

Spatial Cloud Detection and Retrieval System for Satellite Images

Noureldin Laban¹, Ayman Nasr¹

¹Department of Image Processing and its Applications,
Data Reception and Analysis Division, National
Authority for Remote Sensing and Space Sciences ,
P.O.Box 1564, Alf Maskan, Cairo, Egypt

Motaz ElSaban² Hoda Onsi²

²Department of Information Technology,
Faculty of Computers and Information, Cairo University,
P.O.Box 12613 , Orman, Giza, Egypt

Abstract—In last the decade we witnessed a large increase in data generated by earth observing satellites. Hence, intelligent processing of the huge amount of data received by hundreds of earth receiving stations, with specific satellite image oriented approaches, presents itself as a pressing need. One of the most important steps in earlier stages of satellite image processing is cloud detection. Satellite images having a large percentage of cloud cannot be used in further analysis. While there are many approaches that deal with different semantic meaning, there are rarely approaches that deal specifically with cloud detection and retrieval. In this paper we introduce a novel approach that spatially detect and retrieve clouds in satellite images using their unique properties .Our approach is developed as spatial cloud detection and retrieval system (SCDRS) that introduce a complete framework for specific semantic retrieval system. It uses a Query by polygon (QBP) paradigm for the content of interest instead of using the more conventional rectangular query by image approach. First, we extract features from the satellite images using multiple tile sizes using spatial and textural properties of cloud regions. Second, we retrieve our tiles using a parametric statistical approach within a multilevel refinement process. Our approach has been experimentally validated against the conventional ones yielding enhanced precision and recall rates in the same time it gives more precise detection of cloud coverage regions.

Keywords-Satellite images; Content based image retrieval; Query by polygon; Retrieval refinement; cloud detection; geographic information system.

I. INTRODUCTION

Satellite images have become a common component of our daily life either on the Internet, in car driving and even in our hand-held mobile handsets. There is huge image content appearing every second through multiple competing satellite systems [1]. Manual interaction with this large volume of data is becoming more and more inappropriate, which creates an urgent need for automatic treatment to store, organize and retrieve this content [2].

Traditional textual meta-data such as geographic coverage, time of acquisition, sensor parameters, manual annotation, etc., are now insufficient to retrieve images of interest when we target a specific visual concept such as desert, rock, crops, clouds or others [3]. In many fields, we need specific contents from the satellite images as specific crops, geology structures or climate changes.

Manual annotation needs to annotate every region by human where users enter descriptive word after image download from satellite. However it is a labor intensive and tedious process [4]. Therefore we need to retrieve images that contain our intended contents automatically. The content based image retrieval (CBIR) approach challenge is how to fill the gap between the low level features that describe the scenes and our human understandable semantic concepts. This gap of understanding is called the semantic gap [5] [6]. In addition, these semantic concepts themselves may be defined differently, e.g. each one of us interprets what he sees from his point of view.

The most commonly used features include those reflecting color, texture, shape, and salient points in an image. For instance, in a color layout approach, an image is divided into a small number of sub-images and the average color components (e.g. red, green, and blue intensities) are computed for every sub-image [7]. Texture features are intended to capture the granularity and repetitive patterns of surfaces within an image.

The traditional satellite cloud image search method was based on the file name and the sensor parameters of every image. The disadvantages of this method are that it cannot describe the image contents such as cloud shape [8] and also leads to the inconvenience in retrieving images [9].

We have done statistics for Spot4 satellite observation on the Middle East from NARSS archive to determine the percent of clouds on these scenes in the period starts from January 2006 to December 2009. There were about 170000 scenes covering the receiving station area. Normally for each scene; an expert has to decide manually the percentage of cloud coverage.

The different percentages of clouds coverage during each year are shown in figure 1 and table I.

