

Hierarchical Data Model in Content-based Image Retrieval

Quoc Ngoc Ly, and Anh Duc Duong

VietNam National University – Ho Chi Minh City, University of Natural Sciences
227 Nguyen Van Cu Street, 5 District, Ho Chi Minh City, VietNam

lqn@hcm.vnn.vn

Abstract

Nowadays, we are living in the content-based visual information retrieval age. The users would like the query data with the description close to human beings and the resulting images should have the same semantics meanings with query image. Having resolved this problem, we used the high-level feature as regions to fill the gap between low-level features and semantics meanings of images. We used an Hierarchical Agglomerative Clustering algorithm [HAC] to segment images in database into the regions and to classify them into the clusters with their representative that we called the words of the images. The query data can be the whole image, some regions of query images, some cluster's representatives in the image database or the concept related to some description text of the regions in the region's clusters. Retrieving is performed on the hierarchical tree structure of the region's clusters or hierarchical tree structure of the images. Our experiment results have shown that the performance of our system is better in the meanings of precision and recall than the traditional systems only based on using the whole image with the low-level features as query data and linear search or image retrieval system based on automated text annotation.

Keyword: hierarchical agglomerative clustering, low-level features, semantics meaning, hierarchical tree structure..

I. Introduction

Visual perception is the act of sensing a scene, of recognizing it and of describing it with symbol. While humans perform visual perception effortlessly and robustly, visual perception is still a major challenge for artificial vision systems. We deal with the visual perception problem under the artificial vision point of view. The aim of the artificial vision is not to reproduce the mechanisms of human vision but rather to use its proper mechanisms to be close to the results and the performance of human vision. When using the content-based image retrieval system, human beings always wish for easily representing contents that they would like to retrieval. The first approach is that using the query text based on manual image annotation, automatic annotation, the frequency of occurrence of the semantic concepts, visual concept ontology, the main shortcomings in these methods are that

they does not explicitly treat semantics as image classes and therefore, provides little guarantees that the semantic annotations are optimal in retrieval sense, it based on a small vocabulary, and many mental images cannot be represented by query text. The second approach is that using the image query based on the low-level features of images, high-level features as regions and following that, a relevance feedback mechanism is used to produce the final query results. The main shortcomings of these methods are that the users must have an whole image, existing the gap between the low-level features and semantic meanings of images and relevance feedback mechanism can create the conflict results during retrieving by many users.

Our approach is that combining the two kind of query, query image and query text. With query image, if the user don't have a desired image, they can choose the region cluster's representatives, it is useful because they help the user preview the compact view of image database, exploit the real content of them. We choose the HAC algorithm to segment and classify the image database because the hierarchical structure represents the essence of things. Retrieving is performed on region cluster's the hierarchical tree structure [RC's HTS] or image cluster's the hierarchical tree structure [IC's HTS]. The resulted images are ranked appropriately for meanings of query. With the query text, the user can choose from the list of text descriptions of some regions that are already manual annotated or the concepts can be expressed by them. At last, the user can combine two kind of query, the query image help to express some visual concept that is difficult to describe by text and the query text help to express some semantics meanings that is impossible to show by image.

This paper is organized into the following sections. The second section represent how to use the hierarchical clustering algorithm to segment the image database into the regions. The third section represent how to cluster the regions into RC's HTS, the fourth section represent how to cluster the images into IR's HTS. The fifth section represent how to annotate the regions. The sixth section represent the retrieving process. The seventh section represent our experiment results and the final section is conclusions and future work

II. Image Segmentation

In this stage, each image is segmented by HAC algorithm with the feature vectors of 4x4 blocks and the distance metric between blocks.

Each image I (in image database or query image) is partitioned into 4x4 blocks. The feature vector of each block consisted of 3 color features, 3 texture features and its centroid.

We used HSV color space and quantized color space into 12 hue components, 3 saturation components and 3 intensity components.

