins, the density estimates of the corresponding image characteristics (e.g. luminance, hue, etc.). Exploiting the homogeneity of descriptors normalization and dimensionality, the dissimilarity of

¹In this context, a dissimilarity measure is a bounded, positive, and symmetric function defined over a subset of $\mathbf{R}^{\mathbf{n}} \times \mathbf{R}^{\mathbf{n}}$.

Image descriptors capacity curves



Figure 3. The plots report the capacity curves, computed using the L_1 distance and 64 bins per histogram, for some of the image descriptors considered in the paper.

two images can be computed by:

$$d(\mathcal{I}, \mathcal{I}') = \sum_{i} \alpha_i \sum_{j} w_{ij} |F_{ij} - F'_{ij}|$$
(4)

where *i* represents the *i*-th descriptor and *j* the value of the *j*-th bin of the descriptor. This distance introduces a metric structure in the derived descriptors space and can be used to compute the distance of the query set E_R from each database item using the formula reported in Eq. 2.

The default set $\{\alpha_i\}$ is $\{1, \ldots, 1\}$ and can be modified by the user to assign different weights to the image descriptors, possibly excluding some of them from the computation of d. The set $\{w_{ij}\}$ is computed by the system and is used to incorporate relevance feedback into the comparison metric. Relevant images should be similar to each other for some of the components of their descriptors F_{ij} . This means that the standard deviations σ_{ij} computed over set E_R should be small for the components capturing the similarity of the images and larger for the components which are *not relevant*. A method to emphasize distances along the relevant directions is to use the following set of weights [7]:

$$w_{ij} = k \frac{1}{\sigma_{ij}} \tag{5}$$

In this paper a family of weighting schemes is derived from

Effect of descriptors ordering on retrieval efficiency



Figure 4. The plot reports (in double logarithmic scale) the expected gain in speed resulting from properly sequencing the image descriptors before comparing them. Note the significant advantage over the worst case, where the order in which the descriptors are used is inversely proportional to their capacity.

the previous equation

$$w_{ij} = k_\beta \frac{1}{\sigma_{ij}} \tag{6}$$

where $k_{\beta} = \left(\sum_{lm} \sigma_{lm}^{-\beta}\right)^{-1}$ is a normalizing factor, while β is a parameter which modulates the weighting effect and can be varied to optimize image comparison results. The effect of feature weighting on the computation of distances in descriptor space is presented in Figure 5. There are some major drawbacks to the use of Equation 6:

- the use of σ_{ij} tacitly assumes that the images in the query represent a compact set with ellipsoidal shape;
- the comparison metric is modified in the same way all over the descriptor space;
- the time necessary for the computation of *D* depends on the number of images in the query set *E*.

Furthermore, the amount of weighting specified by β is expected to be query dependent and should be optimized on a case-by-case basis. A way to overcome these drawbacks is presented in Section 5.

An interesting addition to relevance feedback for image retrieval comes from the introduction of *negative exam*-

Distance map: β=0



Distance map: $\beta=1$



Figure 5. The effect of feature weighting on the computation of distances. The plots report the value of the distance for each point in a two dimensional descriptor space from a set of five *query* images whose location corresponds to the minima of the plot. The introduction of feature weighting increases the distance values along the direction of lower dispersion of the query set, i.e. the horizontal axis. ples. An approach based on the generalization of boolean searches has been presented in [8] relating image distances to fuzzy predicates. In this paper a new approach is introduced using negative examples as perturbation in the metric used to compute the image distances. The user, by providing positive examples, implicitly defines a generalized boolean query whose value is given by normalized image distances. However, when the images looked for are organized in complex arrangement in the descriptor space, computing image similarity with Equations 4 and 6 may result in persistent irrelevant images, the negative examples. Knowledge of which of the retrieved images are not relevant to the current query, can be used to better characterize the regions of the descriptor space which contain relevant images by creating negative regions that can be used to carve complex geometries in feature space. The set $\boldsymbol{E}_{\hat{R}}$ of irrelevant images retrieved by the system can be used to introduce a modified dissimilarity function D'