TABLE I: AVERAGE CLOUD COVERAGE THOUGH 2006 TO 2009 ON MIDDLE EAST

COVERAGE	2006	2007	2008	2009
0% (A)	0.44	0.39	0.43	0.40
1%-10% (B)	0.08	0.09	0.07	0.06
11%-25% (C)	0.09	0.08	0.07	0.05
26%-75% (D)	0.09	0.08	0.07	0.06
76%-100% (E)	0.30	0.36	0.36	0.43

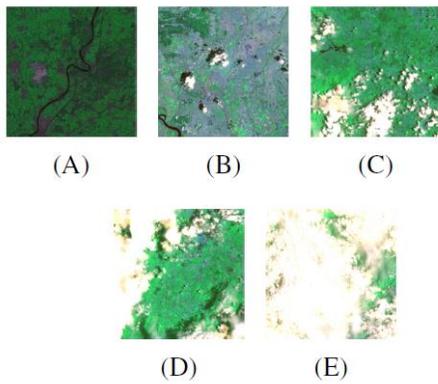


Fig. 1 Clouds coverage percentages

II. REVIEW OF RELATED WORK

During the last decade many approaches have been proposed to retrieve satellite images using their content in general. Specifically less effort has been devoted to cloud despite its importance during satellite image processing or meteorological management and observation. F. Acqua and P. Gamba presented a tool for shape similarity evaluation for query-by shape searching into meteorological image archives based on the point diffusion technique [8]. R. Holowczak et al., reported a system that can automatically determine whether a region of interest is visible in the image, free from cloud, and can incorporate this into the meta-data for individual images to enhance searching capability [10]. T. Nauss et al., have proposed an algorithm based on the analytical solutions of the radiative transfer equations valid for optically thick weakly absorbing cloud layers [11]. D. Fu and L. Xu have used 2D-Gabor wavelet in satellite image classification [12]. D. Upreti has used Gray level Co-occurrence Matrix GLCM and histogram quantization technique to retrieve cloud patterns to discover Tropical Cyclone [13].

The previous approaches were concerned with cloud retrieval. Some observations were found as follow:

- Most of the previous work was directed to meteorological observation images with very low resolution.
- It doesn't care with cloud removal preprocessing operation which is still done manually.
- It doesn't handle spatial distribution of cloud within the scene.

Through our new proposed approach, we covered these missed points of research. It will be very useful to detect and retrieve these clouds and consequently as further process, remove them and replace the cloudy sub-images with other clear ones.

III. SYSTEM OVERVIEW

Our system is composed of two main stages as shown in figure 2.

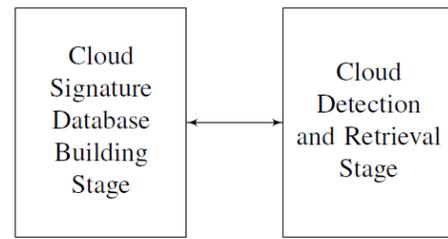


Fig. 2 System Overview

First stage is cloud signature database building stage which is responsible for building up the features vectors for different clouds patterns. Second stage is cloud detection and retrieval stage in each satellite scene, which determines where the clouds in this scene and their percentage are. We have used two strategies in our system [1]. First one is query by polygon strategy where we build our signature database using cloud polygons instead of rectangular shapes. Second one is multiple size tiling strategy where we break down our scene into different sizes followed by features extraction to obtain features vectors. According to these strategies, the two stages have passed through different sub-processes starting by tiling then features extraction to from features vectors. This is done for each level of retrieval.

IV. SYSTEM STAGES

A. Cloud Signature Database Building stage

There are many forms that clouds appear with in satellites images as shown in figure 3.

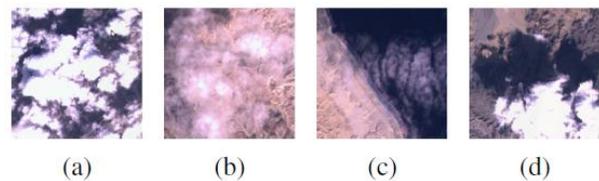


Fig. 3 Some Clouds types

These forms differ depending on altitude and density of clouds [14]. These forms start with low dense water vapor to high dense clouds with different altitude. Beside clouds there are also their shadows which should be taken into account during retrieval. The first stage of the cloud retrieval process is to determine cloud signature as shown in figure 4.

