A NEW METHODOLOGY FOR EVALUATION AND BENCHMARKING OF
SKIN DETECTOR BASED ON AI MODEL USING MULTI CRITERIA
ANALYSIS

QAHTAN MAJEED YAS

THESIS SUBMITTED IN FULFILMENT OF THE REQUIREMENT FOR
DOCTOR OF PHILOSOPHY (ARTIFICIAL INTELLIGENCE)

FACULTY OF ART, COMPUTING & CREATIVE INDUSTRY
UNIVERSITI PENDIDIKAN SULTAN IDRIS

2018

# ABSTRACT

This study aims to develop a new multi-criteria decision analysis methodology for skin detector evaluation and benchmarking based on artificial intelligence models. Two experiments were conducted. The first experiment comprised two stages: (1) Adaptation of the best previous case of skin detection approach utilizes multi-agent learning based on different color spaces. This stage aimed to create a decision matrix of various color spaces, and three groups of criteria (i.e., reliability, time complexity, and error rate within dataset) to test, evaluate and benchmark the adapted skin detection approaches. (2) Performance of multiple evaluation criteria for skin detection engines, this stage included two key stages. First, the correlation between criteria to investigate their relationship and determine their degree of correlation. Second, the performance analysis of criteria to identify the factors that affect the behavior of each criterion. The second experiment utilized a new multi-criteria decision-making by adopting the integration of TOPSIS and AHP to benchmark the results of skin detection approaches. In the validation process, multi-criteria measurement was used to calculate the trade-off for different criteria. Color spaces assessment were conducted to determine the best color spaces with adaptive skin detection engines. Moreover, mean and standard deviation values for thresholds were calculated to select the best color space. Two groups of findings were provided. First, the overall comparison of external and internal aggregation values in selecting the best color space, that is the norm RGB at the sixth threshold. Second, (1) the process proves that the distribution of color spaces with its threshold values affects the behavior of the criteria determined as a trade-off between the criteria according to their weight distribution. (2) The YIQ color space obtains the lowest value and is the worst case, whereas the norm RGB color space receives the highest value and is the most recommended. (3) The best result achieved at the threshold = 0.9. Thus, the implications of this study benefit individuals, research centers, and organizations interested in skin detection applications. Moreover, it provides benefits to software developers working in industrial companies and institutions in developing different techniques and algorithms with different applications.

# METODOLOGI BARU UNTUK PENILAIAN DAN PENYELESAIAN DETEKSI KULIT BERDASARKAN MODEL AI MENGGUNAKAN ANALISIS KRITERIA MULTI

## ABSTRAK

Kajian ini bertujuan untuk membangunkan metodologi baharu bagi menilai dan menanda aras pengesanan kulit berdasarkan model kecedasan buatan menggunakan analisis pelbagai kriteria. Untuk tujuan ini, dua eksperimen telah dijalankan. Eksperimen pertama terdiri daripada dua peringkat: (1) Adaptasi kes terbaik terdahulu dalam mengesan kulit menggunakan pendekatan multi-agen berdasarkan ruang warna yang berbeza. Peringkat ini bertujuan untuk membuat matriks keputusan pelbagai ruang warna dan tiga kumpulan kriteria (iaitu, kebolehpercayaan, kerumitan masa, dan kadar kesilapan dalam set data) untuk menilai dan menanda aras pendekatan pengesanan kulit yang telah disesuaikan. (2) Prestasi kriteria pelbagai penilaian bagi enjin pengesanan kulit, di mana peringkat ini melibatkan dua peringkat kekunci. Pertama, korelasi antara kriteria untuk menyiasat hubungan dan menentukan darjah korelasi. Kedua, analisis prestasi kriteria untuk mengenal pasti faktor kriteria yang mempengaruhi kelakuan setiap kriteria. Eksperimen kedua menggunakan pendekatan membuat-keputusan multi-kriteria baharu melalui integrasi antara TOPSIS dan AHP untuk menanda aras keputusan pendekatan pengesanan kulit. Di dalam proses pengesahan, pengukuran pelbagai kriteria digunakan untuk mengira keseimbangan bagi pelbagai kriteria. Penilaian ruang warna dijalankan untuk menentukan ruang warna yang terbaik dengan enjin pengesanan kulit yang telah diadaptasi. Seterusnya, nilai min dan sisihan piawai dikira untuk memilih ruang warna yang terbaik. Hasil dapatan daripada dua kumpulan adalah seperti berikut. Pertama, perbandingan keseluruhan nilai agregasi luaran dan dalaman dalam memilih ruang warna terbaik, iaitu RGB norma pada ambang keenam. Kedua, (1) proses membuktikan bahawa penagihan ruang warna dengan nilai ambangnya mempengaruhi kelakuan kriteria yang ditentukan sebagai keseimbangan antara kriteria berpandukan pengagihan berat masing-masing. (2) Ruang warna YIQ memperoleh nilai terendah dan merupakan kes terburuk, manakala ruang warna norm-RGB memperoleh nilai tertinggi dan paling disyorkan. (3) Dapatan terbaik dicapai pada ambang = 0.9. Oleh itu, implikasi kajian ini memberi manfaat kepada individu, pusat penyelidikan dan organisasi yang berminat dalam aplikasi pengesanan kulit. Kajian ini turut memberi manfaat kepada pembangun perisian yang bekerja di industri dan institusi dalam membangunkan teknik dan algoritma yang berbeza bagi aplikasi yang berbeza.