$$D'(\mathcal{I}, \boldsymbol{E}_{R}, \boldsymbol{E}_{\hat{R}}) = D(\mathcal{I}, \boldsymbol{E}_{R}) \left[\frac{D(\mathcal{I}, \boldsymbol{E}_{R})}{D(\mathcal{I}, \boldsymbol{E}_{\hat{R}})} \right]^{\gamma}$$
(7)

where γ represents the intensity of the action of $E_{\hat{R}}$ and E_R represents the set of relevant images (i.e. the *positive examples*). As set $E_{\hat{R}}$ is close to E_R , for points lying far from $E_{\hat{R}}$ and E_R in feature space $D \approx D'$, for points nearer to E_R than to $E_{\hat{R}}$ the original dissimilarities are reduced, and for points nearer to $E_{\hat{R}}$ than E_R they are increased. A visual presentation of this effect is shown in Figure 6 for different values of γ .

5 Query Optimization

The effect of the drawbacks associated to the use of Equation 6 on the effectiveness and efficiency of relevance feedback can be minimized by

- determining whether the specified query set Q while not being compact itself, is composed by two or more compact sets: the query could then be split into simpler sub-queries, each of them better suited to the use of Equation 6. Let us note that this splitting also introduces *local* modification in the metric structure of the descriptors space;
- condensing the query set using a smaller number of images while preserving the effectiveness of the original set;
- adaptive choice of β, the parameter that modulates the amount of metric change.

A block structure of the resulting query optimization module is reported in Figure 7, while the following sections introduce the necessary pattern analysis techniques.





Distance map: $\gamma=2, \beta=1$



Figure 6. The effect of the introduction of negative examples on the computation of distances when different γ values are used. Negative examples define repulsive regions in pattern space by modifying the metric used in the computation of image distances.

5.1 Query subdivision

The cloud of points representing the query images in the descriptor space may exhibit local grouping, i.e. clusters, suggesting the splitting of the original query set into multiple subsets, each of them characterized by the images belonging to one of the clusters.

From a data analysis perspective, the relevant issue is whether the structure of the point distribution supports the presence of multiple clusters or not. There are no completely satisfactory methods to determine the number of clusters for any type of cluster analysis [9, 10]. The situation analyzed by the current paper presents additional difficulties due to the small number of images used to define the query: no asymptotic results can be used, and methods relying on density estimates can not be applied. The chosen strategy is based on two steps:

- establish whether the original query should be split or not;
- 2. if the original query should be split determine the number of clusters into which it should be split.

The first step is based on the use of a statistic originally proposed by Duda and Hart [11]. Let us denote with $d(\mathcal{I}, \mathcal{I}')$ the distance between the descriptors of two images $\mathcal{I}, \mathcal{I}'$ and with J(c) the clustering criterion function for c clusters C_1, \ldots, C_c :

$$J(c) = \sum_{i=1}^{c} \sum_{\mathcal{I} \in C_i} d(\mathcal{I}, \boldsymbol{m}_i)$$
(8)

where m_i is the *central* image of the *i*-th cluster. The quantity J(c) is a random variable whose average value decreases monotonically with c. In particular, if data are organized into \hat{c} compact, well separated clusters, the value of J(c) is expected to decrease rapidly until $\hat{c} = c$, and much more slowly thereafter. Knowledge of the distribution of J(2)/J(1) under the null hypothesis that all samples belong to a single cluster forms the basis for a test to reject or accept the null hypothesis. Unfortunately, analytical results are often not available. An approximate result is derived in [11] when the distance used in the comparison is the Euclidean norm. As the comparison metric used by COM-PASS is the L_1 norm for which results are much harder to obtain, a Monte Carlo approach was chosen [12]. As detailed in Section 3, each image is represented by histograms of several low level visual features, normalized to unit. In order to determine the distribution of $\mathcal{J} = J(2)/J(1)$ for different sample sizes n (from 6 to 16), 10000 random samples were generated that satisfied the image descriptors constraints: number of features, number of bins, and normalization to unit. For each random sample the Linde-Buzo-Gray clustering algorithm [13] using the L_1 metric was applied

10 times to find the optimal two cluster partition. The corresponding values of \mathcal{J} were then used to compute the required distributions which are summarized in Figure 8. As anticipated, the distributions for different values of n are markedly different, n being too small to ensure an asymptotic regime.