This is done using query by polygon approach where we first determine different type of clouds, then we draw geo-reference polygons that contain these clouds. These different types of clouds are used to form signature databases according to the type of tiling size used. Using our proposed feature extraction algorithm we compute features vectors of cloud polygon tiles

B. Cloud Detection and Retrieval Stage

After building our cloud signature database, we have to build the features vectors for each scene as shown in figure 5.

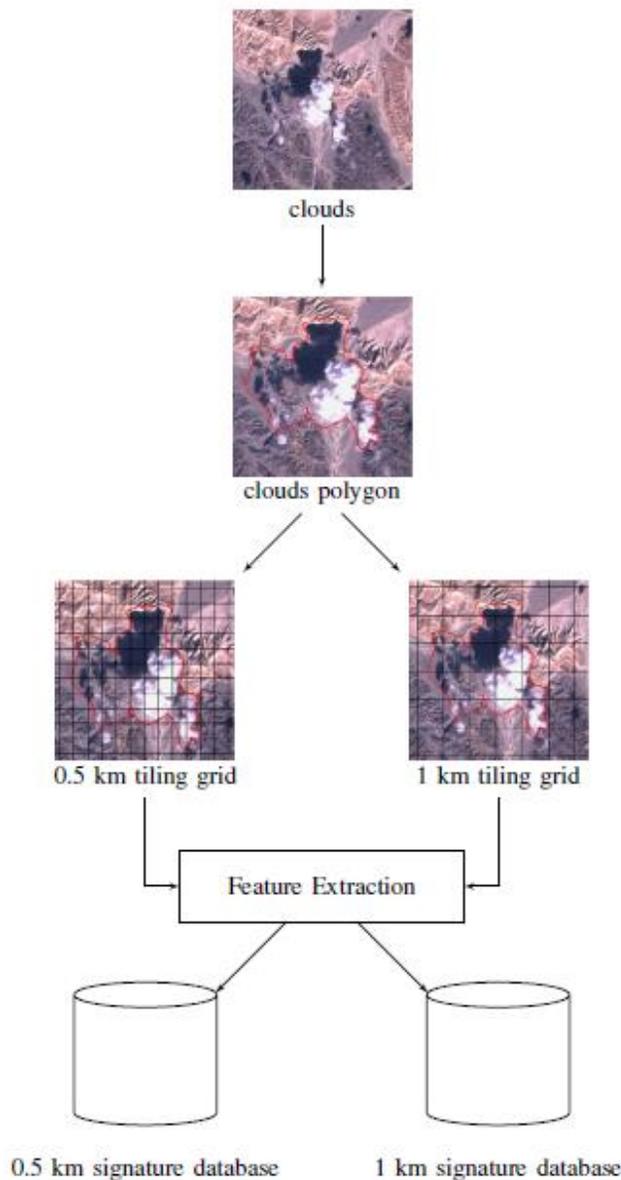


Fig. 4 Building cloud signature databases

Our approach based on breaking down the whole image into small sections of sub-images called tiles. The number of resulted tiles is determined by their sizes.

According to the two stages hierarchy used in [1] for the retrieval process, we have rebuilt the system. Instead of starting with features databases and get query features for each semantic, we have reverse the order which begins with building cloud signature database then the input scene is treated as query image. The two stages hierarchy, candidate selection stage and refinement stage, are used. In candidate selection stage, we define the primary candidate's area for clouds. In refinement stage we refine the first stage areas using its neighborhoods with smaller tile size.

