**STATISTICAL APPROACH FOR IMAGE RETRIEVAL**

**KHOR SIAK WANG**

**DOCTOR OF PHILOSOPHY**
**UNIVERSITI PUTRA MALAYSIA**

**2007**

# STATISTICAL APPROACH FOR IMAGE RETRIEVAL

By

## KHOR SIAK WANG

**Thesis Submitted to the School of Graduate Studies, Universiti Putra Malaysia, in Fulfilment of the Requirement for the Degree of Doctor of Philosophy**

**January 2007**

**Family, wife & sons**

*Abstract of thesis presented to the Senate of Universiti Putra Malaysia in fulfilment of the requirement for the degree of Doctor of Philosophy*

**STATISTICAL APPROACH FOR IMAGE RETRIEVAL**

**By**

**KHOR SIAK WANG**

**January 2007**

**Chairman:**     **Associate Professor Fatimah Bt. Dato' Ahmad, PhD**

**Faculty:**       **Computer Science And Information Technology**

*Since the emergence of Internet, a gigantic volume of images have been uploaded into the Internet from time to time. Relying on the traditional text-based search approach to locate the required images could no longer meet the diverse needs of users. This persistent trend has demanded a more sophisticated search algorithm on these images.*

*One of the popular and common approaches for image search is Content-based Image Retrieval or CBIR for short, i.e. retrieval of images based on their visual contents such as shapes, colours, textures etc.*

*Of all the visual contents identifiable from an image, colour is considered to be the commonest visual attribute that aids in image retrieval. Works on colour-based image retrieval systems are largely based on the use of colour histogram, which has been noted*

*to suffer from a major drawback, i.e. absence of spatial information, which is also an important requirement for an accurate retrieval result.*

*In this thesis, a novel method based on the modified generic framework of CBIR is proposed. This technique, formally known as Image Retrieval Using Statistical-based Approach is based on the idea of grouping pixels with similar colour codes within an image. From these grouped pixels, they are sorted in descending order of pixel count, which intuitively identifies dominant colours within an image. Statistical information, i.e. means and standard deviations will then be derived from these sorted groups. The extracted statistical information will be stored in both text files and matrixes, which will be used to aid in the image retrieval process. The system has also included some adjustable parameters, such as window size, CC percentage similarity, which can be used to improve retrieval accuracy. This statistical-based approach has been tested on the standard UCID image collection where it has shown improved results, with an average precision value of about 70% as compared to an approximate value of 25% using the histogram-based approach, in term of retrieval accuracy.*

*Abstrak tesis yang dikemukakan kepada Senat Universiti Putra Malaysia sebagai memenuhi keperluan untuk ijazah Doktor Falsafah*

**DAPATAN SEMULA IMEJ YANG BERDASARKAN KAEDAH STATISTIK**

**Oleh**

**KHOR SIAK WANG**

**Januari 2007**

**Pengerusi:**     **Profesor Madya Fatimah Bt. Dato' Ahmad, PhD**

**Fakulti:**        **Sains Komputer dan Teknologi Maklumat**

*Semenjak kewujudan Internet, terdapat banyak imej yang dimasukkan ke dalam Internet dari semasa ke semasa. Kaedah mendapatkan semula imej secara tradisional yang berdasarkan teks tidak dapat memenuhi keperluan para pengguna. Tren ini memerlukan kaedah pencarian imej yang sopistikated.*

*Salah satu daripada kaedah yang popular dan biasa untuk mendapatkan semula imej adalah "Content-based Image Retrieval" atau CBIR, iaitu kaedah mendapatkan semula imej berasaskan properti visual seperti bentuk, warna, tekstur dan lain-lain.*

*Dari semua properti visual yang terkandung di dalam imej, properti warna merupakan properti yang sering digunakan untuk mendapatkan semula imej. Kaedah biasa yang digunakan untuk dapatan semula imej berasaskan warna ialah penggunaan histogram. Kelemahan utama kaedah ini adalah kehadiran lokasi objek di dalam sesuatu imej tidak*

*dipertimbangkan. Pertimbangan kehadiran lokasi ini merupakan faktor yang penting untuk mendapatkan semula imej dengan tepat.*

*Dalam tesis ini, model CBIR yang tradisi akan diubahsuai. Kaedah yang dicadangkan dikenali sebagai Dapatan Semula Imej Yang Berdasarkan Informasi Statistik. Kaedah tersebut berdasarkan idea di mana semua pisel yang mempunyai kod warna yang seragam akan dikelompokkan. Kelompok-kelompok pisel ini akan disusun menurut saiznya. Dengan jelasnya, apabila kelompok tersebut telah disusun mengikut saiznya, ia juga memberi gambaran di mana warna dominan mudah ditentukan. Dari kelompok ini, informasi statistic, iaitu min dan penaburan piawai akan diperolehi. Maklumat tersebut akan disimpan di dalam fail dan array untuk membantu proses dapatan semula imej berasaskan warna. Sistem yang dicadangkan juga mempunyai parameter yang boleh digunakan oleh para pengguna untuk memperbaiki keputusan. Eksperimen yang dilaksanakan dengan menggunakan UCID data dapat menunjukkan kaedah yang dicadangkan mampu memberi keputusan ketepatan secara purata 70% ketepatan dibandingkan dengan 25% dengan menggunakan kaedah histogram.*