Given a set of N query images the value of \mathcal{J} is computed: if the null hypothesis of a single cluster can be rejected with the prescribed confidence, the appropriate number of clusters should then be determined. The most appropriate number of sub-queries into which the original query should be split is determined by the so called *silhouette* coefficient \mathcal{S} introduced in [14]. Let us introduce the following quantities:

$$a(i) = \frac{1}{n_{C(i)}} \sum_{j \in C(i), j \neq i} d(\mathcal{I}_i, \mathcal{I}_j)$$

$$\Delta(i, C) = \frac{1}{n_C} \sum_{j \in C} d(\mathcal{I}_i, \mathcal{I}_j)$$

$$b(i) = \min_{C \neq C(i)} \Delta(i, C)$$

where n_C is the number of elements in cluster C; the *silhouette* of element *i* is then defined as

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$
(9)

When a cluster contains a single object, s(i) = 0. The higher the value of s(i) the stronger the membership of i to its corresponding cluster. Elements that can not be clearly assigned to any cluster have a silhouette value near to zero. The silhouette coefficient S is then defined as

$$S = \frac{1}{N} \sum_{i=1}^{N} s_i \tag{10}$$

The value of S is bound to the closed interval [-1, 1]: the higher the value the better the overall classification of data for the given clustering. Furthermore, S is a dimensionless quantity that does not change when the distances between samples are multiplied by a constant factor. The knowledge of the silhouette coefficient can be used to choose an appropriate number of clusters \hat{k} so that

$$\mathcal{S}(\hat{k}) = \max_{k=2,\dots,K} \mathcal{S}(k) \tag{11}$$

The above computations are used to subdivide the original query images into several, simpler queries, each of which is better conditioned for the application of relevance feedback mechanisms (see Figure 9). The resulting simplified queries are then submitted to the image databases. For each simplified query a new comparison metric is computed according to Eq. 6. As a result, the metric used for image comparison is no longer a uniform modification of the unweighted distance: each sub-query locally modifies the comparison metric, overcoming one of the limitations of the original feature weighting approach.

5.2 Strategy optimization

Splitting the original query into smaller ones does not impact directly on the complexity of the computation of Dand does not provide any hint on the optimal value of β . However, the system, upon receiving user feedback, can automatically compare different query strategies by looking at the ranking of the user selected images in the corresponding answers: the lower the average rank, the better the strategy. Note that this is quite different from the approach introduced in [15] where the user is required to rank all the image returned by the system.

Two aspects characterize the choice of the optimal search strategy: the determination of the best query representation and the selection of the optimal β value. In the following analysis only two representations for the (positive) query set are considered:

- the original set E_R ;
- a condensed set obtained replacing the images in E_R with a *virtual* image represented by the arithmetic average of the descriptors in set.

For each representation of the query set, different values of β provide different strategies: the resulting set of dissimilarity functions constitutes the *optimization* space from which the best comparison method must be chosen.

Replacing the original query set with a single, *virtual* average image reduces the amount of computation required to estimate the distance of each database image from the query images. However, this simplification may not be always appropriate. As an example, if the original query set is not compact, the results may be meaningless, as the average image could be located in a region of feature space which is not representative of the original set. While query subdivision reduces this kind of problems, the possibility of using the condensed representation should be assessed more directly. As the purpose of using such a representation is to speed up the computation while obtaining essentially the same results associated to the original, complete set, it is necessary to verify that the two representations yield very correlated answers.

At each interaction, the system returns an image set A with the N database images most similar to the submitted query. Using this restricted number of images, it is possible to decide which representation of the query set is most efficient for the given query. This can be done by simulating system response using the restricted set A as image database. If the responses using the full and condensed representations are strongly correlated, the condensed representation is to be preferred being faster. Let us see how this

correlation can be computed.