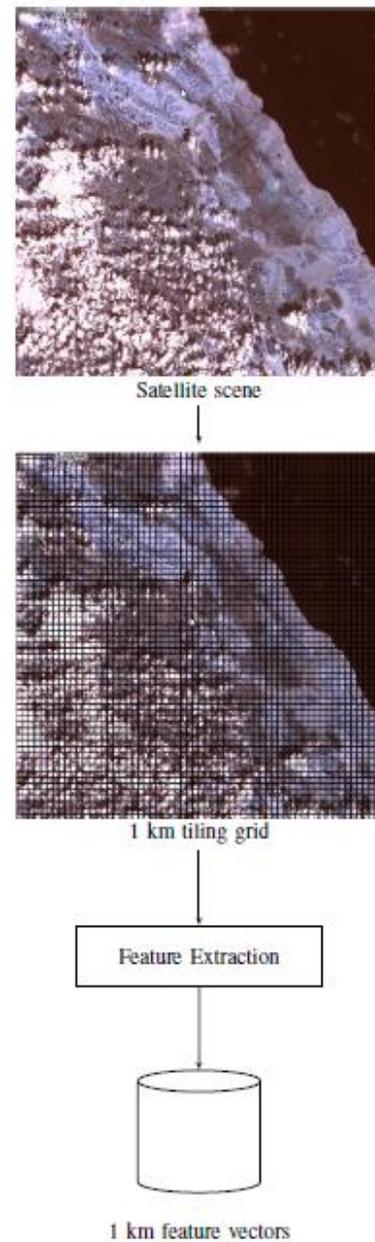


Fig. 5 Building satellite feature vectors

V. MAIN SYSTEM PROCESSES

The main system processes; features extraction, retrieval and evaluation have some key points to be included into the two levels hierarchy to enhance cloud detection and retrieval system.

A. Features extraction Process

We have depended on various domains to get the tile signature either for the cloud example dataset or for input satellite image. These domains extract the spectral and textural characteristics of images.

To build our feature vector database $DB_{features}$, we start by determining the components of our feature vector V_i for each tile i and its length. For each multispectral tile image T_i with number of bands, we form the feature vector V_{ib} of each band b depending on different spectral and textual characteristics of the image. We used the mean μ and standard deviation σ statistics of feature domain for each band. The features we used are histogram H , Daubechies wavelets transform coefficients DWT, Discrete cosine transform coefficients DCT and Discrete Fourier transform Coefficient DFT [15]. Using these domains, we build various feature vectors V_H , V_{DCT} and V_{DFT} . For each multispectral tile with n number of bands, we build domain feature vector V_d for each domain d as in equation 1.

$$V_d = [V_{d1}, V_{d2}, \dots, V_{dn}] \quad (1)$$

We then use these domain feature vectors to form domain feature database DB_d for m number of tiles as in equation 2.

$$DB_d = [V_{d1}, V_{d2}, \dots, V_{dm}] \quad (2)$$

Using all feature vectors for all tiles; we formulate our cloud signature database or input scene feature vectors using all domains as in equation 3.

$$DB_{feature} = [DB_1, DB_2, \dots, DB_d] \quad (3)$$

B. Retrieval Process

The retrieval process, as shown in figure 6, has two sub stages as mentioned in [1], the candidates selection stage and the refinement stage.

In the candidates selection stage, we use 1 km tile size features to get the most appropriate matching tiles similar to cloud. In the refinement stage, we use the 0.5 km tile size features of the first stage results and their neighborhoods to get our final results.

We have used a retrieval engine that based on statistical parametric paradigm using normal distribution [16] rather than the traditional nearest neighbor approach. The statistical parametric paradigm aimed to determine the parameters of the statistical distribution that the data follows as mean μ and standard deviation σ . We define the the training dataset $D_{training}$ that represent cloud example tiles set D_{cloud} and non cloud example tiles set $D_{non\text{ cloud}}$ as in equation 4.

$$D_{training} = [D_{cloud}, D_{non\text{ cloud}}] \quad (4)$$

This is done for every tile size. Therefore, our global signature data D_{global} is formed from all sizes used in our system as in equation 5.

$$D_{global} = [D_{size\ 1}, D_{size\ 2}, \dots, D_{last\ size}] \quad (5)$$

After we have built our statistical model using $D_{training}$, SCDRS is now ready to receive the satellite images as an input.