# ACKNOWLEDGEMENTS

*First and foremost, I would like to thank Associate Professor Fatimah Bt. Dato' Ahmad for giving me an opportunity to start off this project. I 'm indeed obliged by her enthusiastic support of the project from the very early stages. Without her tireless assistance and guidance, this project would never be completed on time. Also, without her constant monitoring and supervisions on the progress of my project, I believe that the contents of this project would still be bits and pieces stored in my hard disk. Her cooperation and contributions are indispensable.*

*Being a part-time student, I could hardly devote my precious time to my wife, Ms. Kwang Wai Ching, my 3-year old son Khor Hoong Yik and my new-born baby, Khor Hoong Yang, who have been very supportive and patiently waiting for me to complete my study.*

*Being one of the key persons in the supervisory committee team, Associate Professor Ramlan bin Mahmod is always tight with his schedules and daily events. He really looks serious but approachable. Without his serious-looking face, I would not be able to ensure my work is of the required quality and standard. Thanks, once again.*

*Associate Professor Hamidah bt. Ibrahim, who is delightful to work with, and always replies me with very short mail on my requests but straight to the point, has been helpful in giving me concrete and constructive comments of my work. I would like to gratefully acknowledge her contributions and her immense help and vast knowledge in database.*

*Finally, many thanks also go to some of my peer colleagues, where they prefer themselves not to be named, who have given me constructive comments and ideas in certain parts of my research work.*

*I certify that an Examination Committee has met on 26/01/2007 to conduct the final examination of Khor Siak Wang on his Doctor of Philosophy thesis entitled "Image Retrieval Using Statistical-based Approach" in accordance with Universiti Pertanian Malaysia (Higher Degree) Act 1980 and Universiti Pertanian Malaysia (Higher Degree) Regulations 1981. The Committee recommends that the candidate be awarded the relevant degree. Members of the Examination Committee are as follows: -*

**ALI MAMAT, PhD**
*Associate Professor*
*Faculty of Computer Science and Information Technology*
*Universiti Putra Malaysia*
*(Chairman)*

**RAHMITA WIRZA, PhD**
*Faculty of Computer Science and Information Technology*
*Universiti Putra Malaysia*

**SHYAMALA DORAISAMY, PhD**
*Faculty of Computer Science and Information Technology*
*Universiti Putra Malaysia*

**TENGKU MOHD TENGKU SEMBOK, PhD**
*Professor*
*Faculty of Information Science and Technology*
*Universiti Kebangsaan Malaysia*

**HASANAH MOHD. GHAZALI, PhD**
*Professor/Deputy Dean*
*School of Graduate Studies*
*Universiti Putra Malaysia*

*Date:*

*This thesis submitted to the Senate of Universiti Putra Malaysia and has been accepted as fulfillment of the requirement for the degree of Doctor of Philosophy. The members of the Supervisory Committee are as follows: -*

**Fatimah Dato' Ahmad, PhD**
*Associate Professor*
*Faculty of Computer Science and Information Technology*
*Universiti Putra Malaysia*
*(Chairman)*

**Ramlan Mahmod, PhD**
*Associate Professor*
*Faculty of Computer Science and Information Technology*
*Universiti Putra Malaysia*
*(Member)*

**Hamidah Ibrahim, PhD**
*Associate Professor*
*Faculty of Computer Science and Information Technology*
*Universiti Putra Malaysia*
*(Member)*

---

*AINI IDERIS, PhD*
*Professor/Dean*
*School of Graduate Studies*
*Universiti Putra Malaysia*

*Date:*

# DECLARATION

*I hereby declare that the thesis is based on my original work except for quotations and citations, which have been duly acknowledged. I also declare that it has not been previously or concurrently submitted for any other degree at UPM or other institutions.*

_____
**KHOR SIAK WANG**

*Date:*

# TABLE OF CONTENTS

# LIST OF TABLES

# LIST OF FIGURES

# LIST OF ABBREVIATIONS

| | |
|---|---|
| *2D* | *Two Dimensions* |
| *3D* | *Three Dimensions* |
| *ATM* | *Asynchronous Transmission Mode* |
| *CAD* | *Computer-aided Design* |
| *CBIR* | *Content-based Image Retrieval* |
| *CBVIR* | *Content-based Visual Information Retrieval* |
| *CC* | *Colour Code* |
| *CCV* | *Colour Coherence Vector* |
| *CD* | *Compact Disk* |
| *CIE* | *Commission Internationale de l'Êclairage* |
| *CMY* | *Cyan (C), Magenta (M), and Yellow (Y)* |
| *CRT* | *Cathode Ray Tube* |
| *DC* | *Dominant Colour* |
| *FE* | *Feature Extraction* |
| *GUI* | *Graphical User Interface* |
| *HIS* | *Hue-Intensity-Saturation* |
| *HSV* | *Hue-Saturation-Value* |
| *IR* | *Information Retrieval* |
| *ISDN* | *Integrated Services Digital Network* |
| *MIR* | *Multimedia Information Retrieval* |
| *MARS* | *Multimedia Analysis and Retrieval System* |
| *MPEG* | *Moving Picture Experts Group* |

QBIC                                  *Query By Image Content*

*QBIC*                    *Query By Image Content*

*RGB*                    *Red-Green-Blue*

*RF*                      *Relevance Feedback*

*SCH*                    *Spatial-Chromatic Histogram (SCH)*

*SI*                      *Statistical Information*

*SMAT*                *Sequenced Multi-Attribute Tree*

*SONET*             *Synchronous Optical Network*

*SQL*                    *Structured Query Language*

*UCID*                 *Uncompressed Colour Image Database*

*VLSI*                  *Very Large-Scale Integration*