# A NEW APPROACH FOR RETRIEVAL OF NATURAL IMAGES

Mislav Grgić [*]— Sonja Grgić [*]— Mohammed Ghanbari [**]

The field of image retrieval is an active research area, particularly under the new MPEG-7 multimedia standard. According to anticipated MPEG-7 visual descriptors, image parameters that were obtained using multiscale decomposition with Gabor filters can be used for retrieval of images of similar textures. This approach has been used for texture image retrieval but we used it for the retrieval of natural images, which is a more real case for the end-user. We show that the performance of this method can be enhanced if the parameters are defined on the subimage level and then performing intersection between retrieval results.

K e y w o r d s:  image retrieval, image features, MPEG-7, Gabor filters

## 1 INTRODUCTION

The growth of multimedia information available in online-digital form is enormous in recent years. Therefore, it has become important to develop systems that will be able to search for this information, in order that some useful information can be derived. To enable search through the multimedia information, a good description of multimedia content becomes inevitable. The emerging new MPEG-7 standard [1], formally known as the "Multimedia Content Description Interface" has the objective of specifying a standard set of description schemes and descriptors to describe various types of multimedia information. Among many requirements, MPEG-7 will support visual descriptors that allow different types of queries. One of these are visual description query by shape, texture and colour. This means that for the content-based image retrieval, user will be able to specify any type of the query image content (shape, texture, colour) and the search engine will retrieve similar images from the database [2, 3]. In our research we are particularly interested in texture [4-6], and in this paper, a method for matching natural images according to their texture [7], together with a novel intersection method will be introduced and implemented for image retrieval through heads-database.

## 2 TEXTURE FEATURES

To support image retrieval or browsing, an effective representation of texture, as one of the three important parameters, is required. Numerous approaches towards texture and image features extraction have been made [8]. We used a set of Gabor filters for texture representation because it provides a perceptual characterisation of texture similar to a human characterisation. The computation of this descriptor proceeds as follows: first, the image is filtered with a bank of orientation and scale tuned filters (modelled using Gabor functions); from the filtered outputs, two dominant texture features are identified which are used for similarity image retrieval. So, our approach is well suited to the MPEG-7 requirements for image retrieval in multimedia systems.

Gabor filters can be seen as two-dimensional wavelets [9]. When applied to an image $I(x, y)$, the discretisation of a two-dimensional wavelet is given by

$$W_{mlpq} = \iint I(x,y)g_{ml}(x - p\Delta x, y - q\Delta y)\mathrm{d}x\mathrm{d}y \quad (1)$$

where $\Delta x$, $\Delta y$ is the spatial sampling rectangle (in our case $\Delta x = \Delta y = 1$), $p$, $q$ are image position, and $m$ and $l$ specify the scale and orientation of the wavelet, respectively, with $m = 0, \ldots, M-1$ and $l = 0, \ldots, L-1$. $M$ is the total number of scales and $L$ is the total number of orientations. The notation

$$g_{ml}(x,y) = a^{-m}g(\tilde{x}, \tilde{y}) \quad (2)$$

where

$$\begin{aligned}
\tilde{x} &= a^{-m}(x\cos\theta + y\sin\theta)\,, \\
\tilde{y} &= a^{-m}(-x\sin\theta + y\cos\theta)
\end{aligned} \quad (3)$$

denotes dilation of the mother wavelet $g(x, y)$ by $a^{-m}$ ($a$ is the scale parameter), and rotation by $\theta = l\Delta\theta$, where $\Delta\theta = 2\pi/L$ is the orientation sampling period.
Function $g_{ml}(x, y)$ is defined so that all the wavelet filters have the same energy

[**] Faculty of Electrical Engineering and Computing, University of Zagreb, Unska 3/XII, HR-10000 Zagreb, Croatia
[**] Department of Electronic Systems Engineering, University of Essex, Colchester, Essex, CO4 3SQ United Kingdom
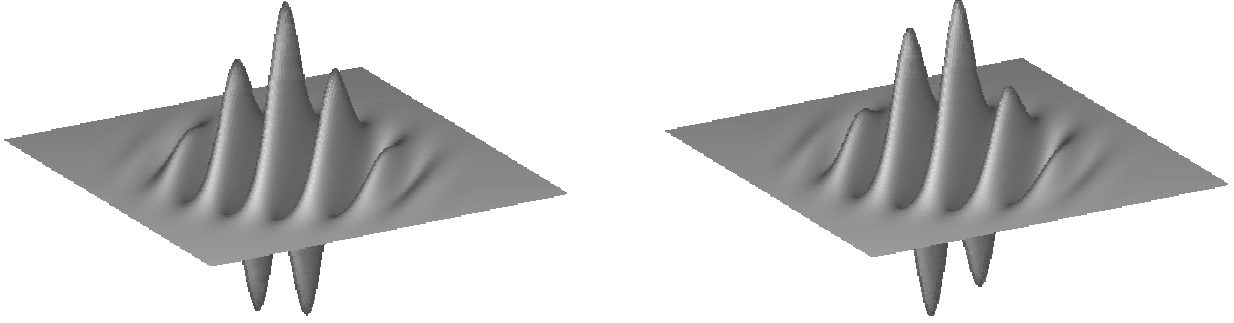
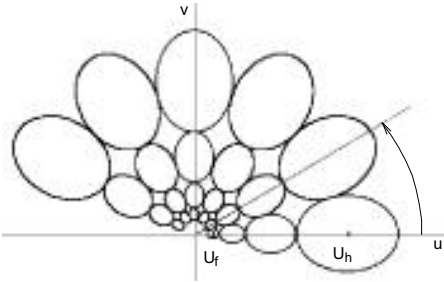**Fig. 1.** Real and imaginary parts of $g(x,y)$