Each image in A can be characterized with a couple of values (D_i^c, D_i^f) representing its dissimilarity to the query using the condensed and full representations respectively. The set $\{(D_i^c, D_i^f)\}_{i=1,...,N}$ can be considered as a random sample from a population with a bivariate distribution function. Let A_i be the rank of D_i^f among D_1^f, \ldots, D_N^f when they are arranged in descending order, and B_i the rank of D_i^c among D_1^c, \ldots, D_N^c defined similarly to A_i . The amount of correlation of the two answers can be quantified by their Spearman rank correlation coefficient:

$$r_s = \frac{\sum_i (A_i - \overline{A})(B_i - \overline{B})}{\sqrt{\sum_i (A_i - \overline{A})^2} \sqrt{\sum_i (B_i - \overline{B})^2}}$$
(12)

where \overline{A} and \overline{B} are the average values of $\{A_i\}$ and $\{B_i\}$, respectively. An important characteristic of rank correlation is its non-parametric nature. To assess the significance of correlation it is not necessary to know the bivariate distribution from which (D_i^c, D_i^f) are drawn. Given a confidence level, it is possible to decide whether the complete and condensed query representations results are sufficiently correlated or not.

The next step in choosing the optimal strategy is the selection of β . In order to adapt its value to the query more information is needed. The required additional data are provided by the user him(her)self with the selection of relevant images from set A. Let us restrict to a discrete set $\{\beta_j\}$ of possible values.

For each value β_j a query can be performed on A: the optimal value of β is chosen by minimizing the average rank of the newly added relevant images and of the original ones if they were selected from the queried databases. As in the complete representation the query images obtained from the queried databases always appear in the first positions (having a zero distance), this procedure is only useful to optimize the condensed query. However, the value of β can also be optimized for the complete representation in the following way: for each image in the query set obtained from the queried databases, a synthetic query is created by removing it and its rank in the resulting system answer stored. The average rank over the synthetic queries is then used for the optimization.

Data from two different queries are reported in Figure 10 and show how the optimal value for β changes for different queries. It is important to note the increase in correlation between the full and condensed queries with increasing β and the shape of the average rank curve which exhibits well defined minima. The condensed query representation is then employed using the lowest value of β for which the rank correlation of the condensed and complete representation results satisfy the required confidence level.



Figure 7. The figure reports the flow chart of the proposed query supervisor agent.





Figure 8. The plot reports the distribution of the J(2)/J(1) statistic for L_1 clustering using sample points generated taking into account the characteristics of the image descriptors used in COMPASS.











Splitting criteria

Figure 9. The figures report two different queries and the way they would be split by the system.

Figure 10. The plot reports the average rank of the query images when a condensed search is employed with different weights. The results from two different queries are reported. The rank correlation with the results obtained using all the images in the query bag is also reported.

6 Database Organization and Browsing

Browsing an image database is substantially different from querying it and presents specific interaction problems. These problems become more evident when multiple databases are browsed simultaneously. As large image databases can not be presented in their entirety to the user, useful abstractions should be developed, presenting to the user a limited number of key images which can be used to pivot the search. It would also be important to present the images in such a way that their visual presentation reflect the notion of similarity used by the system (possibly modified by the interaction with the user): nearby images (in a list or on a screen) should be visually similar and similarity relations among images should be preserved as much as possible. The task is further complicated by the dynamic nature of image similarity due to the adoption of relevance feedback techniques which change the metric structure of the descriptor space.

The solution adopted by COMPASS relies on cluster and multidimensional scaling techniques. Each image database is clustered into groups of images similar to each other according to the L_1 norm. Each group is then represented by a key image, the image closest to the cluster center. The set of cluster representatives provides the required abstraction for database browsing. Each representative acts as an hyperlink to the complete cluster population that can be shown to the user on demand.