Each block has 3 color features, they are denoted as $\{C_1, C_2, C_3\}$:

$$C_1 = \left(\sum_{i=0}^{15} H[i] \right) / 16, C_2 = \left(\sum_{i=0}^{15} S[i] \right) / 16, C_3 = \left(\sum_{i=0}^{15} I[i] \right) / 16,$$

where $H[i], S[i], I[i]$ is respectively Hue, Saturation, Intensity components of the i -th pixel of block

Each block has 3 texture features, they are denoted as $\{T_1, T_2, T_3\}$.

Apply a Haar wavelet transform[7] to the value component of the image, each 4x4 block is decomposed into four frequency bands, each band contains 2x2 coefficients. We used 3 frequency bands : HL band, LH band and HH band.

Denoted 4 coefficients of block in HL band as $\{hl_0, hl_1, hl_2, hl_3\}$, T_1 is computed as $T_1 = \left(\frac{1}{4} \sum_{i=0}^3 hl_i\right)^{1/2}$.

Denoted 4 coefficients of block in LH band as $\{lh_0, lh_1, lh_2, lh_3\}$, T_2 is computed as $T_2 = \left(\frac{1}{4} \sum_{i=0}^3 lh_i\right)^{1/2}$

Denoted 4 coefficients of block in HH band $\{hh_0, hh_1, hh_2, hh_3\}$, T_3 is computed as $T_3 = \left(\frac{1}{4} \sum_{i=0}^3 hh_i\right)^{1/2}$

The i -th block has its centroid Cen_i with coordinates $\{x_i, y_i\}$, we use position features to form the connected region after segmentation.

Denoted f_i as a feature vector of block K_i of the image I ,

$$f_i = \{C_{i1}, C_{i2}, C_{i3}, T_{i1}, T_{i2}, T_{i3}, x_i, y_i\},$$

$C_{ij}, j \in [1..3]$ are color features after normalization,

$T_{ij}, j \in [1..3]$ are texture features after normalization,

x_i, y_i are position features after normalization.

Distance metric between two blocks in image is denoted as :

$$d_B(K_i, K_j) = (w_c \sum_{l=1}^3 (C_{il} - C_{jl})^2 + w_t \sum_{l=1}^3 (T_{il} - T_{jl})^2 + w_p ((x_i - x_j)^2 + (y_i - y_j)^2))^{1/2},$$

where

C_{il} are color features of block K_i ,

T_{il} are texture features of block K_i ,

x_i, y_i are coordinates of centroid of block K_i ,

w_c, w_t, w_p are weights reflecting the importance role of color, texture, position features

Image segmentation is really clustering pixels by its features. We used HAC algorithm because it could partition image into regions with the number of regions are not predefined and could manage the region by hierarchical structure from coarse to fine. After this stage, we have the regions carrying some meanings of image.

III. Region Classification

After segmenting image into the regions which are the high-level features which should help us narrowing the gap between the low-level features and semantics meanings of image. We have clustered the regions to extract the most common concept of the regions for support efficiently retrieval. We have clustered the regions based on the color features, texture features and shape features, dissimilar metric distance and HAC algorithm.

Color features

Region $Re g_i$ consisted of many blocks K_i , we haven't used the average color features as block's color features but we have used region's autocorrelogram [5] so we would like to capture both global occurrence statistics and local spatial organization of colors in region.

We used HSV color space and quantized color space into 12 hue components, 3 saturation components and 3 intensity components.

Let $[D]$ denote a set of D fixed distances d_1, \dots, d_D , which are measured using L_∞ norm, AutoCorrelogram of region $\text{Re } g$ is defined for color c and distance d as :

$$\alpha_c^d(\text{Re } g) = \Pr[p_1, p_2 \in \text{Re } g_c \mid |p_1 - p_2| = d],$$

where

$\text{Re } g$ is a region of image,

$\text{Re } g_c = \{p \in \text{Re } g \mid \text{Color}(p) = c\}$, $\text{Color}(p)$ denote color c of pixel p in $\text{Re } g$.

