# Automatic Estimation of Malaria Parasitemia Based on Microscopy Images

by

**Budi Sunarko**

*Thesis*
*Submitted to Flinders University*
*for the degree of*

## Doctor of Philosophy

College of Science and Engineering
19th December 2017

# Table of Contents

# List of Figures

# List of Tables

# Summary

Malaria is a blood-borne disease and a major cause of mortality in many parts of the world. Malaria is curable if diagnosed in time. Malaria can be diagnosed by microscopy analysis. Automatic malaria diagnosis has the potential to play an important role in reducing mortality due to malaria. The term parasitemia is used to reflect the severity of the malaria disease. This thesis aims to develop computer methods for automatic malaria diagnosis and parasitemia estimation.

According to the literature on biomedical image processing and an understanding of blood and malaria, colour intensity is a strong indicator for malaria parasite identification. In addition, natural characteristics of malaria, such as size, shape and appearance of parasites, and blood components are also considered to play an important role in malaria identification and parasitemia estimation. However, the role of these natural characteristics in automatic malaria diagnosis is unknown and the contribution of automatic estimation of malaria parasiatemia has received relatively little attention.

The focus of this thesis is to include the natural characteristics of blood components and malaria parasites in automatic malaria diagnosis and parasitemia estimation based on microscopy images. Parasitemia can be described as the number of malaria parasites in one microlitre (µl) of blood fluid. The location and the colour of malaria parasites in blood are known to characterize infected erythrocytes in visual images. In addition, the maximum parasitemia suffered in the human body has been clinically determined. Biologically, the composition and size of blood components have been also recognized. The role of these natural characteristics is to transform the knowledge of microscopy malaria diagnosis to automatic malaria diagnosis and parasitemia estimation.

A number of experiments were conducted to determine the best parasitemia estimation based on erythrocyte classification in thin blood film images and parasite classification in thick blood film images. This study made use of the natural characteristics of blood components and malaria parasites to set adaptive thresholds for segmenting leukocytes and parasites. Colour intensity and location features of leukocytes, erythrocytes, and parasites were extracted and measured. Discriminant analysis was used to classify the leukocyte footprints as leukocytes or phagocytes, the parasite footprints as parasites or non-parasites, and identify erythrocyte footprints as normal or infected. Classification performance was evaluated by informal readers using confusion matrices. Subsequently, the

percentage of infected erythrocytes in thin blood films and the number of parasites and leukocytes in thick blood films were converted to estimate parasitemia scores. Parasitemia estimation was validated by parasitemia scores from expert readers.

Results indicate that the parasitemia estimation fitted well with parasitemia scores from expert readers both in thin blood films (r = 0.97, p-value = 0.54) and thick blood films (r = 0.79, p-value = 0.40), at $\alpha = 0.05$ level. In thin blood films, the performance of erythrocyte classification based on the combination of the colour intensity features of parasites and location features perfomed better than that based on the only grayscale intensity features of erythrocytes. Meanwhile, morphological features may not be optimal for an automatic parasitemia estimation in thick blood films.

In addition, a number of discoveries were made in the course of the study. The combination of the natural characteristics of blood components and malaria parasites is an essential feature to set an adaptive threshold for segmenting parasites in thin blood film images. Based on the fact that thrombocytes are naturally located outside of erythrocytes, this location feature is another essential feature to distinguish infected erythrocytes from thrombocytes in thin blood film images. In other words, better erythrocyte identification and parasitemia estimation were obtained by involving the natural characteristics of blood components and malaria in the parasite segmentation and erythrocyte identification in thin blood film images. In thick blood film analysis, the presence of leukocytes is vital for estimating parasitemia scores with leukocytes generally having the highest intensity in inverse thick blood film images. Accordingly, leukocyte intensity may be well utilized as a reference in setting adaptive thresholds for leukocyte and parasite segmentation.

These discoveries have potential contributions to the fields of automatic malaria diagnosis and parasitemia estimation based on both thin and thick blood film images, and so form natural seeds for future work.

# Publications arising from the study

xviii

## Referred Conference Paper

1. Sunarko, B., Williams, S., Prescott, W. R., Bottema, M. J., & Byker, S. M., 2016, 'Correlation between Automatic Detection of Malaria on Thin Film and Experts' Parasitemia Scores', Engineering International Conference 2016, AIP Conference Proceedings 1818, pages 020054-1-020054-10.
2. Sunarko, B., Williams, S., Prescott, W. R., Bottema, M. J., & Byker, S. M., 2013, 'Comparative Study of Two Methods for Blood Cell Segmentation', Engineering International Conference 2013 Proceedings, pages II-96-II-100.

## Conference Abstract

3. Sunarko, B., Williams, S., Prescott, W. R., Bottema, M. J., & Byker, S. M., 2014, 'Red Blood Cell Classification for Automatic Malaria Diagnosis', Australian and New Zealand Industrial and Applied Mathematics 2014, ANZIAM Conference Book, page 92.

# Declaration

I certify that this thesis does not incorporate without acknowledgment any material previously submitted for a degree or diploma in any university; and that to the best of my knowledge and belief it does not contain any material previously published or written by another person except where due reference is made in the text.

Budi Sunarko, Candidate

Murk J. Bottema, Principal Supervisor

Simon Williams, Co-supervisor

# Acknowledgments

I would like to thank my supervisor, A/Prof. Murk J. Bottema for his support and encouragement throughout my PhD study. Importantly, he helped me improve my research skills in medical image analysis and academic writing for my project. Specifically, I really thank him for giving me the opportunity to study in Flinders University.

I also would like to thank my co-supervisor, Dr. Simon Williams for his supervision during my PhD study. Specifically, he helped me in using MATLAB for programming research projects and LaText for writing papers.

I thank Dr. Gobert Lee, Dr. Mariuz Bajger, and Dr. Darfiana Nur for providing literatures and advices to this projects.

Special thanks to thank Hydas World Health for providing thin and thick blood films and Microscopes International for scanning the films to get blood film images for my study.

I also would like to thank expert readers for estimating parasitemia scores of the blood films and informal readers for identifying malaria parasites and blood components in the blood film images.

Gratefully, I also wishe to acknowledge the support of Australian Government which provided the Australian Awards Scholarship for my whole study.

Finally, I thank all my family for their great support and love all the way through.

# Chapter 1:     Introduction

## 1.1  Overview

Malaria is a blood-borne disease that can be diagnosed by microscopy analysis. Malaria is caused by a parasite. For ease of exposition, the term "malaria parasite" or "parasite" will be used to refer to the parasite. Manual malaria diagnostic methods have potential problems due to the limitation of the human visual system in consistently inspecting large numbers of samples, especially in pandemic areas. Instead, this thesis tries to capture the knowledge of microscopy malaria diagnosis including the characteristics of blood components and malaria parasites, in a computer-aided malaria diagnosis system. An automatic system has the potential to estimate malaria parasitemia consistently and reliably.

This introductory chapter is intended to provide context for this thesis. The disease of malaria is reviewed in Section 1.2 and the role and impact of malaria in an example country, Indonesia, is discussed in Section 1.3. The life cycle of malaria in the human body is reviewed in Section 1.4 and the components of blood fluid and their roles are reviewed in Section 1.5. Clinical malaria diagnosis in current practice is presented in Section 1.6, and the characteristics of malaria parasites and parasitemia estimation are discussed in Section 1.7 and 1.8, respectively. The characteristics of colour images and the appearance of the components of blood in thin and thick blood films are presented in Section 1.9. This discussion leads to the motivation and objectives of the thesis in Section 1.10. This is followed by an overview of the data sets used in this thesis in Section 1.11 and an overview of the structure of the remainder of this thesis in Section 1.12.

## 1.2  Malaria

Malaria is a major cause of mortality and is a serious public health burden in many parts of the world. It was reported that around 40% of people in the world are prone to malaria (Suradkar 2013) and that malaria is particularly prevalent in lower-income countries subject to malaria endemics (Sachs & Malaney 2002). Furthermore, in April 2016, the World Health Organisation (WHO) updated a recent fact sheet reporting that the number of people, at risk of malaria was 3.2 billion - almost half of the world's population, covering 95 countries and territories (WHO 2016).

Malaria is caused by a parasite, called *plasmodium* (*P*), which infects human red blood cells (erythrocytes). There are four different malaria plasmodia: *Plasmodium*

*falciparum*, *Plasmodium vivax*, *Plasmodium oval*, and *Plasmodium malariae*. Among the four plasmodia, *P. falciparum* and *P. vivax* are the most prevalent and dominant malaria parasites. These plasmodia are responsible for most malaria infections and have the greatest potential to infect the human body. Marten et al. estimated that the number of malaria-infected people will increase by 300 and 150 million for *P. falciparum* and *P. vivax*, respectively, by 2080 (Martens, P et al. 1999).

The malaria parasite is spread from a malaria sufferer to another healthy person by a female Anopheles mosquito, the primary vector for malaria. The Anopheles mosquito sucks in the infected erythrocytes. The parasite then multiplies in the mosquito's stomach. The malaria parasite is transmitted when an infected Anopheles mosquito bites another healthy person.

The transmission of malaria depends on the number of live mosquitoes, which is in turn are affected by climatic factors including precipitation, temperature and humidity. In many tropical areas, malaria is seasonal and peaks during or immediately following monsoon. It may become epidemic when occurring in locations inhabited by people with low or no immunity against the disease, or when a large number of such people, seeking work or refuge, enter a site intensively infected by malaria (Jagessar & Rampersaud 2014).

Similar to other mosquitos, Anopheles mosquitos need water to breed and so they are mostly found in wet areas, such as swamps and ponds. In addition, Anopheles mosquitos grow well in tropical climate regions, but they also can be found in temperate regions. By mathematical modelling, Marten et al. assessed the correlation of malaria to a range of factors, including climate change and geographical distribution of malaria mosquitoes (Martens, W et al. 1995). In the context of climate change, results from simulations indicated that epidemic mosquitos, including Anopheles, have the potential to increase two-fold in tropical regions and more than 100-fold in temperate regions due to a several-degree global mean temperature increase by the year 2100. This effect of increasing global temperature presents a real risk of reintroduction of malaria mosquito transmission into malaria free regions, such as parts of Australia and the United States.