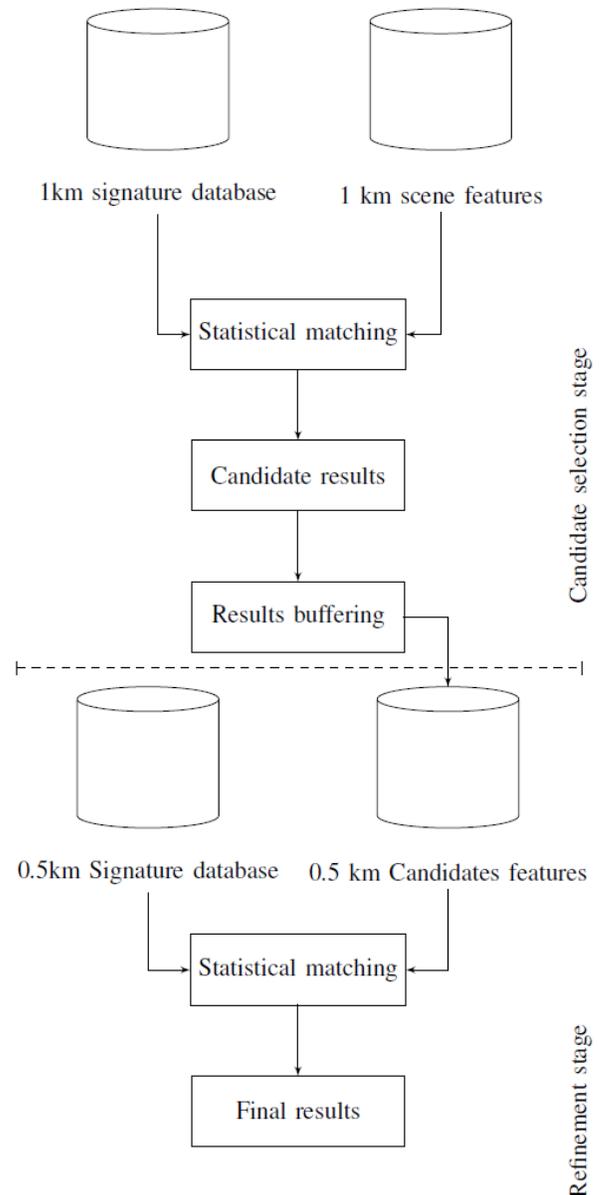


Fig. 6 Two Stages Retrieval Process

C. System Evaluation Process

Our evaluation process is carried out in terms of recall and precision (equations 6, 7 respectively) using relevant areas in the database.

$$recall = \frac{\text{correctly retrieved cloud area}}{\text{actual cloud area}} \quad (6)$$

$$precision = \frac{\text{correctly retrieved cloud area}}{\text{retrieved cloud area}} \quad (7)$$

We use the map coordinates (i.e. Latitude and Longitude) instead of using file coordinates (pixels). As the map coordinates is universal and continuous where the file coordinates is file specific. The global coordinate system is independent from the pixel size whatever the scanning satellite or stored file. So the percent of cloud area in the input scene is as shown in equation 8

$$\text{retrieved cloud percent} = \frac{\text{retrieved cloud area}}{\text{whole scene area}} \quad (8)$$

where the actual cloud percent retrieved is calculated as shown in equation 9

$$\text{actual cloud percent} = \frac{\text{correctly retrieved cloud area}}{\text{whole scene area}} \quad (9)$$

VI. EXPERIMENTAL RESULTS

On our experiments we have used Spot4 satellite scenes with different cloud cover percents which cover about 10800 km^2 . Each scene covers $60 \text{ km} \times 60 \text{ km}$ of earth surface in Egypt with pixel size of 20 m . We used also Landsat archive images database with different cloud coverage percentages. There scenes cover about 22400 km^2 with 30 m pixel size.

Each scene has been divided into sub images of $1 \text{ km} \times 1 \text{ km}$ and $0.5 \text{ km} \times 0.5 \text{ km}$. The experiment scenes have formed more than 100,000 sub-images which are pre-classified clouds images. We have used samples of different clouds types to form our cloud signature database which is composed of 110 sub images acting as clouds examples.