**Fig. 2.** Gabor filter spectrum (contours indicate the half-peak magnitudes of filter responses)

$$\iint \left| g_{ml}(x,y)^2 \right| \mathrm{d}x\mathrm{d}y = \iint \left| a^{-m}g(\tilde{x},\tilde{y}) \right|^2 \mathrm{d}x\mathrm{d}y$$

$$= \iint \left| a^{-m}g(\tilde{x},\tilde{y}) \right|^2 \left| J^{-1} \right| \mathrm{d}\tilde{x}\mathrm{d}\tilde{y}$$

$$\left( J = \begin{bmatrix} a^{-m}\cos\theta & a^{-m}\sin\theta \\ -a^{-m}\sin\theta & a^{-m}\cos\theta \end{bmatrix} \right)$$

$$= \iint \left| a^{-m}g(\tilde{x},\tilde{y}) \right|^2 a^{2m}\mathrm{d}\tilde{x}\mathrm{d}\tilde{y} = \iint \left| g(\tilde{x},\tilde{y}) \right|^2 \mathrm{d}\tilde{x}\mathrm{d}\tilde{y} \quad (4)$$

which is obviously independent of the choice of $m$ and $l$.

We use the family of our filter bank with the following Gabor function as the mother wavelet

$$g(x,y) = \frac{1}{2\pi\sigma_x\sigma_y} \cdot e^{-\frac{1}{2}\left(\frac{x^2}{\sigma_x^2}+\frac{y^2}{\sigma_y^2}\right)} \cdot e^{i2\pi f_0 x} \quad (5)$$

where $\sigma_x$ and $\sigma_y$ are the spreads of the Gaussian and $f_0$ is the spatial frequency of the harmonic wave (frequency bandwidth of the filters). We use $W = 0.5$ which corresponds to a half-amplitude bandwidth of 1 octave and is consistent with neurophysiological findings [9]. An example of real and imaginary parts of $g(x,y)$ is shown in Fig. 1.

The Gabor function is a two-dimensional Gaussian modulated with a complex exponential. Therefore, its frequency domain representation is a two-dimensional Gaussian with an appropriate displacement along $u$-axis, where we use $(u,v)$ to index the frequency domain

$$G(u,v) = e^{-2\pi^2(\sigma_x^2 u^2 + \sigma_y^2 v^2)} ** \delta(u - f_0)$$

$$= e^{-2\pi^2(\sigma_x^2(u-f_0)^2 + \sigma_y^2 v^2)} = e^{-\frac{1}{2}\left(\frac{(u-f_0)^2}{\sigma_u^2}+\frac{v^2}{\sigma_v^2}\right)} \quad (6)$$

where $\sigma_u = 1/(2\pi\sigma_x)$ and $\sigma_v = 1/(2\pi\sigma_y)$. Symbol $**$ represents two-dimensional convolution.

Hence, using $g(x,y)$ as the mother Gabor wavelet, by appropriate dilations and rotations of $g(x,y)$ [7], it is possible to design Gabor filter dictionary, Fig. 2. $U_l$ and $U_h$ are the lowest and highest frequency of interest, respectively. We used $U_l = 0.05$ and $U_h = 0.4$.

For each Gabor filter, the input image can be decomposed into $M \times L$ filtered images. We have chosen $M = 4$, $L = 6$, resulting in 24 different filters which decompose the input image into 24 filtered images. These parameters are selected as a balance between the computational complexity and texture representation accuracy. For every filtered image the mean value

$$\mu_{ml} = \iint |W_{ml}(x,y)|\mathrm{d}x\mathrm{d}y \quad (7)$$

and the standard deviation

$$\sigma_{ml} = \sqrt{\iint \left(|W_{ml}(x,y)| - \mu_{ml}\right)^2 \mathrm{d}x\mathrm{d}y} \quad (8)$$

as the two most important features of the texture are extracted. Hence, the image texture is described by a feature vector with 48 feature parameters as

$$\text{Feature factor} = \big[\mu_{00}\mu_{01}\mu_{02}\mu_{03}\mu_{04}\mu_{05}\mu_{10}\mu_{11}\ldots\mu_{35}$$

$$\sigma_{00}\sigma_{01}\sigma_{02}\sigma_{03}\sigma_{04}\sigma_{05}\sigma_{10}\sigma_{11}\ldots\sigma_{35}\big]_{1\times48} \quad (9)$$