Clinically, many diseases produce similar symptoms. Typically, malaria exhibits flu-like symptoms (Karcheva, Atanasova & Rainova 2017; Maichomo et al. 1999) the incubation period, between eight and 14 days (Baird et al. 1995). However, such period may differ according to the parasite's species. When no suitable medication is available or the

parasite has become resistant to available medication, the infection may develop fast and endanger the subject's life. Without proper treatment, this can lead to death. Malaria infects and eradicate erythrocytes, prompting anaemia, while blocking capillaries and thus preventing blood from reaching the brain (Aikawa et al. 1990).

Due to the natural resilience of the human body to malaria, individuals with less than 10% of infected erythrocytes are likely to survive. Those with parasitemia greater than 10% may need transfusion and have a high mortality rate (Garcia & Bruckner 1997). In other words, the human body can survive malaria only if less than 10% of erythrocytes are infected.

Although the number of people at risk of malaria is still high, malaria is a preventable and curable disease (WHO 2016). Between 2001 and 2015, it was estimated that the number of malaria cases around the world decreased by 18% from 262 million to 214 million and the number of malaria deaths decreased by 48% from 839 thousand to 438 thousand (WHO 2014b). *P. falciparum* malaria was effectively cured through the use of medicines Quinine and Fansidar (Hall et al. 1975). In the frame of research, cysteine proteinase was effective against *Plasmodium vinckei* in murine malaria (Rosenthal, Lee & Smith 1993). This means that medicines have the ability to cure people infected with malaria. Thus, diagnosis enables effective intervention in the successful treatment of malaria.

## 1.3   Malaria Cases in Indonesia

Indonesia has a long history of epidemics caused by malaria. It was estimated that malaria epidemics have appeared throughout the Indonesian archipelago since the first humans dwelled there (Elyazar, Hay & Baird 2011). Furthermore, in the same study, Elyazar et al. reviewed malaria cases and reported that Indonesia has been hit by a malaria outbreak every year. For example, malaria outbreaks occurred in 1999 in five provinces covering six districts with a case fatality rate of 1.12% (2,407 cases and 27 deaths) (Marwoto & Sulaksono 2003).

### 1.3.1   Factors of Malaria Outbreaks

Several factors have contributed to malaria outbreaks in Indonesia. Firstly, the unique human factor is transmigration, a special Indonesian government program. The program promotes the movement of populations from crowded islands, namely Java and Bali, to other less crowded islands, such as Kalimantan, Papua, Sumatra, and Sulawesi. The hometowns of

transmigrants are generally malaria free or have low malaria cases. However, the destinations of transmigrants are hyper malaria endemic areas (Figure 1.3). As a result, the more non-immune malaria transmigrants, the more new malaria sufferers. For instance, in 1988 and 1992, between two and six months after new transmigrants arrived in Arso, Papua, local malaria epidemics occurred (Baird et al. 1995). As many as 14% of the new transmigrants were exposed to malaria within 14 days of arrival. The prevalence reached over 70%, and 10% of these were severely parasitemic and had to be hospitalized.

Secondly, El Niño, a climate change effect, also participates in rising malaria outbreaks. El Niño causes a chain effect in a community. El Niño-affected regions receive reduced rainfall resulting in water and food shortages, especially in hilly areas. These water and food shortages forced people living in hilly regions, with high parasitemia, to migrate to lowland or coastal areas where water and food sources were relatively available. These human movements caused a significant increase in malaria transmission (Bangs & Subianto 1999).

### 1.3.2 Distribution of Malaria Disease

Malaria cases have been found in almost all provinces of Indonesia (Figure 1.1), but the distribution in each province is not uniform (Figure 1.2) and average stratification between provinces is different (Figure 1.3). The stratification is determined by the *annual parasite incidence* (API), an indicator of malaria cases (the number of malaria sufferers per 1000 local people). High numbers of malaria cases occurred in eastern Indonesia, especially West Papua and East Nusa Tenggara. In general, districts with high API are hilly and forested regions in remote areas. However, coastal and lowland areas are not free from malaria epidemic even though they have lower API (Bangs & Subianto 1999).

Figure 1.1: Map of Indonesia consisting 33 provinces in 2009. Modified and reprinted from www.indonesiamatters.com/images/indonesia-map.gif.



| 0 | 0-1 | 1-5 | 5-49 | 50 – 100 | > 100 |
|---|---|---|---|---|---|
| Free | Low | Moderate | | High | |

Figure 1.2: Distribution of malaria outbreaks throughout Indonesia in 2009. Reprinted from Buletin Jendela Data dan Informasi Kesehatan.

Figure 1.3: Malaria stratification in Indonesian provinces. Reprinted from Buletin Jendela Data dan Informasi Kesehatan.

### 1.3.3 Economic Burden of Malaria Disease

Malaria is estimated to costs Indonesia hundreds of millions (in US dollar) each year. The cost covers malaria prevention and treatment as well as income loss due to hospitalisation and mortality.

The Indonesian ministry of health (Kemenkes) reported that in the year 2012 there were 1,237,389 malaria cases and 9,899 deaths (case fatality rate 0.8%) in Indonesia (Kemenkes 2015). The calculated economic loss due to malaria cases in 2012 was around US$ 505 million (Kemenkes 2015). Seventy-five per cent of the deaths caused by malaria were productive people. The income loss due to the mortality and hospitalisation were 84.78% and 7.26% of the total malaria cost, respectively. Another load was the cost of malaria treatment, around 7.96% of the total cost.

### 1.3.4 Obstacles for Preparing Blood Films in Remote Areas

Microscopy diagnosis of malaria requires blood films. High quality blood films must be prepared from fresh peripheral blood taken directly from malaria-suspected patients. However, there are problems in preparing blood films from people suffering malaria in remote areas, far away from a malaria laboratory (more than a single day's travel). This is because of the nature of human blood fluid and microorganisms. Without special tools, it is difficult to carry the human blood samples for a long time in good condition because the human blood fluid easily clots. Clotted blood fluid cannot be made into blood films properly in a laboratory. One method to overcome this is to add an anticoagulant to the blood samples,

6

but adding anticoagulant to blood fluid may alter cellular morphology and result in staining characteristics (Bain 2014). As a result, erythrocytes and leukocytes are difficult to recognize and count. Furthermore, the United Kingdom National External Quality Assessment Service (UK NEQAS) Parasitology reported that in some cases gametocytes and schizonts forms (host erythrocytes and the contained parasites) may be destroyed when exposed to anticoagulant (EDTA) for several hours (NEQAS 2016).

Conversely, if blood films are prepared in the field, the blood films are prone to contamination by fungus and other microorganisms because the blood films are good environments for such growth. Consequently, malaria diagnosis can be destroyed by these organisms appearing on blood films.

## 1.4 The Malaria Lifecycle in the Human Body

In the human body, malaria parasites develop in two main stages: the liver stage and the blood stage. Malaria parasites are also able to dwell in the brain, in severe cases, and white blood cells (leukocytes) when parasites have been phagocytised by neutrophils (Section 1.7). The malaria parasite life cycle in the human body is shown in Figure 1.4. For malaria diagnosis based on microscopy, the blood sample is taken from peripheral blood (blood streaming under skin) and so the blood stage is of primary interest in this thesis.



Figure 1.4 Malaria life cycle in the human body, modified and reprinted from the Centers for Disease Control and Prevention (CDC 2016).

During the blood stage, malaria parasites undergo various form changes. From the liver stage, merozoites infect erythrocytes and start growing in the blood stage. In the

beginning of the blood stage, a merozoite develops to form an immature trophozoite (ring form). The diameter of the ring form of *P. falciparum* is between 1–2 µm (Murray, Rosenthal & Pfaller 2015; Palakuru 2016). The immature thropozoite, then progresses to a mature thropozoite or a gametocyte. Subsequently, the mature trophozoite starts to regenerate asexually and progress to a schizont form. The schizont form consists of a group of merozoites and has a diameter of between 7–8 µm. Finally, the infected erythrocyte bursts and releases new merozoites into the bloodstream. At this step, new merozoites are independent (outside of erythrocytes) before infecting other erythrocytes. Meanwhile, the gametocyte develops and regenerates sexually. The size of gametocyte is between 7–14 µm.

## 1.5   Components of Blood Fluid

Blood fluid consists of three main components: erythrocytes, leukocytes, and platelets (thrombocytes) (Ruberto, C. et al. 2000). Erythrocytes are red blood cells. Their main function is to carry oxygen from the lungs to tissue and carbon dioxide from tissue back to the lungs. Leukocytes are white blood cells responsible for counteracting foreign material including diseases. Thrombocytes, also known as platelets, play a role in clotting. The number of erythrocytes and leukocytes per microlitre of blood is normally about 5,000,000 cells and 8,000 cells, respectively (Moody 2002; Vander, Luciano & Sherman 2001; WHO 2010). Meanwhile, the number of thrombocytes is normally is around 250,000 (Vander, Luciano & Sherman 2001). This means that the number of thrombocytes and leukocytes are much fewer than that of erythrocytes. The ratio of the number of thrombocytes and leukocytes over that of erythrocytes is around 5% and 0.16%, respectively.

Morphologically, the average diameters of normal erythrocytes and thrombocytes are approximately 7-8µm (Löffler & Rastetter 2000; WHO 2010) and 2µm (Freitas 1999), respectively. Thus, thrombocytes are smaller than erythrocytes with diameters around 25% the diameters of erythrocytes. Leukocytes are larger than erythrocytes and are divided into many classes: neutrophils, eosinophils, basophils, monocytes, and lymphocytes, which fall within a size range of 8-12 µm. The most dominant leukocyte class is neutrophils, approximately 50-70% of leukocytes with a size range of 8-10µm, followed by lymphocytes (20-40%), monocytes (2-8%), eosinophil (1-4%), and basophils (0.1%) with size ranges of 6-12 µm, 15-30 µm, 10-12 µm, and 12-15 µm (Freitas 1999; Wheater, Burkitt & Daniels 1979).

Figure 1.5 Erythrocyte views. Modified and reprinted from Biomed 108-Human Physiology.



(a)

(b)

(c)

(d)

Figure 1.6: Original images from thin blood films (Section 1.11). Th is a thrombocyte. (a) An image from a positive slide containing infected erythrocytes (Ip). (b) An image from a normal slide. (c) An image containing a leukocyte (Lu). (d) An image containing a phagocyte (Pg) as described in Section 6.4.