# TABLE OF CONTENTS

# LIST OF TABLES

# LIST OF FIGURES

# LIST OF ABBREVIATION

| | |
|---|---|
| ANN | Artificial Neural Network |
| AHP | Analytic Hierarchy Process |
| ANP | Analytic Network Process |
| CPU | Central Processing Unit |
| CIE | Commission International de L'Eclairage |
| CR | Consistency Ratio |
| DM | Decision Matrix |
| EM | Evaluation Matrix |
| FP | False Positive |
| FN | False Negative |
| GH | Grouping Histogram |
| GDM | Group Decision Making |
| HAW | Hierarchical Adaptive Weighting |
| IT | Information Technology |
| KNIME | Konstanz Information Miner |
| KEEL | Knowledge Extraction based on Evolutionary Learning |
| LUT | Lookup Table |
| MCDM | Multi- Criteria Decision Making |
| MADM | Multi- Attribute Decision Making |
| MCDA | Multi-Criteria Decision Analysis |

| MEW | Multiplicative Exponential Weighting |
| RI | Random Index |
| SVM | Support Vector Machine |
| SAN | Segment Adjacent-Nested |
| SAW | Simple Additive Weighting |
| TP | True Positive |
| TN | True Negative |
| TOPSIS | Technique for Order Preference by Similarity to Ideal Solution |
| WEKA | Waikato Environment for Knowledge Analysis |
| WSM | Weighted Sum Model |

# CHAPTER 1

# INTRODUCTION

## 1.1 Introduction

This chapter introduces the research direction, research background, and a statement of the problem. This chapter also presents the ambitions, motivations, and objectives of this research are also presented.

Section 1.2 presents a brief background of the research components. Section 1.3 introduces the statement of the problem, which is the basis of the research direction. Section 1.4 discusses the scope of the research. Section 1.5 describes the research objectives.  Section 1.6 presents a general view of the research. Finally, Section 1.7 briefly outlines the main structure of the research.

## 1.2 Research Background

Decades ago, the skin detection approach has been considered an important platform for various fields, such as medical and several scientific disciplines (L. Huang et al. 2015). In other words, skin detection has gained an important function in a wide range of image or video processes for various applications. A few factors that directly impact skin appearance include illumination, background, camera characteristics, and ethnicity (Kakumanu, Makrogiannis, and Bourbakis 2007). Elgammal, Muang, and Hu (2009) defined the skin detection approach as a process of finding skin-colored pixels and regions in an image or a video into a specific region. This process is typically used as a preprocessing step to finding regions in images that potentially detect the human face and limbs. The skin detection approach includes various applications, such as face detection, (Zhipeng, C., Junda, H., & Wenbin 2010), face tracking (Tsai 2012), gesture analysis (Hussain, I., Talukdar, A. K., & Sarma 2014), Internet pornographic image filtering (Lee, Kuo, and Chung 2010), surveillance systems (Zui Zhang 2009), content-based image retrieval systems (Patil, C. G., Kolte, M. T., Chatur, P. N., & Chaudhari 2014), and various human–computer interaction domains (Hollender et al. 2010). The most practical and effective techniques are used in developing skin detector artificial intelligence (AI) algorithms according to the literature on skin detection for skin pixel and non-skin pixel features based on color features. On the contrary, many researchers have applied hybrid algorithms in AI models (Singh Sisodia and Verma 2011; (Shruthi, M. L. J., & Harsha 2013; Zaidan et al. 2014b). However, with the current rapid development of the skin detection approach in various applications, finding an evaluation and benchmarking