## 3 DATABASE CREATION AND SIMILARITY MEASURE

We created a database which contains 300 head images of $512 \times 512$ pixels [10], Fig. 3. It is in fact 30 sets, each containing 10 images of the same person with different facial expressions, illumination conditions and occlusions (glasses, hairstyle, *etc*). For each of the 300 images, 48 texture features are extracted, and a database with $300 \times 48$ features is created. Following the conventional retrieval approaches, each image in the database is identified with its 48 texture features. To search for a desired image (query image) in the database, the 48 texture features of this query image are also extracted. We then calculate the distance between the set of 48 features of the
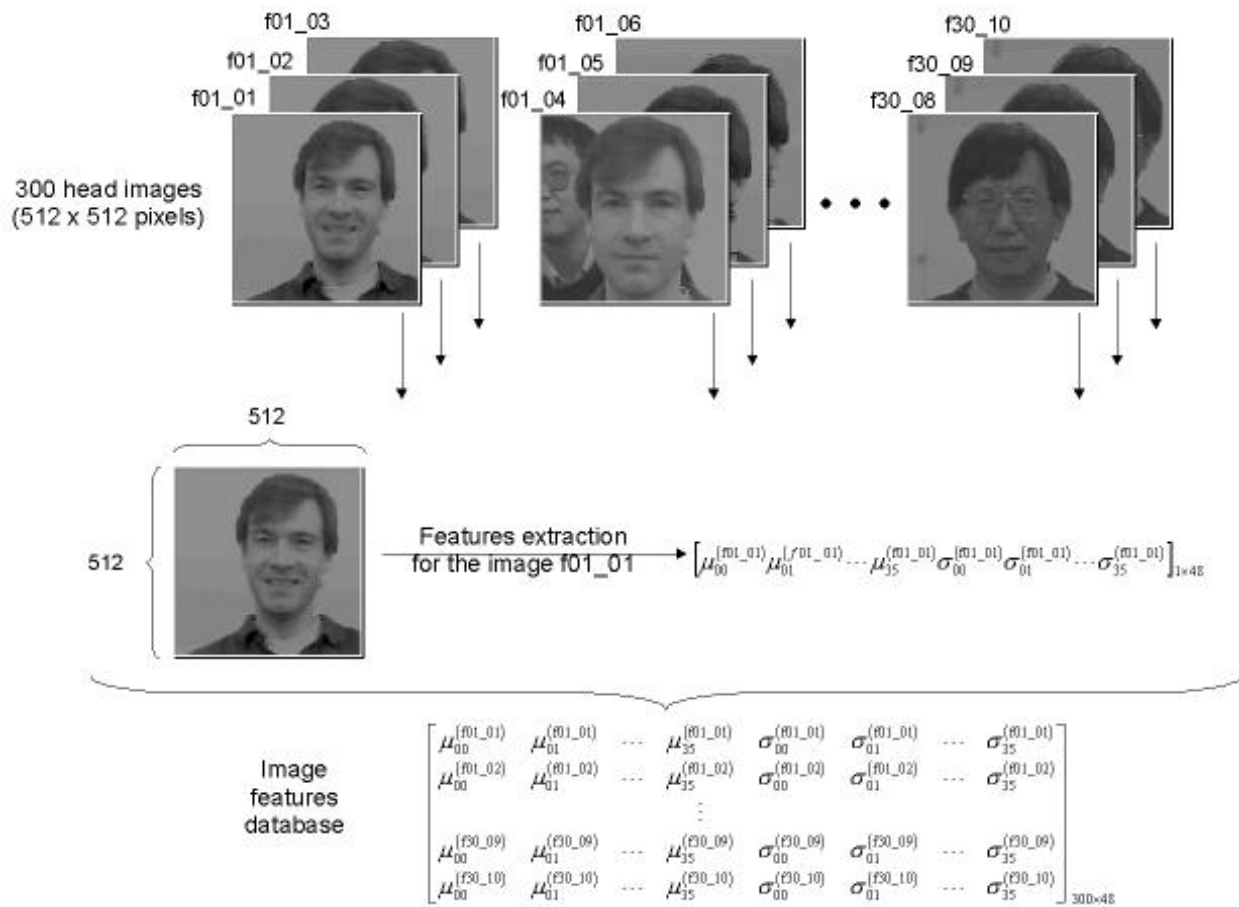
**Fig. 3.** Creating the image features database

query image and each set of 48 features of the database images. The distance is determined by measuring the normalised Euclidean distance between two feature vectors as

$$\mathrm{d}_{ml}(i,j) = \left| \frac{\mu_{ml}^{(i)} - \mu_{ml}^{(j)}}{\sigma(\mu_{ml})} \right| + \left| \frac{\sigma_{ml}^{(i)} - \sigma_{ml}^{(j)}}{\sigma(\sigma_{ml})} \right| \qquad (10)$$

where $\mathrm{d}_{ml}(i,j)$ represents the distance between features for image $i$ and image $j$ on a specific scale $m$ and orientation $l$. The total distance between feature vector for image $i$ and image $j$ can be derived as

$$\mathrm{d}(i,j) = \sum_m \sum_l \mathrm{d}_{ml}(i,j). \qquad (11)$$

Similarity measure of two images can be defined as an inverse value of the total distance. If the total distance between feature vectors increases, similarity measure of images will decrease. The described procedure produces a set of 300 different similarity measures, which are sorted in an increasing order of similarity. Since every similarity measure is associated with an appropriate database image number, then the most similar image(s) can be retrieved. It should be noted that the retrieved image can not be

exactly identical to the wanted image, unless it is itself in the database. For this reason one might ask to search for most similar images, preferably to be rank ordered on the degree of similarity. Hence a good image retrieval method is the one that can find the $N$ most similar images.