Figure 1.7: Original images from thick blood films (Section 1.11). Rm is residual erythrocyte membrane. Lu and Pg are as in Figure 1.6. (a) An example of image, from a positive slide, containing leukocytes and parasites (Pi). (b) An example of image containing leukocytes and phagocytes.

Erythrocytes are relatively regular in shape. From the top, erythrocytes are circular in shape, but from the side they are biconcave discs (Figure 1.5). Unlike erythrocytes, thrombocytes are irregular shaped bodies, without a nucleus but with fine red granules on blueish background (WHO 2010). Like thrombocytes, leukocytes are relatively irregular in shape, but each leukocyte has a nucleus surrounded by cytoplasm; sometimes, the cytoplasm is granular. Some leukocytes have a multi-lobe nucleus, which forms an irregular shape and uncertain volume or area. Some erythrocytes, leukocytes, and thrombocytes are shown in Figure 1.6 and Figure 1.7.

## 1.6 Methods for Clinical Diagnosis

Currently, the most common malaria diagnostic technique relies on observing blood films of the subject under a microscope and then identifying the parasite. This microscopy based malaria diagnosis is the gold standard (WHO 2010). The method needs trained microscopists. To analyse malaria, generally, microscopists examine blood films after a preparation and staining process (WHO 2010) and view the blood films under a microscope in a laboratory. Microscopists find an appropriate area manually by moving the stage up and down and left and right and adjusting the fine focus and coarse focus and noting the number of parasites to determine parasitemia. Details of the process for determining parasitemia will be described in Section 1.8.

Microscopy slides usually include two blood films from the same subject, a thin blood film and a thick blood film. These two blood films have different characteristics and different uses. Therefore, these are recommended for malaria diagnosis. A thick blood film is usually used to determine if there is an infection and to estimate parasitemia, especially if the parasitemia is low. This is because a thick blood film contains a larger amount of blood

per unit area of slide. A thin blood film is used to determine species and to estimate parasitemia if the parasitemia is high (Garcia & Bruckner 1997).

An accurate and efficient diagnosis is possible in ideal conditions and by well-trained staff; however, in developing countries, especially in rural areas, the method is often inaccurate and inefficient. Potential weaknesses of the technique arise from (a) inaccuracy from the limitation of the human visual system in identifying huge samples due to high work load, (b) time consuming for inexperienced microscopists who must frequently refer to the references, (c) considerable cost for training microscopists (Oaks Jr et al. 1991; Ross et al. 2006). Incorrect observation due to the lack of technique may potentially result in misdiagnosis or delayed diagnosis, possibly leading to a more severe disease state or death.

Other methods for malaria diagnosis are rapid diagnostic tests (RDTs) based on antigen and quantitative real-time polymerase chain reaction (qPCR) diagnostics based on nucleic acid amplification. Although these tests use new technology, overall they are not yet as reliable as microscopy. For example, qPCR yielded misdiagnosis in cases of low parasitemia (Dakić et al. 2014). Accordingly, the role of these tests is as a method complementary to microscopy tests and must be confirmed by the microscopy method (Moody 2002). Because they are portable, easy, quick, and cost-effective to perform, RDTs are valuable and practical to use in remote areas, where microscopy tests are not available, or in emergency conditions, such as the early detection of malaria outbreaks (Mboera et al. 2013).

However, the RDTs do not replace the use of malaria microscopic examination and must be confirmed by microscopy (CDC 2016) for the following reasons. Firstly, the assay can only detect if the subject is infected or not, but is not able to determine the percentage of infected erythrocytes, which is an important prognostic indicator. Secondly, after being stored a year at room temperature, the performance of RDTs is relatively poor (Mboera et al. 2013). Thirdly, the currently approved RDT can only detect two different malaria antigens: one is specific for *P. falciparum* and the other is found in all four malaria species. In addition, the sensitivity of RDTs was low in diagnosis of *P. ovale* (Tanizaki et al. 2014). Thus, microscopy is needed to confirm all negative RDTs and determine the malaria species of positive RDTs.

## 1.7 Recognising the Malaria Parasite in Blood

With proper Giemsa staining, it is possible to distinguish an infected erythrocyte from normal erythrocytes (Figure 1.8(a)) (WHO 2010). Each infected erythrocyte consists of a host erythrocyte, parasite(s), and vacuole. Each component of an infected erythrocyte responds to the stain differently. At the ring form, a parasite comprises one or two red dot chromatins, blue cytoplasm, and a clear vacuole (Figure 1.8(b)) (Dluzewski et al. 1992). The presentation of cytoplasm depends on the stage of development. In early stages, cytoplasm is small or even absent. At the developed stage, the pigment –a granular by-product of parasite growth– appears in colour from golden-brown to dark-brown or even black inside the parasite and with stippling, pink dots, which are the effect of parasite, on the host erythrocytes.

In the schizonts form, an infected erythrocyte may contain full parasites (merozoites). Generally, the size of infected erythrocytes in the schizont form are bigger than normal erythrocytes (WHO 2010). Malaria parasites also appear in leukocytes in some cases. Neutrophils, a kind of leukocyte, may contain malaria pigment, which is a by-product of parasite metabolism and is all that remains of parasites that have been phagocytosed (engulfed or eaten) by neutrophils.



(a)

(b)

Figure 1.8: Zoomed staining erythrocytes. (a) left cell is a normal erythrocyte and righ cell is an infected erythrocyte, zoomed and printed from an image of thin blood films (Section 1.11), (b) infected erythrocyte in detail, reprinted from WHO (WHO 2010).

## 1.8 Parasitemia

Parasitemia is the degree of malaria parasitic infection. There are several systematic measurements to assess parasite load in the blood based on thin and thick blood films. In thin blood films, parasitemia may be presented in two ways: the percentage system and parasitemia scores. The percentage system is calculated by noting the number of infected erythrocytes for every 100 erythrocytes (Equation (1.1)) and the parasitemia score is found by multiplying the percentages of infected erythrocytes by the standard count of erythrocytes

in one µl of blood fluid 5,000,000 erythrocytes/µl (Equation (1.2)) (Moody 2002; WHO 2010). For ease of exposition, the term "standard correlation" (Equation 1.2) will be used to refer to the conversion of the percentage of infected erythrocytes to parasitemia scores.

$$\text{Parasitemia (\%)} = \frac{\text{number of infected erythrocytes}}{\text{total number of erythrocytes}} \text{ x } 100 \qquad (1.1)$$

$$\text{Parasitemia score} = \frac{\text{parasitemia (\%)}}{100} \text{ x } 5\ 000\ 000 \qquad (1.2)$$

In thick blood film, parasitemia may be presented in two ways: the "plus system" and the parasitemia score. The plus system is simple but far less accurate for establishing parasite density. The plus system can be explained as below:

+       = 1-10 parasites per 100 oil-immersion thick blood film fields
++     = 11-100 parasites per 100 oil-immersion thick blood film fields
+++    = 1-10 parasites per single oil-immersion thick blood film field
++++   = more than 10 parasites per single oil-immersion thick blood film field

The plus system may lead to confusion because "many workers forget the finer details of the system and mix up the code (the number of plus signs) and the count (the number of parasites per field or per 100 fields)" (WHO 2010), this leads to confusion and provides unreliable information on parasite density. In conclusion, the "plus system" is unreliable and as a consequence, the system is not recommended (WHO 2010, 2014a). The WHO encourages microscopists to use parasite quantitation instead. The parasitemia score is calculated by summing the number of parasites per number of leukocytes, multiplied by 8000, a standard count of leukocytes/µl of blood fluid (Frean 2009; WHO 2010).

$$\text{Parasitemia score} = \frac{\text{number of parasites}}{\text{number of leukocytes}} \text{ x } 8000 \qquad (1.3)$$

## 1.9 Blood Film Images

In normal practice, blood films are viewed under a microscope and decisions regarding the presence of malaria are made immediately. The slides (blood films) may be stored for future

reference, but no images are taken of the blood films. For automatic systems, images of the blood films must be acquired. The images must be of sufficient resolution for the image analysis methods to distinguish between erythrocytes, parasites, thrombocytes, leukocytes and artefacts. On the other hand, a sufficiently large portion of the blood film has to be imaged in order to examine sufficient amounts of blood to reach statistically sound decisions. Together, these two requirements dictate that hundreds of images are needed per blood film.

## 1.9.1   Image Acquisition

Since hundreds of images are needed, manual selection of imaging fields is not practical. Typically, a dedicated microscope fitted with a suitable camera scans a fixed rectangular shaped region of the blood film. Since the blood film is not entirely uniform, some images contain no information and some contain a prohibitively dense accumulation of stained material. Thus, not all images are useful for diagnosis and must be discarded. Another possibility is that fields are chosen automatically using additional artificial intelligence steps.

Modern microscopes have a fixed illuminator to produce colour blood images; however, in some cases, the movement of the controller in shifting the glass is so fast that the camera cannot capture images properly. As a result, the images might be bluer or have different illumination due to differences in angle.

## 1.9.2   Thin Blood Film Images

Images of thin blood films from malaria-free subjects consist of background, erythrocytes, thrombocytes and possibly leukocytes. Thrombocytes are clearly located outside of erythrocytes (Figure 1.6(a), (b), (c), and (d)). Meanwhile, images of thin blood films from malaria infected subjects also contain background, erythrocytes, and thrombocytes, with the difference that some erythrocytes contain parasites (Figure 1.6(a) and (d)).

A well-prepared thin blood film comprises a single layer of erythrocytes. Erythrocytes are clearly visible and dominant in the sense that the total area of the image (number pixels) associated with erythrocytes is greater than the area of the image of all other components, except possibly the background. In general, thrombocytes occupy the smallest percentage of the total area. This is in line with the fact that the number of thrombocytes is lower than the number of erythrocytes and they are smaller in size. Leukocytes are much less frequent than either erythrocytes or thrombocytes and are not necessarily present in every thin blood film image.

Generally, leukocytes have the highest intensity value, followed by malaria parasites (if present). Even though the size of a malaria parasite is much smaller than that of an erythrocyte, the intensity profile of malaria parasites is much greater than that of erythrocytes (Figure 4.8(c)). If a parasite lies within an erythrocyte (as happens most frequently), the intensity profile of the parasite is clearly visible over the profile of the erythrocyte. Commonly, thrombocytes have the lowest intensity histogram above background intensity and erythrocytes have a significantly different intensity from the background. In thin blood film images, thrombocytes may present on top of an erythrocyte and may appear as if it lies within the erythrocyte.