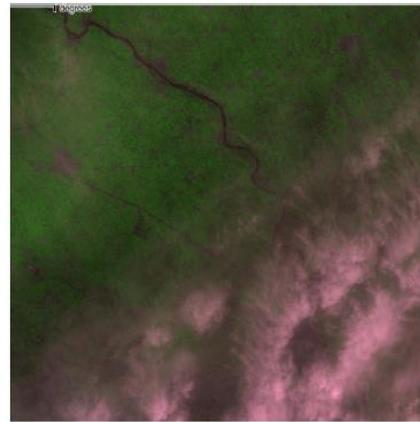
VII. RESULTS ANALYSIS

For our semantic concept which is cloud; first we have used two categories of polygons shapes, one used for building cloud signature database and the other is tied with each input scene used for evaluation. An example result of our system is depicted in figure 7.

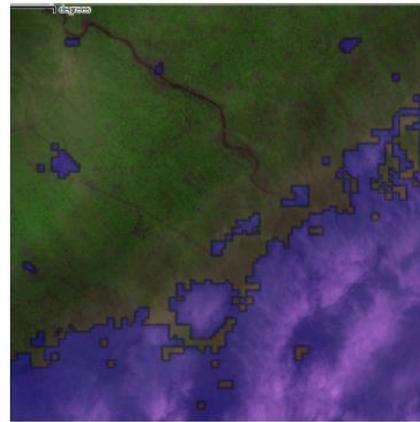
The results of each input scene could be evaluated by two ways. First, the input test polygon for cloud; which determines exactly the positions of clouds in this scene and the area of clouds compared to the whole scene area. Second, the expert's estimation used in ground station which estimate the range of cloud cover as explained in table I.

As shown in table II, results of the two successive stages of the system are presented. It shows how the different types of features domains affect the results.

To determine cloud percentage coverage, we have calculated the total area of output results cloud tiles with respect to the whole scene area which is 3600 km^2 as in equation 8. We have put into consideration that the most important parameter is precision as we should guarantee that the output results have to be more accurate and decrease the non clouds tiles resulted. So, when we select the cloud examples, it should be purely determined.



(a) Clouds as seen in the real scene



(b) Detected Clouds regions retrieved by SCDRS

Fig. 7 SCDRS result example

TABLE II : DIFFERENT RECALL AND PRECISION FOR TWO STAGE HIERARCHY

	first stage		second stage	
	Precision	Recall	Precision	Recall
Histogram	76	88	72	91
Wavelets	74	91	73	92
DCT	74	90	72	92
FFT	72	87	70	93

Table III shows the recall and precision results using the different feature domains. The accuracy of different features is very comparable. The results explain that the key point here is the processing time, which is recorded to histogram features as it is the least complex than the others. As the tile becomes more smaller the spectral characteristics become more sufficient than textural characteristics to distinguish between tiles.

TABLE III : DIFFERENT RECALL(R) AND PRECISION (P) FOR DIFFERENT TYPES OF FEATURES USING 0.5 KM TILE SIZE AND PROCESSING TIME (PT)

	Histogram		Wavelets		Discrete cosine		Fourier	
	P	R	P	R	P	R	P	R
A	88	94	89	100	84	100	89	100
B	66	90	68	90	67	86	68	90
C	86	75	75	75	80	75	79	69
D	97	82	97	76	97	79	97	79
E	100	97	100	97	100	97	100	100
PT(MIN)	22		45		28		24	

VIII. CONCLUSIONS

In this paper, a new approach was developed to detect the percentage of clouds and retrieve their positions within the satellite images using two stages; Cloud Signature Database Building stage and Cloud Detection and Retrieval Stage. The two stages used multilevel framework hierarchy of candidates selection and candidates refinement processes. This is done using spatial and textural features and parametric statistical approach for retrieval process. The capability of the developed system was tested using a dedicated satellite images and assessed in terms of cloud percentage coverage with the traditional precision and recall measurements. Results show that the developed system enhanced the precision and recall and in the same time it gives a closer assessment for cloud coverage to the real area calculations. They also show that the spectral features have higher accuracy than textural features. We propose as future work to represent a system for detecting different types of clouds using more robust retrieval algorithms which integrated with GIS systems.