## 4 RETRIEVAL RESULTS

To determine the precision of retrieved images for particular $N$ top matches (number of results per query image), we have measured the ratio between the number of retrieved relevant images (that are similar to the query image) and the number of all retrieved images (output images).

Naturally, the most similar retrieved image will be the image itself, the maximum number of retrieved relevant images is 10 (all associated with the same person), and the maximum number of retrieved images is 300. In the precision determination, the first retrieved image (image itself) will not be counted because in a real case, if the query image is not included in the database, surely, it will not be retrieved. So, precision determination will start from $N = 2$. Let us call the above described retrieval
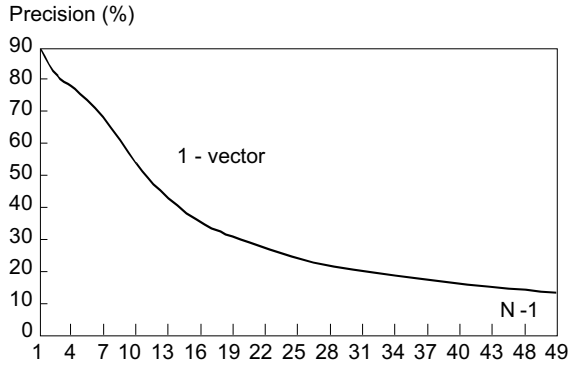
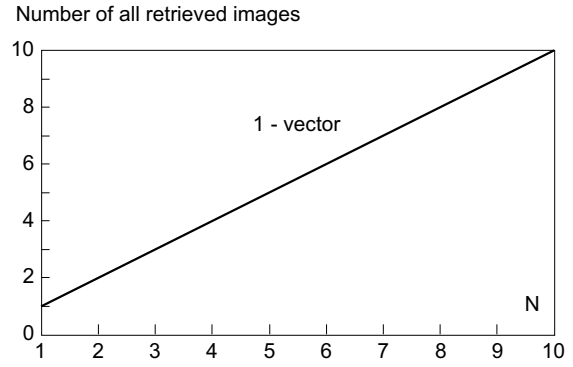**Fig. 4.** Precision of retrieved images of 1-vector approach



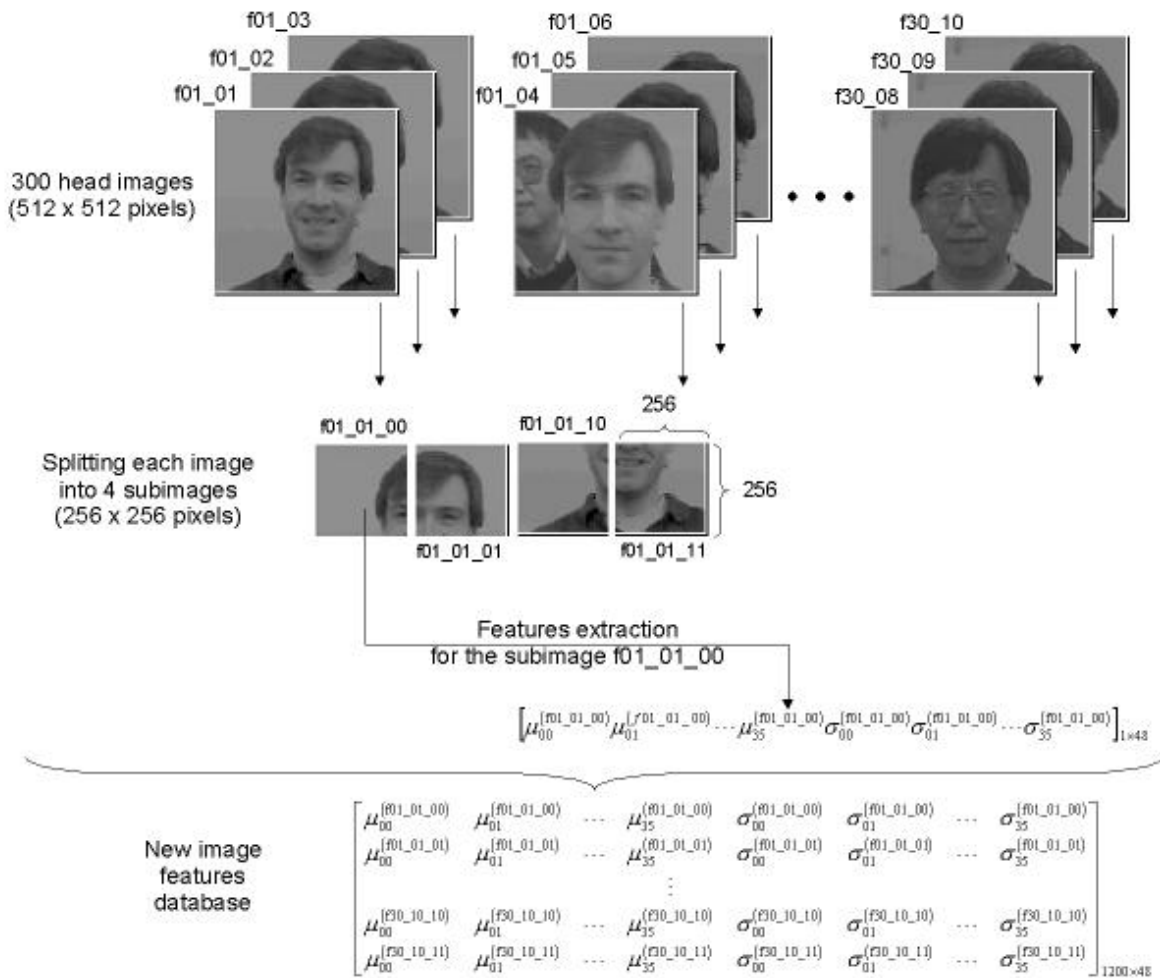**Fig. 5.** Number of all retrieved images