Denote $\alpha_{\text{Re } g_i}$ as AutoCorrelogram of $\text{Re } g_i$

Texture features

Suppose each region is consisted of N blocks.

Denoted texture features of block K_l as $\{T_{l1}, T_{l2}, T_{l3}\}$,

Texture features of region $\text{Re } g_i$ as $T_{\text{Re } g_i} = \{T_{\text{Re } g_i 1}, T_{\text{Re } g_i 2}, T_{\text{Re } g_i 3}\}$,

Texture features of region $\text{Re } g_i$ are computed as

$$T_{\text{Re } g_i 1} = \frac{1}{N} \left(\sum_{l=1}^N T_{l1} \right), T_{\text{Re } g_i 2} = \frac{1}{N} \left(\sum_{l=1}^N T_{l2} \right), T_{\text{Re } g_i 3} = \frac{1}{N} \left(\sum_{l=1}^N T_{l3} \right),$$

where

T_{l1}, T_{l2}, T_{l3} are texture features of block K_l ,

$T_{\text{Re } g_i 1}, T_{\text{Re } g_i 2}, T_{\text{Re } g_i 3}$ are texture features of region $\text{Re } g_i$

Shape features

Shape features are the high-level features having an important role in retrieving image more appropriately for semantics meanings.

The first feature described the area of region as

$$\text{Area}_{\text{Re } g_i} = \frac{|\text{Re } g_i|}{|I|},$$

$|\text{Re } g_i| = \sum_{\{x,y\} \in \text{Re } g_i} 1$, is the number of pixels in region $\text{Re } g_i$,

$|I| = \sum_{\{x,y\} \in I} 1$, is the number of pixels in image I .

The second feature described the location of the region as its centroid, this point could be not stayed in the region, however this matter is not important because distance metric between two regions is not dependent on this matter. Denote the location of the region as

$$\text{Pos}_{\text{Re } g_i} = \left(\sum_{(x,y) \in \text{Re } g_i} x / |\text{Re } g_i|, \sum_{(x,y) \in \text{Re } g_i} y / |\text{Re } g_i| \right),$$

x, y are the coordinates of pixel in region $\text{Re } g_i$,

$|\text{Re } g_i|$ is the number of pixels in region $\text{Re } g_i$

The third feature described compact property of the region as

$$Com_{Re\ g_i} = Per_C / (Per_{Re\ g_i}),$$

$Per_{Re\ g_i}$ is the perimeter of region $Re\ g_i$,

Per_C is the perimeter of a circle with the same area as the region Reg_i

Dissimilarity metric distance

We used HAC algorithm to cluster the regions. At first we clustered the regions by shape features, then by color features and finally by texture features.

Clustering the regions by shape features, we used the following dissimilarity metric distance as

$$d_{Re\ g} (Re\ g_i, Re\ g_j) = \alpha_{Area} |Area_{Re\ g_i} - Area_{Re\ g_j}| + \alpha_{Com} |Com_{Re\ g_i} - Com_{Re\ g_j}|$$

Clustering the regions by color features, we used the following dissimilarity metric distance as

$$d_{Re\ g} (Re\ g_i, Re\ g_j) = \sum_{c=1}^C \sum_{d=1}^D |\alpha_c^{(d)}(Re\ g_i) - \alpha_c^{(d)}(Re\ g_j)|$$

Clustering the regions by texture features, we used the following dissimilarity metric distance as

$$d_{Re\ g} (Re\ g_i, Re\ g_j) = |T_{Re\ g_i} - T_{Re\ g_j}|_{L_1}$$

After clustering the regions, we have the clusters $\{Clus_i, i = 1..NClus\}$ with its representative denoted as $\{Re\ p_i, i = 1..NClus\}$. This RC's HTS should be used for retrieving stage. It reduced so much time for retrieving because we must not exhausted all region database with distance metric on all the region's features.