### 1.9.3    Thick Blood Film Images

A thick blood film has a thickness of many erythrocyte diameters and so an image of a thick blood film represents a much larger volume of blood than an image of the same resolution of a thin blood film. However, in thick blood films, erythrocytes have been dissolved and so individual erythrocytes are not visible in thick blood film images. The image contains residual erythrocyte membranes, thrombocytes, leukocytes, possibly artefacts and parasites if the subject is malaria positive. Because the volume of blood is larger, many more parasites can be seen in a thick blood film image compared to a thin blood film image of the same resolution from the same subject.

Similar to thin blood film images, leukocytes generally have the highest intensity profile, followed by malaria parasites (if present) and background. However, thick blood film images are significantly different from thin blood film images. Significantly differences in thick blood film images are that erythrocytes disappear and residual erythrocyte membranes resulted from haemolysis, the process of releasing cytoplasm of erythrocytes into surrounding blood fluid, present in thick blood film images (Figure 1.7(a) and (b)). Furthermore, the background is completely absent in some cases.

An important difference between a thick blood film image and a thick blood film viewed under a microscope is that thick blood film images are 2-dimensional. This means that the images do not allow the user to change focus as is done in practice for direct microscopy (Section 1.3). Therefore, some objects are difficult to interpret because they are out of focus or overlap with other objects.

## 1.10 Motivation and Objective of the Thesis

From Section 1.2 and Section 1.6, the reduction of mortality and morbidity due to malaria disease and current malaria diagnosis barriers is vital. Accordingly, consistent accuracy of malaria diagnosis is crucial to avoid fatalities. Reliable and practical malaria diagnosis tools are important in malaria endemic countries, especially in remote areas.

In malaria-infected human blood, malaria parasites infecting erythrocytes are also present in the images (Section 1.7). These parasites have different colour intensity and size characteristics from the main components of blood images. Based on morphology and colour intensity, it is possible to identify the parasites. Computer vision has the potential to diagnose malaria consistently and accurately; however, automatic diagnosis of malaria has received relatively little attention in the literature.

The purpose of this study is to develop an automatic image analysis method for consistent and accurate malaria diagnosis and estimating parasitemia. The long-range goal is to produce a device that will be suitable for use in remote areas by persons having minimal training so that the barriers of malaria diagnosis in those areas, such as eastern Indonesia (Section 1.3), can be overcome. Such a device must be affordable for general use in developing countries. In this thesis, image analysis methods are developed that could be used in a reliable and automatic device.

The contribution of this thesis is to provide a method for estimating parasitemia and to test the estimates on cases spanning a wide range of known levels of parasitemia. As described above, other studies have focussed on counting parasites or on estimating parasitemia without comparison to ground truth.

## 1.11 Overview of Data Sets

Slides with thin and thick blood films were provided by Hydas World Health (Pennsylvania, USA). The anonymized slides were made from blood samples collected under an Institutional Review Board sanctioned minimal-risk human-use protocol, treated with Giemsa stain and cover slips attached with mounting glue. The slides were scanned by Microscopes International using an automatic scanning microscope, fitted with a 10x ocular and a 40x objective lens, equivalent to a numerical aperture of 650. JPEG images at the resolution of the camera, 540x960, were produced.

A total of ten slides were provided for this study. Seven slides were positive for malaria and three slides were free of malaria. Each slide consisted of one thin and one thick blood film. Within each blood film, a 5 mm x 1 mm rectangle was imaged as a tessellation of 800 individual images of 540 x 960 pixels. Several images from each blood film were excluded due to faults, such as being blank, being too blue due to accumulation of stain, or being messy due to faults in the staining process. Some examples of thin and thick blood film images were shown previously in this Chapter in Figure 1.6 and 1.7. The slides came complete with parasitemia scores based on the evaluation of thick blood films from between 14 and 24 expert readers per slide (Table 1.1). Expert readers are certified microscopists with expertise in reading malaria slides.

Table 1.1: Parasitemia by experts. $N$ is the number of experts for each positive slide and $SD$ is standard deviation of their parasitemia.

| Slides | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|--------|---|---|---|---|---|---|---|
| Mean | 103 | 634 | 1340 | 5583 | 11442 | 32730 | 319393 |
| Median | 90 | 340 | 1440 | 5191 | 10885 | 32730 | 301996 |
| Minimum | 19 | 39 | 320 | 48 | 3051 | 15736 | 30180 |
| Maximum | 370 | 2987 | 2476 | 13500 | 36000 | 78300 | 976500 |
| *SD* | 77 | 716 | 518 | 2592 | 6740 | 18972 | 178080 |
| *N* | 20 | 24 | 24 | 24 | 24 | 14 | 24 |

In addition to expert readers, this thesis involved two "informal" readers. The informal readers do not have formal microscopic training, but have gained substantial experience in viewing malaria parasites. Malaria parasites are not hard to identify in many cases. In difficult cases, the variation between trained microscopists, and even for the same microscopist viewing the images at different times, is similar to the performance of informal readers. In this thesis, experts' parasitemia scores were used to validate parasitemia estimation of each slide in final testing. The informal readers assigned erythrocytes and non-erythrocytes for testing erythrocyte identification, normal and infected erythrocytes for training and testing erythrocyte classification, and parasites and non-parasites for training and testing parasite classification.

## 1.12 Structure of the thesis

This thesis contains seven chapters including the current introduction chapter (Chapter 1). Chapter 2 reviews theories and techniques used in this thesis for developing the final methods of automatic malaria diagnosis. The remaining chapters describe the methods developed for automatic malaria diagnosis.

Methods for automatic malaria diagnosis based on thin blood film images are developed in Chapters 3 and 4. In malaria diagnosis, the methods presented in Chapter 3 use histograms of normal and infected erythrocytes; whereas, Chapter 4 develops methods that identify malaria parasites separately from erythrocytes. Also, a new feature derived from the natural characteristics of components of blood fluid to determine an adaptive threshold for image segmentation process is introduced.

Methods for automatic malaria diagnosis based on thick blood films are developed in Chapters 5 and 6. Chapter 5 is a preliminary study to determine the extent to which physical appearance and size can be used to distinguish between parasites and thrombocytes. Meanwhile, Chapter 6 develops methods by segmenting leukocytes and malaria parasites. Then, leukocytes and malaria parasites are identified and counted to estimate parasitemia scores. Finally, Chapter 7 concludes the study and outlines further research opportunities.

# Chapter 2: Technical Background and Literature Review

## 2.1 Overview

This thesis comprises extensions and applications of image analysis to estimate parsitaemia from thick and thin blood films. This chapter provides background on image normalization (Section 2.2) and mathematical morphology used to manipulate object shapes for supporting image pre-processing (Section 2.3). Segmentation methods are presented in Section 2.4 and methods for identifying objects resulting from segmentation are presented in Section 2.5. Features used in this thesis are described in Section 2.6 and median filters that can be used to remove small noise are discussed in Section 2.7. Section 2.8 describes Gaussian mixture models used to resolve object histograms into three groups of attributes and Section 2.9 explains classifiers. Methods for selecting useful features and classifier performance are discussed in Section 2.10, 2.11 and 2.12. Section 2.13 provides a review of literature reporting work similar to that conducted in this thesis.

## 2.2 Image Normalisation

Images from different sources or acquired using different acquisition conditions or parameters are often not consistent in terms of overall brightness or colour. This may cause inconsistencies in the performance of automatic image analysis algorithms. The purpose of image normalisation is to establish consistency of luminance or colour over a set of images. In this section, two different methods are described which will be used in the following chapter.

### 2.2.1 Colour Constancy

The aim of colour constancy is to obtain the characteristics of the target objects that are independent of variation of luminance. Historically, this was compensating for varying light sources and so the description bellow is in terms of illuminant. There are different models of luminance transform. The simplest model that offers reasonable results is the diagonal model (Barnard 1999). This model was derived from the von Kries hypothesis (Kries 1878) that "human colour constancy is an independent gain regulation of the three cone signals, through three different gain coefficients".

Illuminate colour normalization is given by $p^o = Dp^u$, where: $p^o$ is the normalised image, $D$ is diagonal-matrix transform, and $p^u$ is a colour image of unknown illuminate

(original image). In other words, the diagonal-matrix transform maps the colour response of unknown illuminate images to those of known illuminate images (Barnard, Cardei & Funt 2002; Tek, F B, Dempster & Kale 2006).

$$D = \begin{bmatrix} D_{rr} & 0 & 0 \\ 0 & D_{gg} & 0 \\ 0 & 0 & D_{bb} \end{bmatrix}$$

The image, $p$, is represented as a matrix of size 3x$N$, where $N$ is the number of pixels and column $k$ lists the R, G, B values of pixel $k$.

The illuminate factors, the non-zero elements of the diagonal matrix $D$, are calculated by the ratios of the average values of each channel of the reference images ($\mu_i^c$) to those of unknown illuminate image ($\mu_i^u$). The reference value(s) can be derived from the average values of the scene under a known illuminate.

$$D_{ii} = \frac{\mu_i^c}{\mu_i^u} \quad i = \{r, g, b\},$$

$$\mu_i^c = \frac{1}{N}\Sigma_1^N I_i^c \quad i = \{r, g, b\},$$

$$\mu_i^u = \frac{1}{N}\sum_1^N I_i^u \quad i = \{r, g, b\}$$

where, $I^c$ is the reference image. This is also called the database grey-world algorithm (Barnard, Cardei & Funt 2002; Hordley & Finlayson 2004). The database grey-world algorithm assumes that under uniform illumination, the reference image has stable average values for the image set.

## 2.2.2  Histogram Normalisation

An alternative to reduce the effects of different luminance in grayscale images is through normalization of the histogram of grayscale values. Histogram normalization is a technique based on viewing the distribution of intensities as a discrete probability distribution. There are several approaches to histogram normalization in grayscale images. "Median normalisation" is a simple and satisfactory method. The normalised image, $I^n$, is the ratio of the difference between the grayscale image and the minimum value to the difference between median and minimum value.

$$I^n = \frac{I^u - \alpha}{M - \alpha} \tag{2.1}$$

where, $I^u$ is the original grayscale image value, M and α are the median and the minimum values of grayscale pixels of the image $I^u$. One variation is to use the mean instead of the median in the formula above.