**Fig. 6.** Creating the subimage features database

method 1-vector approach (conventional approach), because only 1-vector with 48 parameters for the query image is extracted and compared with the database vectors.

We made 300 tests using each database image as the query image and then these 300 results are averaged for particular $N$. Figure 4 shows the precision of retrieved images when the number of top matches varies from 2 to 50. It can be seen that for larger $N$, precision decays rapidly. However, an important issue in most retrieval

processes is not only the precision of retrieved images but also the number of all retrieved images that the end-user needs to deal with. In practice, this number needs to be such that the retrieved images can be displayed, view and scrolled conveniently. With conventional 1-vector approach, it is obvious that the number of all retrieved images is fixed and determined by $N$, what can be seen from the linear curve in Fig. 5. Here should be noted that retrieved images need to be relevant images. The preci-
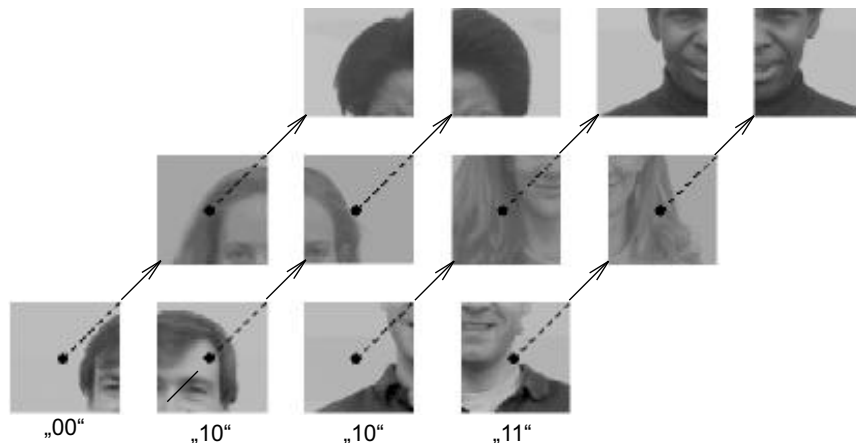
„00"      „10"      „10"      „11"

**Fig. 7.** Comparing the subimages with the same coordinates
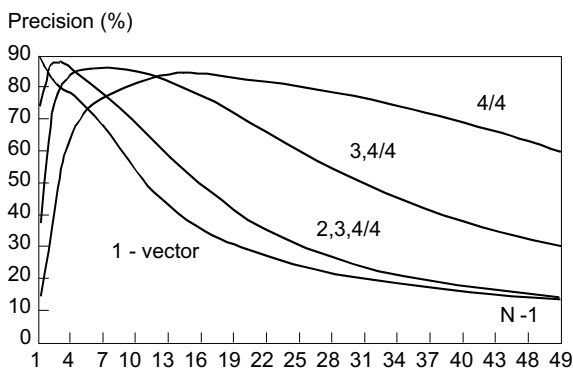


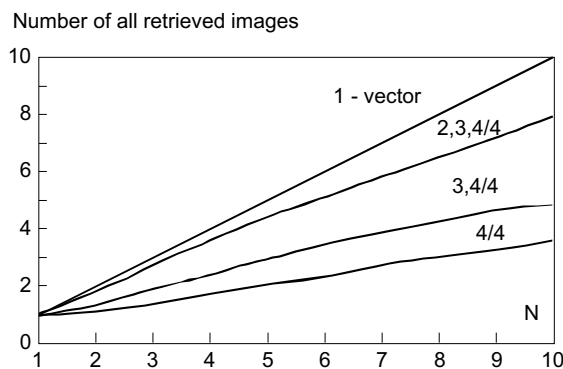**Fig. 8.** Precision of retrieved images