IV. Image Classification

To explicitly treat semantics as image classes and therefore, provides guarantees that the semantic are optimal in retrieval sense, we try to cluster the images into IC's HTS. We replaced the regions in each image by the nearest representatives of the region clusters, clustered the images based on the features of the representatives, dissimilarity metric distance between two images and HAC algorithm.

Supposed that the image I is consisted of some regions as $\{Re\ g_i, i \in [1..Nreg]\}$, after replacing the regions by their representatives, the image I is consisted of some common concepts of regions as $\{Re\ p_j, j \in [1..Nrep]\}$.

We have clustered the images based on HAC algorithm with the metric distance between images as following :

Suppose we have two images I_i, I_j ,

$I_i = \{Re\ p_{i0}, Re\ p_{i1}, \dots, Re\ p_{iN_i-1}\}, Rep_{ik}, k \in [0..N_i - 1]$ are the representatives of I_i

$I_j = \{Re\ p_{j0}, Re\ p_{j1}, \dots, Re\ p_{jM_j-1}\}, Rep_{jl}, l \in [0..M_j - 1]$ are the representatives of I_j

$$d_{\text{img}}(I_i, I_j) = \frac{\sum_{k=0}^{N_i-1} w_{ik} d(\text{Re } p_{ik}, I_j) + \sum_{l=0}^{M_j-1} w_{jl} d(\text{Re } p_{jl}, I_i)}{2},$$

M_j, N_i is the number of representatives in image I_j, I_i ,

$d(\text{Re } p_{ik}, I_j) = \min\{d_{\text{Re } p}(\text{Re } p_{ik}, \text{Re } p_{jl}), l = 0..M - 1, \text{ is the distance between } \text{Rep}_{ik} \text{ of image } I_i \text{ to } I_j,$

$d(\text{Re } p_{jl}, I_i) = \min\{d_{\text{Re } p}(\text{Re } p_{jl}, \text{Re } p_{ik}), k = 0..N - 1, \text{ is the distance between } \text{Rep}_{jl} \text{ of image } I_j \text{ to } I_i,$

$w_{ik} = N_{ik} / NK_i$, N_{ik} is the number of blocks in Rep_{ik} , NK_i is the number of blocks in I_i

$w_{jl} = N_{jl} / NK_j$, N_{jl} is the number of blocks in Rep_{jl} , NK_j is the number of blocks in I_j

$w_{ik} w_{jl}$ reflecting the importance role of region having large area in distance metric.

Distance metric between two representatives is denoted as following :

$$d_{\text{Re } g}(\text{Re } p_i, \text{Re } p_j) = w_c d(\alpha_c^{(d)}(\text{Re } p_i), \alpha_c^{(d)}(\text{Re } p_j)) + w_t |T_{\text{Re } p_i} - T_{\text{Re } p_j}|_{L_1} + \\ + w_a |Area_{\text{Re } p_i} - Area_{\text{Re } p_j}| + w_{com} |Com_{\text{Re } p_i} - Com_{\text{Re } p_j}|$$

where

$$d(\alpha_c^{(d)}(\text{Re } p_i), \alpha_c^{(d)}(\text{Re } p_j)) = \sum_{c=1}^C \sum_{d=1}^D |\alpha_c^{(d)}(\text{Re } p_i) - \alpha_c^{(d)}(\text{Re } p_j)|,$$

$\alpha_c^{(d)}(\text{Re } p_i)$ is the c -th and d -th element of autocorrelogram of representative Rep_i

The representatives are participated in metric distance, the more similarity between the representatives, the more similarity between the two images.

V. Manual Region Annotation

In this stage, we try to attach the semantic words to the regions in RC's HTS. The annotation process is proceed after segmentation process and clustering. In the region cluster, we just only choose the region having the visual perception appropriately to semantic meaning to attach the word. This process don't use much labor as comparison as with the manual annotation in the automatic image annotation and retrieval model. At last we have prepared already a list of words support the retrieving process based on semantics afterward.