As the number of clusters is usually much smaller than the number of images in the database, computationally expensive algorithms can be applied to organize the visual presentation of the key images to reduce browsing stress. Let us consider the situation in which two databases are used for browsing and a weighted L_1 norm is used for comparing the images. Key images should be arranged on the screen in such a way that nearby images are visually similar. The user can choose one (or more) of the image features, e.g. luminance, as a sorting key for the arrangement of the images onto the screen (see Figure 12). The key images from the two databases are then arranged onto a line preserving as far as possible their mutual similarities as quantified by L_1 distance. This can be accomplished by using multidimensional scaling techniques [11] which deal with the following problem:

Given a set of similarities or distances between every pair of N items, find a representation of the items in few dimensions such that the inter item proximities nearly match the original similarities or distances.

While it is not necessary to use the values of the similarity between the samples (the rank orders could be used instead), the following discussion is restricted to the case where the values are explicitly used. The following paragraphs follow the basic introduction of [11] to multidimensional scaling.

Let $X = \{x_i\}_{i=1,...,n}$ be a set of samples represented by points in \mathbb{R}^{K} . Let $Y = \{y_i\}_{i=1,...,n}$ be the projection of the X in \mathbb{R}^{k} , k < K and δ_{ij} , d_{ij} the distances between samples *i* and *j* in $\mathbf{R}^{\mathbf{K}}$ and $\mathbf{R}^{\mathbf{k}}$ respectively. The objective of multidimensional scaling is then to find a configuration of image points $\{\mathbf{y}_i\}_{i=1,...,n}$ for which the n(n-1)/2 distances δ_{ij} are as close as possible to the original distances d_{ij} . A measure of *closeness*, usually called *stress*, must then be introduced: the lower the stress the better the obtained scaling. Three commonly used closeness measures are:

$$J_{ee} = \frac{1}{\sum_{i < j} \delta_{ij}^2} \sum_{i < j} (d_{ij} - \delta_{ij})^2$$
(13)

$$J_{ff} = \sum_{i < j} \left(\frac{d_{ij} - \delta_{ij}}{\delta_{ij}} \right)^2 \tag{14}$$

$$J_{ef} = \frac{1}{\sum_{i < j} \delta_{ij}} \sum_{i < j} \frac{(d_{ij} - \delta_{ij})^2}{\delta_{ij}}$$
(15)

These criterion functions are invariant to rigid transformation of the points and to global (uniform) scaling. Among the above criterion functions, J_{ef} was chosen to perform the reported experiments, being a compromise between J_{ee} , which emphasizes large absolute errors, and J_{ff} , emphasizing large fractional errors. It is important to observe that it is not necessary for the distances d_{ij} and δ_{ij} to be computed in the same way, e.g using the Euclidean distance. This is true in particular in the reported context as the original distances between image descriptors are computed using the weighted L_1 norm while the distances of the projected set Y are computed using the Euclidean norm which corresponds to the user perceived distance in the space where images are to be presented. The necessity of relying on different metric structures in the original and projected spaces somehow limits the choice of the algorithms to be used in finding the desired configuration. For instance, a commonly used technique to project vectors onto a lower dimensional space is given by the Principal Component Analysis (PCA). This algorithm provides good results when the metric structure in the original and reduced space is given by the Euclidean norm and the stress is measured according to J_{ee} .

Direct minimization of the stress value leads to the so called Sammon mapping which does not depend on the type of distances used in the source and destination spaces but suffers from the following disadvantages:

- high computational requirements;
- presence of many suboptimal local minima;
- the map is given as a look-up table that must be recomputed whenever new points are added.

More recently, a fast algorithm for multidimensional scaling, FastMap, was introduced [16]. Given the set of original distances, the algorithm finds a representation of the original points in $\mathbf{R}^{\mathbf{k}}$ computing the set δ_{ij} using the Euclidean distance. When new data points are added, they can be easily projected in the reduced space. In the presented work the three cited algorithms (PCA, Sammon mapping, and FastMap) have been compared in the task of projecting points from \mathbf{R}^{16} to \mathbf{R}^1 using J_{ef} as stress indicator and different metrics in the source and destination spaces. Points in source space are given by the image luminance histograms. The results are reported in Figure 11. Sammon mapping outperforms the competing algorithms in quality and is considered to be a viable choice in spite of its shortcomings. An example of the application of multidimensional scaling to the presentation of images using as the luminance histogram as derived descriptor is reported in Figure 12.