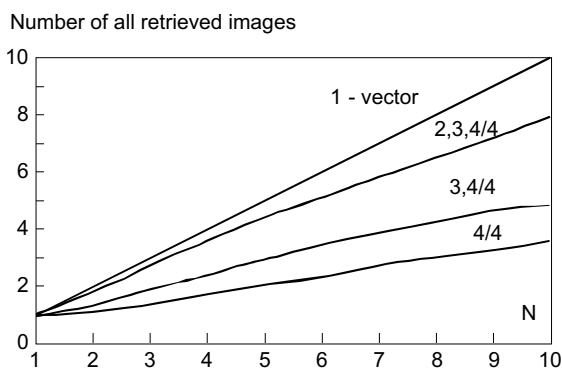

**Fig. 9.** Number of all retrieved images



**Fig. 10.** Number of cases where only the image itself is retrieved

sion for $N = 10$ is $58\%$. It means that $58\%$ of retrieved images will be relevant images. To realise the reduction of the number of all retrieved images, and at the same time, to preserve or even improve the precision, we are performing intersection between results using $256 \times 256$ subimages.

## 5 A NEW APPROACH

We are dividing each $512 \times 512$ image into 4 non-overlapped subimages (format $256 \times 256$ pixels), producing a database of 1200 subimages, Fig. 6. The same features extraction process is performed but now on each of the subimages.
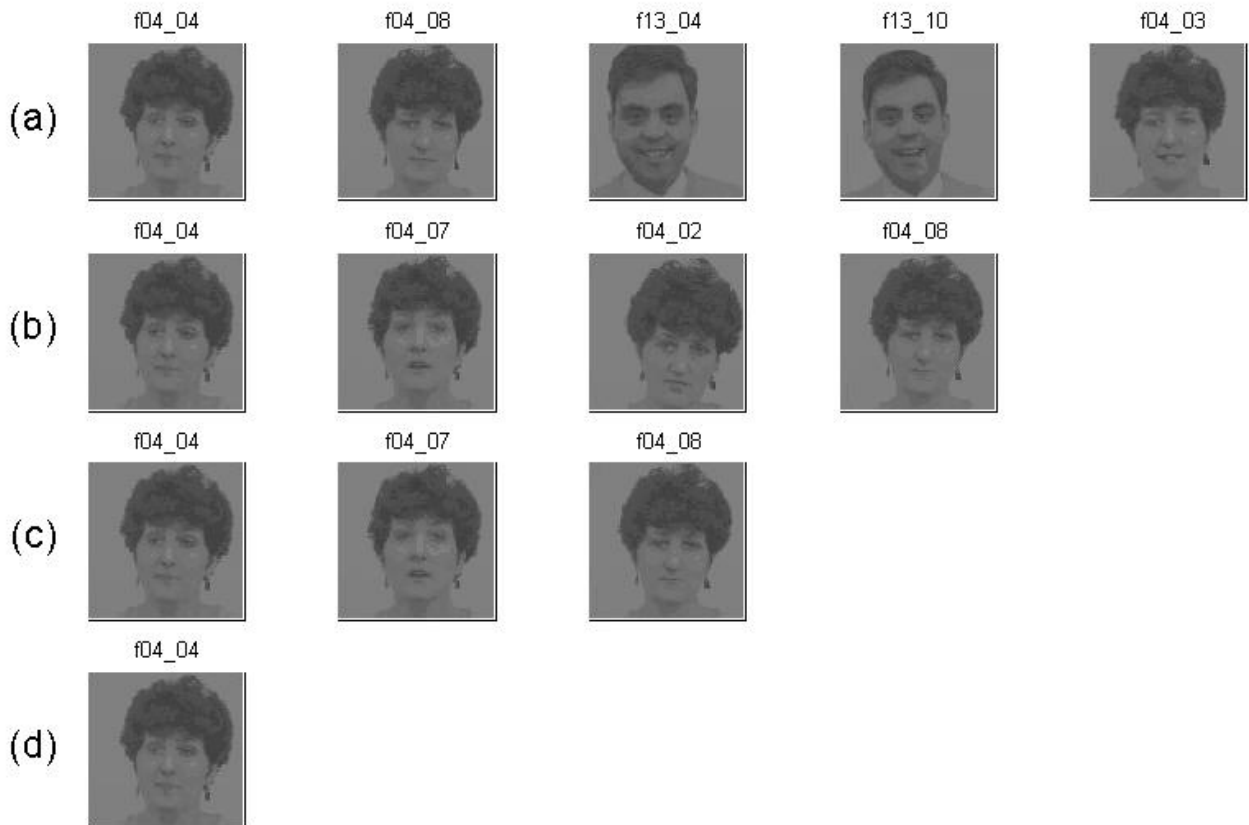
To preserve the equality with conventional 1-vector approach, we are searching for $N$ top matches within 300 subimages (same database length!), according to the following rule: each $256 \times 256$ subimage that belongs to the top-left, top-right, bottom-left and bottom-right quadrant from the $512 \times 512$ query image is compared with the 300 database subimage that belongs to the same top-left, top-right, bottom-left and bottom-right quadrant, respectively, Fig. 7.