VI. Retrieving Process

In this stage, we can retrieve based on RC's HTS and IC's HTS. We can give the query image and query text.

Query by image based on RC's HTS

In this step, the user can input data by three different ways as following :

In the first way, the user can select the query image and then the system should automatically segment this image into the regions. These regions will be used in retrieving stage.

In the second way, the user can select some query images and then the system should automatically segment these images into the regions, and finally the user should select these query regions and combine them by logical operators as or, and, not. The typical query data is represented as following

$$(Sel(Reg_i) OR Sel(Reg_j)) AND (Sel(Reg_k) OR Sel(Reg_l)) AND NOT (Sel(Reg_m) OR Sel(Reg_n))$$

where $Sel(Reg)$ return boolean value True if Reg is selected

In the third way, if the user don't have already the query image, they can select the regions in RC's HTS and combine them by logical operators as or, and, not. The typical query data is represented as following :

$$(Sel(Reg_i) OR Sel(Reg_j)) AND (Sel(Reg_k) OR Sel(Reg_l)) AND NOT (Sel(Reg_m) OR Sel(Reg_n))$$

where $Sel(Reg)$ return boolean value True if Reg is selected

With a given query region Reg_q , we can compare it with the representatives of region's clusters in hierarchical tree structure to choose the nearest neighborhood.

Denoted the representative of region cluster matching with query region as Rep_r , cluster with its representative Rep_r as $Clus_r$.

Denoted the set of images having at least one region belongs to cluster $Clus_r$ as $SIClus_r$. The resulting images from query form as :

$$(Sel(Reg_i) OR Sel(Reg_j)) AND (Sel(Reg_k) OR Sel(Reg_l)) AND NOT (Sel(Reg_m) OR Sel(Reg_n))$$

should be :

$$(SIClus_i \cup SIClus_j) \cap (SIClus_k \cup SIClus_l) \setminus (SIClus_m \cup SIClus_n)$$

Finally, we rank the resulting images based on distance metric between two images.

Based on this approach, we should have a high precision, high recall, high speed of searching because we don't compare the query region with all region database and don't use all the features in distance metric. But the shortcomings is the resulting images can be empty.

Query by image based on IC's HTS

In this step, the user can input data by three different ways as before, then these query regions should be the region cluster's representatives or replaced by the region cluster's representatives based on RC's HTS. We compare the query image with image clusters's representatives on IC's HTS to choose the nearest neighbor. The set of images being in the cluster with the nearest representatives are the candidate images. Finally, we rank the the candidate images by the metric distance between images to have the resulting images. Retrieving based on IC's HTS can overcome the shortcomings of retrieving based on RC's HTS, the resulting images can't be empty.

Query by text based on RC's HTS

In this step, the user can choose the keywords from the predefined list of words. With each keyword, we can find exactly the region having the semantics meanings of keyword and the cluster including it. Then this region is used as query region and the process is proceed as query by region with the candidate regions in this cluster. The resulting regions are most similar regions to query region. The candidate images including the resulting regions are ranked to have the resulting images.

Query by image and text based on RC's HTS

In this step, the user can choose the query image and the keyword as before. Combine two kind of query help the user can represent their idea following the visual features and semantics meanings. The retrieving process are proceed as such as query by image and query by text.

VII. Experiments Results

We have experimented on the image database consisted of 9560 different kind of images as :

Natural scene: landscape, flower, fruit, animals.

Artificial scene : City, car, train, ship, plane.

Man activities : Sports activities

From 9560 images. We extracted 47820 regions, and clustered them in 205 clusters to construct the RC's HTS and clustered the images in 125 clusters to construct the IC's HTS.

The list of words (are manual annotated) are consisted of 120 words as :

Cloud, beach, ...; rose, sunflower, gladiolus,...; blue dragon, durian, banana,...; bear, tiger, monkey,....

We used two value Precision and Recall to evaluate the performance of our system.