methodology that is reliable, effective, and comprehensive has become critical (Jones and Rehg 1999; Phung, Bouzerdoum, and Chai, D. 2005; Gamage, Akmeliawati, and Chow 2009; Taqa and Jalab 2010a).

Considering the basic criteria evaluation of reliability, time complexity, and error rate within the dataset in the design of any skin detector application, (Jones and Rehg (1999) adapted three criteria, namely reliability, computational cost, and error rate of skin detection. In one of the earliest works that highlight the problem of skin detection evaluation and benchmarking, three general requirements for the skin detection approach are reported: adapted reliability (i.e., the obtained skin detection rate and false positives) and datasets (i.e., the obtained equal error rate comparison of AI models) with less time-consuming requirements to process web images.

Despite the importance of the remaining criteria, Phung, Bouzerdoum, and Chai, D. (2005) highlighted the dataset criterion by comparing two algorithms. The dataset is represented by training and testing for skin and non-skin pixels for skin segmentation images. However, the output images created through a classifier are compared pixel-wise with the ground truth of skin segmentation. Gamage, Akmeliawati, and Chow (2009) reported a skin detection algorithm that has been tested with images through independent databases. They investigated the size of the image, which has a significant impact on time complexity. Thus, they proved that increasing image size leads to low accuracy than increase time complexity of the experiment. Finally, Taqa and Jalab (2010a) stated that reliability is a prerequisite for

skin detection evaluation. They highlighted a reliability criterion based on accuracy, precision, and recall of the image color despite the importance of the remaining criteria. However, the quality assessment of skin detection requires attention.

Consequently, two key problems are encountered by skin detection developers. One is the evaluation of skin detection approaches based on the abovementioned evaluation criteria and benchmark new skin detection approach versus existing approaches. Therefore, the evaluation and benchmarking process need to consider these requirements. Despite the tradeoff among various criteria, (Jones and Rehg (1999); Phung, Bouzerdoum, and Chai, D. (2005); Gamage, Akmeliawati, and Chow (2009); Taqa and Jalab (2010a) have adopted each of the proposed criteria. They attempted to evaluate the reliability criterion for a given time complexity based on different datasets. However, the term "reliability" is unclearly defined in the literature. According to the preceding studies mentioned, the percentage of reliability varies depending on different adapted algorithms and thus exhibit an inconsistent level. Meanwhile, Fernandes, Cavalcanti, and Ren (2013) reported time complexity variation between the algorithms, which depend on the CPU time. Consequently, the processing time of an image is affected, but this aspect is excluded in the scope of the present research. Therefore, the calculation should be the highest percentage of reliability compared with the lowest time complexity of the output image. Kawulok (2013) mentioned that the dataset can be divided into two classes, namely training and validation data, to find the minimum detection error. In general, all these studies have

proven the evaluation and benchmarking process of each of these criteria based on independent guidelines.

Therefore, conducting further investigations and developing a clear methodology for testing, evaluation, and benchmarking are necessary to standardize basic and advanced requirements for the skin detection approach. Redefining the problem of evaluation and benchmarking need is also necessary. Moreover, the new evaluation methodology must be flexible to handle the conflicting criteria problem and must have the capability to maintain the current criteria.

## 1.3 Significance of Study

The evaluation and benchmarking of skin detection approaches are important areas for many researchers and organizations interested in their applications. Many individuals and organizations are interested in the applications of skin detection approaches, such as researchers working in scientific research centers, developers working in industrial companies and institutions, and graduate students enrolled in schools that develop various applications of skin detection approaches. Thus, the importance of the study is the development of multiple applications of skin detection approaches, including face detection, face tracking, gesture analysis, Internet pornographic image filtering, surveillance systems, content-based image retrieval systems, and various human-computer interaction domains. Moreover, this study