After that, within the retrieved $4 \times N$ images, we are performing intersection. Intersection analysis is made according to three different rules:

1. image will be considered as the intersection result if it appears 4 times in each of the 4 retrieved sets (4/4-method);

2. image will be considered as the intersection result if it appears 3 or 4 times in each of the 4 retrieved sets (3,4/4-method);

**Fig. 11.** Example of retrieval process with $N = 5$: (a) retrieval results using 1-vector approach, (b) retrieval results using intersection: (2,3,4/4)-method; (c) retrieval results using intersection: (3,4/4)-method; (d) retrieval results using intersection: (4/4)-method. First image from the left in each row represents the query image.



**Fig. 12.** Example of retrieval process with the same notation like on the Fig. 11.

3. image will be considered as the intersection result if it appears 2, 3 or 4 times in each of the 4 retrieved sets (2,3,4/4-method).

Repeating the same procedure for each database subimage as the query one (300 tests), we can draw the similar curves like for the 1-vector approach, Fig. 8 and Fig. 9. It can be seen that the retrieval process with the intersection method is significantly improved. Precision of retrieved images is much better than for the 1-vector approach, and at the same time, number of all retrieved images is reduced.

Comparing four precision curves, it can be seen that all three intersection cases outperform the 1-vector approach. The highest precision will be achieved by going along the (2,3,4/4)-curve from $N = 2$ until $N = 5$, then, by going along the (3,4/4)-curve from $N = 6$ until $N = 13$, and finally, by going along the (4/4)-curve, from $N = 14$ onwards. These results are logical because for the higher $N$ it is necessary to use a stronger condition. So, by performing image retrieval using the described intersection method on our head-database, it is advisable simply to change the intersection rule parallel with increasing the $N$ in order to achieve better precision.

Comparing four curves from Fig. 9 it can be seen that all three intersection cases are better than the 1-vector approach, because the number of all retrieved images on the output is reduced and, at the same time, precision is increased. However, from these intersection curves new conclusions about three intersection cases can be derived. (2,3,4/4)-method is too close to the 1-vector approach. It means that for the same $N$, number of all retrieved images will be very similar to the 1-vector approach, and for the higher $N$, precision is falling rapidly in comparison to other two intersection methods, Fig. 8.

On the other hand, the (4/4)-method is reducing the number of all retrieved images too much. It means that there will be many cases (out of 300 tests) in which only the image itself will be retrieved on the output, Fig. 10. For $N = 1$, all methods have the same number of such cases because the image itself is always retrieved if it is included in the database. Therefore, for $N = 1$ there are 300 described cases for all methods. For $N = 2$, the 1-vector approach will always retrieve 2 images, and the number of cases where only the image itself is retrieved is equal to zero. For $N = 2$, the (3,4/4)-method will have 192, and the (4/4)-method 256 described cases. For $N = 10$, there are 26 and 60, and for $N = 100$, 0 and 3 cases, for (3,4/4)-method and (4/4)-method, respectively. The case where only one image is retrieved on the output is the worst because in real applications, where query image is not inside the database, number of retrieved images will be zero. Based on these conclusions, it can be seen that the (4/4)-method is not acceptable for real applications, because even with a high number of top matches per subimage, there are still some cases where only the image itself is retrieved. However, for higher $N$ this strong intersection rule proves to be good if one wants to preserve good precision. So, to achieve a compromise between good precision and the optimal number of all retrieved images, the (3,4/4)-method should be used. Precision with this method is more than 80 % until $N = 13$, what is acceptable number of top matches per query subimage, and the average number of all retrieved images for $N = 13$ is 6, what is very convenient for displaying, viewing and scrolling.

To justify the mentioned conclusions, an example of retrieval process is shown in Fig. 11. Results of the above mentioned methods can be seen in this figure. The first image from the left in each row represents the query image (image itself; similarity measure is equal to zero). The 1-vector approach retrieves many images that are completely different from the query image, Fig. 11(a). The (2,3,4/4)-method is good but retrieves one non-relevant image, Fig. 11(b). The (4/4)-method is too strong, especially for the smaller $N$, like $N = 5$ in this example, where only one result (image itself) is retrieved, Fig. 11(d). Finally, the (3,4/4)-method proves to be the best because the precision for this example is the highest and the number of all retrieved images is optimal, Fig. 11(c). Similar conclusions can be derived from the examples in Fig. 12. If we increase the number of top matches ($N$), 1-vector approach will retrieve too many images with a lower precision (there will be many retrieved results different then the query image). The intersection methods will preserve the precision with the lower number of retrieved images.

## 6 CONCLUSION

We have presented an approach for significant improvement of image retrieval. Improvement is based on performing the intersection between retrieval results for 4 query subimages. We introduced three different rules for intersection analysis. We compared precision results and the number of all retrieved images for conventional method and three novel intersection methods. We proved that all intersection methods are better than the conventional 1-vector approach. We found that the (3,4/4)-method should be used to achieve a compromise between good precision and the optimal number of all retrieved images. Our approach has several advantages: 1. It does not effect the database length (300), and hence, can be correctly compared with the conventional 1-vector approach, 2. It outperforms the conventional 1-vector approach according to the precision of retrieved images, 3. It reduces the number of all retrieved images.

REFERENCES

[1] ISO/IEC JTC1/SC29/WG11, N3752 (ed. Martínez, J.M.), Overview of the MPEG-7 Standard (Version 4.0), La Baule, France, October 2000.