The median and mean of a normal distribution are not significantly different; however, natural data, commonly, has an asymmetric or skewed distribution and (or) outliers. If a distribution is highly skewed, the median may be a more useful measure of the centre and if there are outliers, the median is more robust than the mean (Moore, McCabe & Craig 2012).

## 2.3 Mathematical Morphology

Mathematical morphology is a method for manipulating the shape of objects manifested on binary or grayscale images (Haralick, R. M., Sternberg & Zhuang 1987). Mathematical morphology, as a part of lattice theory, was introduced in 1984 by Serra where objects are visually perceived in the framework of a lattice (Serra 1984). This method was further developed by the International Society for Mathematical Morphology. This method offers many useful tools for image analysis including dilating, eroding, opening, closing, and filling of shapes. Some important definitions and functions, which are used extensively in this thesis, are described below in the context of binary images.

### 2.3.1 Structuring element

The basic principle of mathematical morphology is to manipulate a shape by set addition with a structure element. Let *A* denote a set of 'on' pixels in a binary image and let *B* denote a set of on pixels in another binary array. The basic morphological operation is

$$C = A \oplus B = \{a + b | a \in A, b \in B\}$$

The operation $\oplus$ is known as Minkowski set addition. In the simplest form, a structuring element consists of a small binary array specified by shape and size (Efford 2000). The shape (e.g., disk or diamond) is represented by the pattern of ones and zeros.

The two basic morphological operations are dilation and erosion. To dilate shape $A$ by structure element $B$ means to compute $A \oplus B$. To erode shape $A$ by structure element $B$ means to compute $A \ominus B = (A^C \oplus B)^C$, where $A^C = 1 - A$ is the complement of $A$.

### 2.3.2 Opening and closing

Opening is a morphological image process consisting of two operations: erosion of a given image, $A$, based on a structuring element, $B$, followed by dilation based on the same structuring element as the erosion operation. The opening of $A$ by $B$ is denoted $A \circ B$

$$A \circ B = (A \ominus B) \oplus B$$

Conversely, closing is a combination of dilation followed by erosion.

$$A ¥ B = (A \oplus B) \ominus B$$

In principle, opening and closing operation methods depend on the shape and size of structuring element, and objects of segmentation.

Morphological operations provide useful fundamental characteristics for image filtering and segmentation tasks, such as eliminating small holes and sharp peaks, smoothing the contour of a segmented object, cutting off narrow isthmuses, and removing thin protrusions (Serra 1984; Vincent 1993b). In object identification, Haralick et al. reviewed two methods and concluded that mathematical morphology operations are more beneficial than convolution operations (Haralick, R. M., Sternberg & Zhuang 1987). Furthermore, Vincent used morphological operations to develop the opening by reconstruction algorithm for geodesic image analysis (Vincent 1993b). In the experiment, a disc of radius 2 was used as an opening to separate several connected components in a segmented image and, subsequently, a closing operation was employed to reconstruct the separated components. In industrial application, morphology operations were successfully used to detect defects in ceramic tiles (Elbehiery, Hefnawy & Elewa 2005). In this application, a dilation operation was used to enhance the defects and a fill gap operation was employed to increase the clearance of the cracks. In this paper, the shape and size of the structuring elements were not discussed clearly. Disk shaped structure elements have been used in opening methods for cell segmentation (Dorini, Minetto & Leite 2007; Ruberto, C. et al. 2000).

## 2.4   Segmentation Methods

Image segmentation is a division of an image into objects based on their characteristics (Gonzalez, Woods & Eddins 2001). Every pixel in the image is categorized to one of the objects. Two different segmentation methods involved in this thesis are reviewed below.

### 2.4.1   Segmentation based on boundaries: Canny's Method

One way to segment objects is to identify their boundaries. This requires edge detection as a first step. Canny's algorithm provides optimal localization of edges  in general situations (Canny 1986). Edge localization can significantly reduce data size and remove useless information, but keep important morphological image properties. For these reasons, Canny's method has been used in many image analysis tasks including segmentation in medical images (Huang, Y-L et al. 2008).

An automatic system involving Canny's algorithm was proposed by Huang, Y-L et al. The system used a Canny edge detector for edge detection of cells, the dilation operator to enhance the detected edges, and the opening operation to fill the cell bodies and smooth the cell footprints. The experimental results indicated that the proposed method well determined the outlines of a cell. The dataset consisted of 2573 cells from 45 images comprising six different patterns: 519 cells in diffuse patterns, 482 cells in peripheral patterns, 788 cells in coarse speckled patterns, 634 cells in fine speckled patterns, 64 cells in discrete speckled patterns, and 86 in nucleolar patterns. The system recognized 2130 cells consisting of 444 cells, 389 cells, 688 cells, 479 cells, 54 cells, and 76 cells in diffuse pattern, peripheral pattern, coarse speckled pattern, fine speckled pattern, discrete speckled pattern, and nucleolar pattern, respectively (Huang, Y-L et al. 2008).

### 2.4.2   Segmentation based on intensity regions: Otsu's Method

Another way to segment images is to identify regions of fixed intensity. Thresholding techniques may be used to separate foreground objects from background. Success often depends on choosing an appropriate threshold and many applications require that such a threshold be chosen automatically. Otsu's method for selecting a threshold (Otsu 1979) applies if the distribution of grayscale levels in an image is bimodal with one mode representing objects of interest. Otsu's method works by finding the threshold that separates the image into two regions of maximum internal homogeneity. Since image intensity

variance is a measure of homogeneity, the threshold is chosen to minimize the region variance. Otsu's method returns binary images that may consist of many objects.

## 2.5 Object Identification Methods

Object identification is a procedure for recognizing objects and/or its attributes. Three different algorithms employed in this thesis are presented below.

### 2.5.1 Circle Hough Transform

The circle Hough transform (CHT) is used to detect circles. The Hough transform was introduced in 1962 (Hough 1962) and first used a decade later to find lines and curves in images (Duda & Hart 1972). The CHT measures the response obtained by assuming that the image contains circles of radius $r$ centred at $(a,b)$ for every combination of $a$, $b$, and $r$ (Duda & Hart 1972). The stronger the response, the larger the evidence that such a circle exists.

Liangwongsan et al. investigated the use of CHT in detecting circle defect patterns of hard disk drives. Each input image was converted from Cartesian to parameter space by the CHT. A voting procedure was carried out with the maximum value indicating the corresponding radius for the circle. Then, the maximum value was compared with a threshold representing the boundary of the circle pattern to be classified. The threshold value was not discussed in the study. On a database of 120 defect media images, the system achieved an overall accuracy of 95.8% at five faulty detections (Liangwongsan et al. 2011)

### 2.5.2 Distance Transform

In a binary image, the distance transform assigns to every pixel in the foreground, the distance to the nearest pixel in the background. The distance transform is useful in many situations (Borgefors 1984; Danielson 1978; Huang, CT & Mitchell 1994) including object separation and identification.

In some cases, binary images resulting from Otsu's method (Section 2.4.2) are not properly segmented from the background, especially if there are touching and overlapping objects. To improve separation of the touching objects, Binghan et al. used the distance transform and dilation operator to separate objects in binary hepatitis pathology images (Binghan et al. 2002).

In another implementation, Basalamah used a distance transform-based histogram to identify circle objects (Basalamah 2012). More specifically, the pixels in every cell were

scanned and distances between the pixels and all the edge pixels were calculated. A histogram was used to count the frequency of the distances. The bin with the largest number of distances identified the radius of the circle. The method is robust and able to detect circles and partial circles.

The distance transform has been used to find the centres of erythrocytes (Sunarko et al. 2013). Morphologically, erythrocyte forms resemble circles which means that the centre of the binary footprint of an erythrocyte corresponds to a local maximum (Section 2.5.3) of the distance transform. This means that it is possible to detect erythrocytes by identifying circles. The centroids are determined by measuring distance from border objects. The local maxima of the distance transform are considered to be the centre of the objects. The algorithm could find perfectly circular objects; however, overlapping objects and imperfect objects were not appropriately recognized.

In another study, Yadollahi and Prochazka employed mathematical processing consisting of the Watershed Distance Transform, Gradients and Region Growing Algorithms to separate overlapping objects in binary images by tracing boundary pixels and converting to polar coordinates and then smoothing the curve. The boundary line of occluded objects was detected by connecting local minima pixels. This algorithm is appropriate for separating two overlapping objects in a microscopic image; however, this algorithm has difficulty separating many occluded cells (Yadollahi & Procházka 2011).

### 2.5.3 Regional Maxima

In 2D images that are the output of the distance transform (Section 2.5.2), the regional maxima algorithm is a procedure which finds a local maximum of connected neighbourhood pixels. The algorithm returns binary images: pixels that are the locations of a local maximum value are replaced by 1; other pixels are set to 0.

Vincent conducted a study on using regional maxima in developing area openings and closings. By employing the regional maxima algorithm, the computation of area openings and closings can be both computation and space efficient (Vincent 1993a)

## 2.6 Features

A feature is a measurable property of an object being observed that characterizes and discriminates the object (Bishop 2006). Choosing informative, discriminating and

independent features is a crucial step for effective algorithms in pattern recognition. Features used in this thesis are described below.

### 2.6.1 Statistical Moment Features and Colour Channel Histogram

In general, the recent technology of colour cameras creates and display colour images by mixing the colours: red (R), green (G), and blue (B) in different proportions (Efford 2000). Grayscale space and the other colour spaces can be derived from the primary colours. For examples: the space set of hue (H), saturation (S), and intensity (V) and the space set of luminance (Y) and colour information (Cb and Cr). For purposes of image processing, colour channels may be quantified and represented as a histogram. This allows computation of statistical properties.

Statistical properties from any colour channel (R, G, B, H, S, etc.) may be computed over the image or regions of the image. Examples include central moments such as the mean, variance, skewness and kurtosis. Central invariant moments were introduced in the 1960s (Hu 1962) where the invariant moments were used as features to develop pattern recognition and data simulation. Dudani et al. employed central invariant moments of transformations, elevation angle, and distance to recognize aircraft types (Dudani, Breeding & McGhee 1977). The automatic system was shown to be more accurate than human observers.