Precision = Number of relevance detected images / Total number of detected images.

Recall = Number of relevance detected images / Total number of relevance images

Table 1. Query by image based on RC's HTS

Input method and retrieval method	Number of query images	Number of resulting images	Average Precision	Average Recall
Query image and retrieving based on its global histogram	100	10	78	57
		20	73	59
		50	70	60
All the regions of query image and retrieving based on region features	100	10	80	64
		20	79	65
		50	77	68
Some regions of query images and retrieving based on region features	100	10	85	73
		20	82	76
		50	80	77
Some regions in RC's HTS and retrieving based on region features	100	10	86	73
		20	83	77
		50	80	78

Table 2. Query by text based on RC's HTS

Input method and retrieval method	Number of query text	Number of resulting images	Average Precision	Average Recall
Query text and retrieving based on RC's HTS	50	10	91	80
		20	88	82
		50	84	83

The experiments showed that this approach can narrow the gap between the low-level features and the semantics meanings.

VIII. Conclusions

We have approached the content-based image retrieval by using the high-level feature as region to fill in the gap between low-level features and semantics meanings of images. We have represented our efficient and user-friendly image retrieval system, this system is a part of our system named VIROS (Visual Information Retrieval of Saigon). In the future, we shall develop our system to the concept-based image retrieval system.

References

- [1] Al Bovik, *Handbook of Image and Video Processing* : Academic Press, 2000.
- [2] D.Bimbo, *Visual Information Retrieval* : Morgan Kaufmann, 1999.
- [3] G.Sheikholeslami, W.Chang, A.Zhang "SemQuery : Semantic Clustering and Querying on Heterogeneous features for Visual data" in *IEEE Trans. Knowledge and Data Engineering* vol. 14, no.5: IEEE, 2002, pp. 988-1002.
- [4] I. Kompatsiaris, E. Triantafillou, M.G. Strintzis, "Region-based color image indexing and retrieval" in *IEEE International Conference on Image Processing (ICIC'01)*: IEEE, 2001
- [5] J. Huang, S.R. Kumar, M. Mitra, W.J. Zhu, R.Zabih, "Image indexing using color correlograms" in *IEEE International Conference on computer Vision and Pattern Recognition (CVPR'97)*: IEEE, 1997
- [6] J. Jeon, V.Lavrenko, R.Mammatha, "Automatic Image Annotation and Retrieval using Cross-Media Relevance Models" in *SIGIR '03*: ACM, 2003
- [7] P.Salembier, F.Marques, "Region-based representations of image and video : Segmentation tools for multimedia services" in *IEEE Trans. On Circuits and Systems for Video Technology*, vol.9, no.8 : IEEE, December 1999
- [8] Rafael C. Gonzalez, Richard E.Woods, *Digital Image Processing Second Edition* : Prentice-Hall Inc., 2002.

- [9] Sergios Theodoridis, Konstantinos Koutroumbas, *Pattern Recognition*: Academic Press, 1999.
- [10] Vasileios Mezaris, Ioannis Kompatsiaris, Michael G. Strintzis “An ontology approach to object-based image retrieval” in ICIP’03: IEEE, 2003



My name is Quoc Ngoc Ly, I was born in September, 1st, 1964 in South VietNam. I got the Bachelor degree in Mechanical Mathematics in 1987, the Master degree in Computer Science in 1995 and now I am a Phd.Student at the University of Natural Sciences of Ho Chi Minh City, VietNam. My current research are Image Processing, Computer Graphics, Computer Vision, Artificial Intelligent and Applications in Architecture, Advertising, Animations.



Duong Anh Duc obtained B.S. in CS in 1990. He graduated M.S. in 1995, and Ph.D. in 2002 from University of Natural Sciences, Vietnam National University - Hochiminh city. At present, he is currently the dean of Faculty of IT and the director of Software Engineering Laboratory. His current research are Geography Information Systems, Cryptography and security, Computer graphics and image processing, and Software Engineering.