REFERENCES

[1] N. Laban, M. ElSaban, A. Nasr, and H. Onsi, "System refinement for content based satellite image retrieval," *The Egyptian Journal of Remote Sensing and Space Sciences*, vol. 15, June 2012.

[2] M. Martins, L. Frutuoso Guimaraes, and L. Maria Garcia Fonseca, "Texture feature neural classifier for remote sensing image retrieval

systems," in *Computer Graphics and Image Processing*, 2002. Proceedings. XV Brazilian Symposium on, 2002.

[3] C.-R. Shyu, M. Klaric, G. Scott, A. Barb, C. Davis, and K. Palaniappan, "Geoiris: Geospatial information retrieval and indexing system mdash;content mining, semantics modeling, and complex queries," *Geoscience and Remote Sensing*, IEEE Transactions on, vol. 45, no. 4, pp. 839–852, April 2007.

[4] H. H. Wang, D. Mohamad, and N. A. Ismail, "Semantic gap in cbr: Automatic objects spatial relationships semantic extraction and representation," *International Journal Of Image Processing*, vol. 4, pp. 192–286, July 2010.

[5] I. Gondra and D. R. Heisterkamp, "Content-based image retrieval with the normalized information distance," *Computer Vision and Image Understanding*, vol. 111, no. 2, pp. 219–228, 2008.

[6] H. Min and Y. Shuangyuan, "Overview of content-based image retrieval with high-level semantics," in *3rd International Conference on Advanced Computer Theory and Engineering (ICACTE)*, 2010, vol. 6, Aug. 2010, pp. 312–316.

[7] R. Datta, D. Joshi, J. Li, and J. Z. Wang, "Image retrieval: Ideas, influences, and trends of the new age," *ACM Comput. Surv.*, vol. 40, no. 2, pp. 5:1–5:60, May 2008.

[8] F. Dell'Acqua and P. Gamba, "Query-by-shape in meteorological image archives using the point diffusion technique," *Geoscience and Remote Sensing*, IEEE Transactions on, vol. 39, no. 9, pp. 1834–1843, Sept. 2001

[9] W. ShangGuan, Y. Hao, Y. Tang, and Y. Zhu, "The research and application of content-based satellite cloud image retrieval," in *International Conference on Mechatronics and Automation. ICMA 2007.*, Aug. 2007, pp. 3864–3869.

[10] R. Holowczak, F. Artigas, S. A. Chun, J.-S. Cho, and H. Stone, "An experimental study on content-based image classification for satellite image databases," *Geoscience and Remote Sensing*, IEEE Transactions on, vol. 40, no. 6, pp. 1338–1347, June 2002.

[11] T. Nauss, A. Kokhanovsky, T. Nakajima, C. Reudenbach, and J. Bendix, "The intercomparison of selected cloud retrieval algorithms," *Atmospheric Research*, vol. 78, no. 12, pp. 46–78, 2005.

[12] D. Fu and L. Xu, "Satellite cloud image texture feature extraction based on gabor wavelet," in *Image and Signal Processing (CISP)*, 2011 4th International Congress on, vol. 1, Oct. 2011, pp. 248–251.

[13] D. Upreti, "Content-based satellite cloud image retrieval," Master's thesis, Faculty of Geo-Information Science and Earth Observation of the University of Twente, Enschede, The Netherlands, 2011.

[14] J. Oliver, Ed., *Encyclopedia of World Climatology*, ser. Encyclopedia of Earth Sciences Series. Springer, 2005.

[15] M. Petrou and C. Petrou, *Image Processing: The Fundamentals*. Wiley, April 2010.

[16] A. K. Jain, M. N. Murty, and P. J. Flynn, "Data clustering: a review," *ACM Computer Survey*, vol. 31, no. 3, pp. 264–323, Sept. 1999.