EF stress using luminance



Figure 11. The plot reports the final stress J_{ef} for different distances and dimensionality reduction techniques.

7 Conclusions

In this paper an architecture for a general image retrieval and browsing system featuring relevance feedback was presented and discussed. In particular, the possibility of tuning search strategies and comparison metrics to varying user behavior was investigated and novel solutions presented using pattern analysis techniques. The resulting image retrieval system is able to optimize retrieval speed by reducing the number of query images while preserving retrieval effectiveness. The use of local modification of the image comparison metric, coupled to the use of negative examples further enhances the ability of the system at modeling



Figure 12. The two figures show the effect of *sorting* the cluster representatives of two databases using global luminance information. Note how the lower picture, with sorted representatives, appears more homogeneous and easy to analyze than the upper one where images are unsorted.

user needs on a per query basis. The possibilities offered by data analysis techniques have been adapted to the activity of database browsing, suggesting how clustering techniques and multidimensional scaling can be used to present a database map to the user.

References

- 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. *Computer*, pages 23–32, September 1995.
- [2] R. Brunelli and O. Mich. On The Use of Histograms for Image Retrieval. In *Proceedings of ICMCS'99*, Florence, June 1999.
- [3] I.J. Cox, M.L. Miller, S.M. Omohundro, and P.N. Yianilos. PicHunter: Bayesian Relevance Feedback for Image Retrieval. In *Proc. of Int. Conf. on Pattern Recognition*, Austria, 1996.
- [4] Y. Ishikawa and R. Subramanya C. Faloutsos. MindReader: Querying Databases through Multiple Examples. In *Proceedings of the International Conference on Very Large Data Bases*, 1998.
- [5] 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, 1997.
- [6] K. Porkaew, K. Chakrabarti, and S. Mehrotra. Query Refinement for Multimedia Similarity Retrieval in MARS. In *Proceedings of the ACM International Multimedia Conference*, Orlando, Florida, US, November 1999.
- [7] Y. Rui, T.S. Huang, and S. Mehrotra. Content-based Image Retrieval with Relevance Feedback in Mars. In *Proc. of IEEE Int. Conf. on Image Processing '97*, pages 815–818, October 1997.
- [8] Michael Ortega, Yong Rui, Kaushik Chakrabarti, Kriengkrai Porkaew, Sharad Mehrotra, and Thomas S. Huang. Supporting Ranked Boolean Similarity Queries in MARS. In *IEEE Tran on Knowledge and Data Engineering*, volume 10, December 1998.
- [9] G. W. Milligan and M. C. Cooper. An examination of procedures for determining the number of clusters in a data set. *Psychometrika*, 50:159–179, 1985.

- [10] A. K. Jain and J. V. Moreau. Bootstrap technique in cluster analysis. *Pattern Recognition*, 20(5):547–568, 1987.
- [11] R. O. Duda and P. E. Hart. Pattern Classification and Scene Analysis. Wiley, New York, 1973.
- [12] R. C. Dubes. How many clusters are best? An experiment. *Pattern Recognition*, 6(20):645–663, 1987.
- [13] Y. Linde, A. Buzo, and R.M. Gray. An Algorithm for Vector Quantizer Design. *IEEE Transactions on Communications*, COM-28(1):84–95, 1980.
- [14] L. Kaufman and P. J. Rousseeuw. Finding Groups in Data. An Introduction to Cluster Analysis. John Wiley & Sons, New York, 1990.
- [15] Y. Rui, T.S. Huang, S. Mehrotra, and M. Ortega. A Relevance Feedback Architecture in Content-Based Multimedia Information. In *Proceedings of IEEE* Workshop on Content-Based Access of Image and Video Libraries, pages 82–89, Puerto Rico, June 1997.
- [16] C. Faloutsos and K.I. Lin. FastMap: A Fast Algorithm for Indexing, Data-Mining and Visualization of Traditional and Multimedia Datasets. In *Proceedings* of SIGMOD '95, pages 163–174, 1995.