[2] SMITH, J. R.: Integrated Spatial and Feature Image Systems: Retrieval, Analysis and Compression, PhD Thesis, Columbia University, New York, 1997.

[3] RUBNER, Y.: Perceptual Metrics for Image Database Navigation, PhD Thesis, Stanford University, Stanford, 1999.

[4] BRODATZ, P.: Textures — A Photographic Album for Artists and Designers, Dover, New York, 1966.

[5] GRGIĆ, M.—GHANBARI, M.: Similarity Texture Retrieval, Proceedings of the 2nd International Workshop on Video Processing and Multimedia Communications, VIPromCom, Zadar, Croatia, 2000, 103–108.

[6] TUCERYAN, M.—JAIN, A. K.: Texture Analysis, Chapter 2.1, in Chen, C.H., Pau, L.F., Wang, P.S.P.(Eds.): Handbook of Pattern Recognition & Computer Vision,, World Scientific, Singapore, 1993.

[7] MANJUNATH, B. S.—MA, W. Y.: Texture Features for Browsing and Retrieval of Image Data, IEEE Transactions on Pattern Analysis and Machine Intelligence **18** No. 8 (1996), 837–842.

[8] SMEULDERS, A. W. M.—GEVERS, T.—KERSTEN, M. L.: Computer Vision and Image Search Engines, Proceedings of the Workshop on Image and Video Content Based Retrieval, Milano, Italy, 1998.

[9] LEE, T. S.: Image Representation Using 2D Gabor Wavelets, IEEE Transactions on Pattern Analysis and Machine Intelligence **18** No. 10 (1996), 959–971.

[10] PEIPA (The Pilot European Image Processing Archive): http://peipa.essex.ac.uk/ipa/pix/faces/manchester/.

**Mislav Grgić** received the BSc, MSc and PhD degrees in electrical engineering from the University of Zagreb, Faculty of Electrical Engineering and Computing, Zagreb, Croatia, in 1997, 1998 and 2000, respectively. He is currently a Research Assistant at the Department of Radiocommunications and Microwave Engineering, Faculty of Electrical Engineering and Computing, University of Zagreb, Croatia. His research interests include image and video compression, wavelets image coding, texture-based image retrieval and digital video communications. Dr Grgić has been a member of the organising committees of several international workshops and conferences, as well as a Program Co-Chair of the 3rd International Symposium on Video Processing and Multimedia Communications — VIPromCom-2001. From October 1999 till February 2000 he was on a research study at the Department of Electronic Systems Engineering, University of Essex, Colchester, United Kingdom, working with Professor Mohammed Ghanbari. Dr Grgić is the recipient of four chancellor awards for best student work, he received a bronze medal "Josip Lončar" from the Faculty of Electrical Engineering and Computing in Zagreb for an outstanding BSc thesis work, and a silver medal "Josip Lončar" for outstanding MSc thesis work. He is a member of IEEE, IEEE Signal Processing Society and IEEE Communications Society.

**Sonja Grgić** received the BSc, MSc and PhD degrees in electrical engineering from the University of Zagreb, Faculty of Electrical Engineering and Computing, Zagreb, Croatia, in 1989, 1992 and 1996, respectively. She is currently an Assistant Professor at the Department of Radiocommunications and Microwave Engineering, Faculty of Electrical Engineering and Computing, University of Zagreb, Croatia. Her research interests include television signal transmission and distribution, picture quality assessment, wavelet image compression, and broadband network architecture for digital television. Dr Grgić is a member of the international program and organising committees of several international workshops and conferences. She was a visiting researcher at the Department of Telecommunications, University of Mining and Metallurgy, Krakow, Poland, working with Professor Janus Filipiak and Professor Krzysztof Wajda. Dr Grgić is the recipient of the silver medal "Josip Lončar" from the Faculty of Electrical Engineering and Computing in Zagreb for an outstanding PhD thesis work.

**Mohammed Ghanbari** received the BSc degree in electrical engineering from the Aryamehr University of Technology, Tehran, Iran, in 1970, the MSc degree in telecommunications, and the PhD degree in electronics engineering, both from the University of Essex, UK, in 1976 and 1979, respectively. After working almost 10 years in industry, he started his academic career as a Lecturer in the Department of Electronic Systems Engineering, University of Essex, in 1988, and was promoted to Senior Lecturer, Reader, and then Professor in 1993, 1995, and 1996, respectively. He is currently a Professor in the Department of Electronic Systems Engineering, University of Essex. His research interests include video compression and video networking. He is best known for his pioneering work on two-layer video coding for ATM networks. He has published more than 240 technical papers and registered for 6 international patents, in the areas of video networking. He is the author of Video Coding: An Introduction to Standard Codecs (London, UK: IEE Press, 1999). Dr Ghanbari has been a member of the organising committees of several international workshops and conferences. He was the Chairman of the steering committee for the 1997 International Workshop on Packet Video, and is currently an Associate Editor for IEEE Transactions on Multimedia. He is the co-recipient of the 1995 A. H. Reeves Premium Prize for the year's best paper published in the IEE Proceedings on the theme of digital coding. He is a Fellow of IEE and a Fellow of IEEE.