In another study, Zitova and Flusser used invariant moments (the mean and standard deviation) of translation, rotation, and blurring to estimate camera motion. From an original image, 18 images were generated by rotating over two angles: $6.2^0$ and $10.8^0$, translating one pixel in vertical and horizontal directions, and blurring manually with defocus and two different foreign object insertions. The computed values were compared with the ground truth to evaluate the performance. The estimation accuracy was satisfactory with mean and standard deviation of differences in rotation: 0.06 and 0.05; vertical translation: 0.1 and 0.4; and horizontal translation: -0.3 and 0.5 (Zitová & Flusser 2002).

Also, Avci and Varol used seven invariant moments as features to classify parasite eggs in microscopic images. After a segmentation process, Hu's seven invariant moments of the parasite egg footprints were extracted in order to classify human parasite eggs. The Hu's seven invariant moments are invariant under translation, rotation, and scaling. Based on a sample set of 938 parasite eggs from sixteen human parasite egg types, classification using the multi-class support vector machine classifier resulted in an overall correct classification rate 97.70% (Avci & Varol 2009).

26

### 2.6.2 Ellipticity and Eccentricity

Two useful features for describing shape are eccentricity and elliptical irregularity. Eccentricity is a very well establish concept although there are many variants (Ballard & Brown 1982; Burger & Burge 2016; Jähne 2005; Jain 1989). In principle, eccentricity is a measure of elongated shape. For an ellipse with major and minor semi axis lengths $a$ and $b$, the eccentricity is commonly defined by

$$\text{Eccentricity} = \sqrt{1 - \frac{b^2}{a^2}}$$

Eccentricity has been used for centuries to describe elliptical shaped objects, especially the orbits of planets; however, in 1992, eccentricity was employed to classify a trajectory as a straight line, a curve to the left or a curve to the right (Charayaphan & Marble 1992). Combined with other features, the eccentricity was used to interpret motion in American Sign Language. Five out of nine signs could be classified well by the proposed method.



Figure 2.1: Ellipticity and Eccentricity. a is semi-major axis, b is semi-minor axis, A and B are area inside and outside of the best fitting ellipse.

Elliptical irregularity was introduced in (Kruk et al. 2016) and is used to quantify the extent to which a shape differs from an ellipse. For any given shape in the plane (Figure 2.1), consider all ellipses having the same area as the shape. For each such equal-area ellipse, the area within the shape but outside the ellipse is a measure of how different the shape is from the ellipse. The equal-area ellipse for which this area is minimum is called the best fitting

ellipse. Let A denote the area of the shape. Note that A is also the area of the best fitting ellipse. Let B denote the area within the shape and outside the best fitting ellipse. Then the elliptical irregularity of the shape is given by

$$\text{Elliptical irregularity} = \frac{B}{A}$$

Because the shape and the best fitting ellipse have the same area, B is also equal to the area within the best fitting ellipse lying outside the shape. For an object of any shape, the best fitting ellipse can also be used to define the eccentricity of the shape.

### 2.6.3 Heterogeneity

Heterogeneity is a measure of how decentralised a distribution is. Young used heterogeneity as a feature to characterize chromatin distribution. Based on the histogram of luminance distribution (Figure 2.2), the heterogeneity was defined as

$$\text{Heterogeneity} = \frac{N_B + N_W}{N_B + N_G + N_W}$$

where, $N_B$, $N_G$, and $N_W$ are the number of pixels with value less than 80% of mean, between 80% and 120% of mean, and more than 120% of mean, respectively. Compared with morphological features, heterogeneity performed well in practice (Young, Verbeek & Mayall 1986).

In 2011, Elter et al. included heterogeneity as one of the features of a classification system for detecting malaria parasites in thick blood film images. Based on a sample set of 878 region of interests (ROIs) (266 parasites and 612 non-parasites), the experimental study showed that the classifier was able to identify ROIs with a high detection sensitivity (0.97) and low false-positive detections per image (0.80) (Elter, Haßlmeyer & Zerfaß 2011).

Figure 2.2: An example of histogram of a ROI. Thresholds are automatically selected at 20% above and below the mean value of the histogram.

## 2.7 Median Filter

A median filter is a non-linear method that may be used to remove noise while preserving edges, sharpen signals or images in many situations where an averaging filter could be used but a more robust method is needed. The median filter works as follows. From a moving finite window of real numbers, the neighbouring pixels in an $m$-by-$n$ neighbourhood are sorted (the $i$th order statistic of $N$ numbers $X_{1, ..., } X_N$) and the new value of the central pixel is the median given by $X_{(1+N)/2}$ if $N$ is odd and $(1/2)(X_{N/2} + X_{(N+1)/2}$ if $N$ is even, where $N = mn$.

In 1974, Tukey introduced one dimensional median filtering to smooth signals in time series analysis (Tukey 1974). In the following year, in combination with a linear filter, the median filter was utilized to smooth speech signals and reported quite efficient and competitive results (Rabiner, Sambur & Schmidt 1975). In the same area, Jayant made a comparative study of the use of the median filter and the mean filter in the transmission of digital speech signals and concluded that the two methods had the same performance for independent error occurrences (Jayant 1976).

Meanwhile, Pratt developed two-dimensional median filters by using a two-dimensional window with a size and shape to suppress noise and enhance images in image processing. Based on the examination of a various sizes of windows for the two-dimensional median filter, the 3x3 two-dimensional filter was recommended for significantly reducing

impulse noise in images (Pratt 1978). Furthermore, Huang et al. made a study to improve the speed of the two-dimensional median filter. By using a fast two-dimensional median filter, median filtering can be far more efficient because it is simple to update the histogram from window to window (Huang, TS, Yang & Tang 1979).

Based on some recent work in the actual implementation of the median filter, several benefits in using the median filter have been seen. First, there is no reduction in contrast across steps, since output values available consist only of those present in the neighbourhood. Secondly, for small to moderate levels of (Gaussian) noise, the median filter is demonstrably better than Gaussian blur at removing noise whilst preserving edges for a given, fixed window size (Arias-Castro & Donoho 2009). In addition, the median filter is effective for removing speckle noise and salt-and-pepper noise in medical imaging although the median filter performance is not much better than Gaussian blur for high levels of noise (Arce 2005; Jayaraman, Esakkirajan & Veerakumar 2009).

## 2.8 Gaussian Mixture Model

A Gaussian mixture model (GMM) is a probabilistic model for describing characteristics of groups within a set of data points (Permuter, Francos & Jermyn 2003). The data points are modelled as having been generated from a mixture of a finite number of Gaussian distributions with unknown parameters. Let $X = (X_1, ..., X_L)$ be a data set and $i \in \{1, 2, ..., k\}$ denote the index for group sets. By Bayes theorem, the Gaussian density function is

$$f(x) = \sum_{i=1}^{k} p_i \, N_i(x|\mu_i, \sigma_i^2)$$

$$N_i(x|\mu_i, \sigma_i^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp(\frac{-(x - \mu_i)^2}{2\sigma_i^2})$$

where $p_i$ denotes the mixture proportion (the prior probability distribution) of each group and $N_i(x|\mu_i, \sigma_i^2)$ represents a multi-variate normal distribution with mean $\mu_i$ and standard deviation $\sigma_i$. These prior probabilities satisfy: $\sum_{i=1}^{k} p_i = 1$ and $0 \leq p_i < 1$.

One numerical method for estimating the parameters of the Gaussian distributions is the expectation-maximization (EM) algorithm (Farnoosh & Zarpak 2008). The EM

algorithm is an alternative algorithm for the maximum likelihood estimation in the practical application (Fu & Wang 2012). The EM algorithm consists of two steps: E-step and M-step. The E-step finds the expected value of the complete likelihood given the current parameterization, $\theta^r$, while the M-step looks for the set of parameters $\theta^{r+1}$ that maximizes the expectation from the E-step (McLachlan & Peel 2004). Implementation of the EM algorithm requires the following four steps (Farnoosh & Zarpak 2008).

1. Set $k$ groups and initialize parameters:

$$\theta^{(0)} = (p_1^{(0)}, \dots, p_k^{(0)}, \mu_1^{(0)}, \dots, \mu_k^{(0)}, \sigma_1^{2(0)}, \dots, \sigma_k^{2(0)})$$

2. E-step: compute the probabilities according to

$$p_{ij}^{(r+1)} = \frac{p_i^{(r)} N(x_j | \mu_i^{(r)}, \sigma_i^{2(r)})}{f(x_j)}$$

3. M-step: update the parameter values as

$$\hat{p}_i^{(r+1)} = \frac{1}{n}\sum_{j=1}^{n} p_{ij}^{(r)},$$

$$\hat{\mu}_i^{(r+1)} = \frac{\sum_{j=1}^{n} p_{ij}^{(r+1)} x_j}{n\hat{p}_i^{(r+1)}},$$

$$\hat{\sigma}_i^{2(r+1)} = \frac{\sum_{j=1}^{n} p_{ij}^{(r+1)} (x_j - \hat{\mu}_i^{(r+1)})^2}{n\hat{p}_i^{(r+1)}}.$$

4. Iterate steps 2 and 3 until:

$$\sum_i e_i^2 < \varepsilon,$$

where $e_i$ is an error, e.g. $e_i = \hat{\mu}_i^{(r+1)} - \hat{\mu}_i^{(r)}$ and $\boldsymbol{\varepsilon}$ is a pre-set tolerance.

In 1998, EM-GMM was successfully used to estimate human face skin colour (Yang & Ahuja 1998). The results demonstrate that the estimated GMM fits with skin images from a database. Reynolds et al. employed the modified EM-GMM for speaker verification (Reynolds, Quatieri & Dunn 2000). The results indicated that the method improves verification performance. Zivkovic developed an adaptive algorithm using GMM to improve image segmentation by background subtraction (Zivkovic 2004). The adaptive GMM was able to select an appropriate number of components for each pixel. In the study by Huang et al., the use of GMM showed exceptional performance in classifying multiple limb motion (Huang, Y et al. 2005).

## 2.9 Classifiers

A classifier can be defined as a procedure that assigns an item to one of two or more classes. The classes are known *a priori* and one or more new observations are classified into one of the known groups based on the measured characteristics. In this section, two classifiers are described that will be used in the following chapters.

### 2.9.1 Discriminant Analysis

In this classifier, there are a fixed number of classes and known feature values of training objects. The goal is to find a decision rule (pattern) to classify a query point $x$ into the corresponding class. Let $X$ be the number of data points in group $G$. The discriminant analysis technique works as follows.

For every sample of a group $G$ in the training set, consider $x$ in the feature space. The goal of discriminant analysis classifier is to find a good predictor for the group $G$ of any sample from the training set. Suppose $f_k(x)$ is the density function of $X$ in group $G = k$, and let $\pi_k$ denote the prior probability of group $k$, where $\sum_k^K \pi_k = 1$. Using Bayes' theorem, an estimate of the posterior probability of a sample $x$ being member of group $k$ is

$$P(G = k | X = x) = \frac{f_k(x)\pi_k}{\sum_{k=1}^K f_k(x)\pi_k}$$

Subsequently, the point $x$ in the feature space is assigned to the group $k$ with the largest posterior probability $P(G = k | X = x)$. The theoretical details and more information can be found in texts (Hastie, Tibshirani & Friedman 2001; McLachlan 2004)

Let $\pi_1$ and $f_1(x)$ be the prior probability and the density function of group 1, respectively and $\pi_2$ and $f_2(x)$ be the prior probability and the density function of group 2, respectively. A sample $x$ is assigned to group 1 if $P(G = 1 | X = x) \geq P(G = 2 | X = x)$, otherwise the $x$ is assigned to group 2.

Among the classifiers, discriminant analysis is one of the most commonly used methods. A benefit of a discriminant analysis classifier is that it is quick to train (Lyons, Budynek & Akamatsu 1999). Also, the mathematical form is simple and the method is easy to implement while still being very powerful (Raudys, Š & Young 2004). Linear discriminant analysis is guaranteed to provide optimal classification if the distributions of the feature values in the groups are normal. Originally, discriminant analysis was developed

by Fisher to classify three species of iris flowers based on three different features (Fisher 1936).

In the study by Altman et al. the linear discriminant analysis was found to be promising for corporate distress classification, with classification accuracy of over 0.90 (Altman, Marco & Varetto 1994). Also, the linear discriminant analysis was acceptable for classifying antibiotic resistance patterns of indicator bacteria (Harwood, Whitlock & Withington 2000). In 2001, Hadjiiski et al. employed linear discriminant analysis to classify masses as benign or malignant (Hadjiiski et al. 2001). In classifying benign and malignant masses on a data set of 348 mass ROIs comprising 169 benign and 179 malignant, the classifier achieved the mean $A_z$ score of 0.78.

Lu et al. proposed linear discriminant analysis for human face recognition (Lu, Plataniotis & Venetsanopoulos 2003). Results showed excellent performance with only a very small set of features being used, its misclassification rate was 0.34. In another study in the field of economics, the linear discriminant analysis was used to predict financial distress of public companies listed in the Amman stock exchange (Al-khatib & Al-Horani 2012).

## 2.9.2 Exhaustive Search Classifier

Discriminant analysis automatically finds a discriminant surface that separates the classes in the training set. The surface found is guaranteed to be optimal if the underlying distributions for the considered features are normal. In cases where the features do not have normal distributions, better classification may be possible by considering non-linear discriminant surfaces using more complicated methods such as neural networks or support vector machines. These methods require large amounts of data to avoid over fitting –a natural consequence of allowing very complicated discriminant surfaces.

An example of a method that retains much of the robustness of linear discriminant analysis but is more flexible in situations where the underlying distributions are not necessarily normal is to set thresholds for individual features rather than a single linear surface within the entire feature space. This must be accompanied with a rule for assigning class labels to the many compartments of the feature space generated by the individual thresholds and so requires good understanding of the expected contribution of individual features to the final decision. For this reason, setting thresholds for individual features is not commonly used in very complex classification problems.

In most applications, the best thresholds for each feature are not known ahead of time and so different combinations of thresholds must be considered to train the classifier. If the number of features is not too large, exhaustive search may be used to find the optimal combinations of thresholds. The feasibility of doing so also depends on the number of thresholds per feature that need to be considered –in other words, the resolution required for the threshold values.

The method described above will be referred to as 'exhaustive threshold classification'. Exhaustive threshold classification is closely related to classification based on decision trees (Hastie, Tibshirani & Friedman 2001). In decision trees, an optimal threshold is found for each individual feature using exhaustive search. The feature with the best performing threshold is used to segment the features space into two compartments according to the optimal threshold. Next, the same process is repeated separately on both compartments and these steps are repeated until a stopping condition is reached. The resulting classifier consists of a set of rules for assigning classes. A typical rule is of the form: assign class 1 if $f_1 > T_1$ and $f_3 > T_2$ and $f_2 < T_3$ or if $f_1 > T_1$ and $f_2 < T_4$ and $f_3 > T_5$ and assign class 2 otherwise. Here $f_i$ is the value of feature i and $T_j$ is threshold j. The thresholds are set sequentially as the optimal threshold for the compartment at hand. The rule for assigning classes is determined automatically as part of the process.

In exhaustive threshold classification, all the thresholds are considered simultaneously but the rule for assigning classes is fixed beforehand. In cases where the rule for assigning classes is indeed available as prior information, then exhaustive threshold classification is guaranteed to find the combination of thresholds that provides the best classification. The optimal combination of thresholds is not guaranteed by the decision tree due to the sequential nature of finding the thresholds.

## 2.10 Feature Selection

The aim of feature selection is to select a small set of meaningful features instead of using many features. In other words, feature selection is used to determine a subset of features that generates the optimal classification performance. There are several reasons why feature selection is crucial. Firstly, classification based on a complex decision rule and large number of features may need careful attention (Raudys, SJ & Jain 1991), and most importantly, using a large number of features relative to the number of samples in the data will result in artificially high classification. For any data set, arbitrarily good classification

may be achieved on training data if enough features are used but the results will usually be poor on testing data and in actual use. Secondly, the use of ineffective or redundant features may distract the classifier and result in unstable training. This means that the selected features are unreliable and are likely to lead to poor performance on test observations (Chan, Petrick & Sahiner 2000; Fukunaga 1990; Miller et al. 2003). Thirdly, wide-ranging features may result in high computational cost. On the other hand, careful feature selection can improve the performance of predictors.

The performance evaluation of candidate methods in feature selection commonly requires a search strategy in selecting candidate subsets and an objective function. In the study by Dash and Liu, objective functions, called evaluation functions, are categorised into five groups: distance (also known as separability), information (or uncertainty), dependence (or correlation), consistency, and classifier error rate (misclassification rate) (Dash & Liu 1997). Generally, the misclassification rate is used as the objective function in performance evaluations.

Two feature selection methods used in the following chapters, exhaustive search and sequential feature selection, are reviewed below.

## 2.10.1 Exhaustive Search Feature Selection

Exhaustive search feature selection consists of evaluating all possible combinations of features to find the combination that results in the best value from the objective function.

Given the number of original features is $n$, there are $2^n - 1$ possible feature subsets. This means that exhaustive search feature selection is practical only for relatively small numbers of $n$. However, this method for feature selection is one of the few methods that guarantees finding the optimal feature subset. If there are n features, but attention is restricted to selecting subsets of at most k features, then the number of feature sets to consider is M(n,k) = C(n,1) + C(n,2) + … + C(n,k), where C(n,i) is the binomial coefficient C(n,i) = n!/(n-i)!i!. Even for large n, M(n,k) is reasonable if k is not too big. As an example, for n = 100, there are about $10^{158}$ total subsets, but M(100,3) = 166,750.

## 2.10.2 Sequential Feature Selection

Sequential feature selection methods can be grouped into two main methods: sequential forward feature selection and sequential backward feature selection. Sequential forward feature selection starts with the empty set. All features are tested individually and the one

with the best performance is added to the set. At each step, all remaining features are tested in combination with the features already in the set until no features improve the performance or the improvement is less than a prescribed minimal improvement condition. There is no guarantee that this method finds the optimal feature subset.

On the other hand, sequential backward feature selection starts from a full candidate set and features are sequentially discarded from the candidate set until the removal of further features does not improve the objective function. Sequential backward feature selection generally works best when the optimal feature subset is large. However, sequential backward feature selection is not able to re-evaluate the usefulness of features after they are removed.

## 2.11 Cross-validation

Cross-validation is a technique for estimating how well a classifier is likely to perform in making new predictions on unseen data (Hastie, Tibshirani & Friedman 2001). Cross-validation splits the whole data set into $K$ parts: $K$-$1$ parts are used for training and the remaining part for testing. This process is repeated K times, so that each of the K parts plays the role as the testing data once.

Cross-validation was introduced in the 1931 by Larson (Larson 1931) and further developed by Mosteller and Wallace (Mosteller & Wallace 1963). In 1983, Efron developed a prediction rule involving cross-validation for estimating classification error rate (Efron 1983). The results indicated that cross-validation offers almost unbiased error estimation but has high variance, so may result in unreliable estimates. However, this drawback may be ignored, especially if the data set is small.

Sometimes, to improve the reliability of estimation performance, the entire $K$-fold cross-validation process is repeated multiple times (Refaeilzadeh, Tang & Liu 2008), each time with the data set divided into a new set of $K$-folds. The assignment into folds is usually random, but with consistent representation per class.

## 2.12 Accuracy, Sensitivity and Specificity

In the field of classification, accuracy is the proportion of correct prediction over the total number of cases including positive and negative ones.

$$\text{Accuracy} = \frac{\text{the number of correct assignments}}{\text{the total number of actual cases}}$$

36

Accordingly, accuracy is used to measure the classification performance. Generally, the higher accuracy, the better classification performance.

Table 2.1: Confusion matrix with two groups.

|       |          | Prediction |          |
|-------|----------|------------|----------|
|       |          | Positive   | Negative |
| Actual | Positive | *TP*       | *FN*     |
| Cases  | Negative | *FP*       | *TN*     |

However, in the framework of medical analysis, misdiagnosis due to an error in diagnosing a patient as having disease when the patient is actually free disease has different consequences compared with a misdiagnosis as being free from disease when the patient actually does have the disease (Metz 1978). For this reason, the utilization of accuracy is not enough to be a performance indicator alone. Better measures of performance are given by sensitivity and specificity. These are based on the confusion matrix.

Two useful measures of performance are sensitivity and specificity are given by:

$$\text{Sensitivity} = \frac{\text{the number of } TP}{\text{the total number of } TP \text{ and } FN}$$

$$\text{Specificity} = \frac{\text{the number of } TN}{\text{the total number of } TN \text{ and } FP}$$

A general confusion matrix is a square matrix, $L \times L$, containing the values of actual and predicted classification (Kohavi & Provost 1998). $L$ denotes the number of groups. In a two group classification problems, there are four possible fractions: true positive prediction (*TP*), false positive prediction (*FP*), true negative prediction (*TN*), and false negative prediction (*FN*) (Table 2.1).

Classification schemes may be adjusted to maximize sensitivity or specificity or, more usually, a linear combination of the two, according to the relative consequences of errors due to false positive reports compared to errors due to false negative reports. In the context of screening mammography programs, for example, high sensitivity is desirable since missing a true cancer may be fatal while erroneously classifying an anomaly as cancerous will likely be rectified by follow-up procedures. However, in the context of deciding if a severely invasive and dangerous intervention is required, high specificity is desirable.

## 2.13 Literature Review

Many of the techniques of image normalisation and pattern recognition reviewed in Sections 2.2 to 2.12 have been used for computer-aided malaria diagnosis based on thin and thick blood film images. Automatic estimation of malaria parasitemia based on microscopy images is the ultimate objective of this thesis. This section reviews previous work on identification of malaria parasites and parasitemia estimation on thin and thick blood film images.

In the last two decades, many studies have reported on the use of image processing for medical diagnosis. The fundamental problems in medical diagnosis based on image analysis are segmentation and classification. Both fields are very large and only aspects related to this thesis will be reviewed.

### 2.13.1 Detecting Malaria Parasites

The scope of this section is to review previous work in the field of the identification of infected erythrocytes.

Colour and morphological features have been developed for the identification of infected erythrocytes and normal erythrocytes from an image. Makkapati and Rao evaluated dominant hue range in the hue, saturation, and value (HSV) colour space of images to segment erythrocytes from the background and identified an optimal saturation threshold to detect malaria parasites. Erythrocytes were segmented by thresholding using Otsu's method (Section 2.4.2) and a saturation value of 0.34 was manually determined as an optimum cytoplasm threshold. To check its sensitivity and specificity, the method was applied to 55 images taken from Leishman-stained blood films. The segmentation yielded 88% sensitivity and 95% specificity (Makkapati & Rao 2011). The parasitemia level was not investigated.

Ruberto et al. utilized granulometric functions on grayscale images based on size and shape, and employed colour transformation from red, green, and blue (RGB) to HSV. Parasite nuclei were identified by the intersection of regional maxima of image $H$ and $S$. Then, the mean grayscale of the parasite nuclei computed on the image $H$ and $S$, $\mu H$ and $\mu S$, were used as threshold values to segment the objects of interest in image $H$ resulting in image $TH$, and in image $S$ to obtain $TS$, respectively. Finally, leukocyte and parasite footprints were obtained from the intersection of $TH$ and $TS$, yielding image $THS$. Morphological erosion with a disk-shaped structuring element of size 22 (Section 2.3.1) was employed to remove

leukocyte footprints in *THS* images and a morphological smoothing was used to remove stray points and close small holes. Infected erythrocytes were identified and isolated by applying morphological reconstruction based on dilation of *THS* (Soille 2013). On a database of 12 images, a microscopist found 1,910 erythrocytes with 479 of the erythrocytes being infected. Parasitemia was 25.08% for infected erythrocytes. The algorithm found 1,953 erythrocytes, of which 506 were labelled as infected. The parasitemia was estimated to be 25.91% (Ruberto, C. et al. 2000). However, a parasite classification was not conducted to distinguish parasites or artefacts. This means that all parasite candidates, *THS* minus leukocyte footprints, were assumed as parasites. Moreover, the parasitemia scores (the number of infected erythrocytes per µl of blood) were not reported.

In another granulometry-based study, Ross et al. utilized morphological, colour intensity, and texture features of erythrocytes (natural feature candidates) for erythrocyte classification as normal or infected. Ross et al. used thresholding for erythrocyte and parasite classification. Otsu's method (Section 2.4.2) was applied to segment erythrocytes from the background. Subsequently, the first minimum after the principal mode of the erythrocyte histogram was taken to be the parasite threshold. Six features were measured on the histogram of erythrocytes. The six features were size, eccentricity, smoothness, colour, texture, shape of parasites, and the number of parasites per erythrocyte. On a database of 2361 infected erythrocytes, the system achieved a sensitivity of 92.07% and a positive predictive value of 39.64% with 1378 false positives and 78 false negatives. After testing, the method correctly identified 85.13% of the total number of erythrocytes and the accuracy of the classification system was 73% (Ross et al. 2006). However, the parasitemia level was not studied and thrombocytes were not considered.

Kumar et al. developed a malaria identification scheme by applying adaptive thresholding with Otsu's algorithm (Section 2.4.2) on the blue channel of RGB colour images in order to segment erythrocytes and parasites. Parasites and erythrocytes were segmented and counted separately. Erythrocytes were segmented by the Otsu threshold value and the Otsu threshold value plus 0.25 was used to segment the blue colour channel to obtain parasite footprints. In this study, there was no classification of parasite and erythrocyte footprints. After applying morphology to fill the holes of the erythrocyte footprints and dilation or erosion with disk shaped structure elements to remove noise, parasite and erythrocyte footprints were counted. The method was examined using a set of images containing 87 parasites and 311 erythrocytes from eight malaria-infected images. Following implementation, the scheme incorrectly recognized 18 of 311 erythrocytes and 14 of 87

parasites in the testing set. The study reported that the classification results on parasites were close to the manual counts even though some differences were observed regarding erythrocyte counts. In the frame of diagnosis, the parasitemia score was not computed in the study. In addition, in the classification strategy, thrombocytes, which naturally present in thin blood film images, were not considered. The result suggested that Otsu's algorithm should be implemented on the green channel instead of the blue channel (Kumar et al. 2012).

Savkare and Narote investigated support vector machines to classify erythrocytes into two classes: normal erythrocytes and infected erythrocytes. Ten features were extracted from erythrocyte footprints to estimate parasitemia. The ten features were geometric (radius, perimeter, area, compactness, and metric), colour saturation, and statistical features (skewness, kurtosis, energy, and standard deviation). To evaluate the performance, the automatic parasitemia was compared to manual reader parasitemia values. On a database of 15 images, the system had sensitivity of 93.12% and specificity of 93.17% (Savkare & Narote 2011). However, the validations were performed on parasitemia score per image and parasitemia score per slide was not discussed.

Kim et al. identified and classified erythrocytes and leukocytes in blood images automatically. The watershed transform was used to segment the images and then the k-means algorithm was implemented to merge the nearest regions based on colour features. After identifying cells, a neural network model, with inner edges and counter features as input, was developed to classify erythrocytes and leukocytes. Furthermore, Kim et al. compared performance to two other algorithms: learning vector quantization-3 (LVQ-3) and k-nearest neighbour (KNN). Tested on a database of 680 erythrocytes and 410 leukocytes, the overall classification rate was 91% and 81% for erythrocytes and leukocytes, respectively compared with KNN (82%) and LVQ-3 (89%). The accuracy of the neural network classifier was better in this experiment (Kim et al. 2001). However, thrombocytes, that may resemble erythrocytes in size and intensity, were not considered and occluded erythrocytes, which naturally present in thin blood film images, were not discussed.

Diaz et al. proposed the use of a genetic algorithm for pattern recognition in identifying infected erythrocytes. Five different histogram features were measured to distinguish normal or infected erythrocytes. The five histogram features were: colour, saturation level, grayscale, Tamura texture, and Sobel histogram. Based on a sample set of 450 malaria images, classification using mean, standard deviation, skewness, kurtosis, and entropy of the histograms yielded a sensitivity of 94% and a specificity of 99.7% for

identification of infected erythrocytes (Díaz, González & Romero 2009). However, thrombocytes were not considered in the erythrocyte classification strategy.

Tek et al. reviewed works on automatic malaria diagnosis based on microscopic thin and thick blood film images. The two main emerging areas reported by Tek et al. in which image processing was most likely to have an important role in malaria diagnosis were segmentation and classification. In thin blood film images, segmentation aims at separating the foreground, commonly consisting of erythrocytes, leukocytes, and parasites (if present), and the background. The typical problems of segmentation were identified as under segmentation or over segmentation. Under segmentation produces two or more blood cells in one region (occluded cells). Conversely, over segmentation results in a single blood cell segmented into two or more regions. Several methods have been presented to overcome these problems, such as morphological area closing and distance transform. However, none of these methods is applicable to high-occluded cells. In terms of classification, few classification studies were conducted to distinguish parasites and other stained objects and the proposal was to differentiate between normal and infected erythrocytes. In the frame of thick blood film images, only a preliminary study was reported. Overall, studies reported by Tek et al. resulted in useful systems to identify malaria parasites (Tek, F. B., Dempster & Kale 2009). However, these studies in thin blood films were not compared to parasitemia estimated by microscopists in thick blood films.

Furthermore, Tek et al. investigated the use of image processing and pattern recognition methods to diagnose malaria based on thin blood film images. Probability density was used to determine stained cells, local area granulometry was applied to estimate size, and a modified KNN classifier was employed to detect malaria-infected erythrocytes in thin blood films. A set of 630 colour images from nine slides consisting of 669 infected erythrocytes and 3431 normal erythrocytes was used. A subset of 336 infected erythrocytes and 1645 normal erythrocytes was used for training. The other 333 infected erythrocytes and 1645 normal erythrocytes were used for testing. In detection experiments, an accuracy of 93.17% was achieved, with sensitivity of 72.37% and specificity of 97.45%. In the per slide detection performance, the algorithms needed to observe at least 45,889 erythrocytes to achieve a sensitivity of 72.37% in a slide of 500 parasites/µl. In the same slide, the classifier yielded a specificity of 97.45% with 114 false detections (Tek, F. B., Dempster & Kale 2010). The parasitemia estimation was not discussed.

In thick blood film images, Toha and Ngah developed a system for computer aided malaria diagnosis. After converting to grayscale images, a threshold value *T2* was used to segment the grayscale images for obtaining parasite footprint candidates (binary images). Then, cluster analysis using the Euclidean distance algorithm and size and locality was used to determine parasites. The system was tested by counting the number of malaria parasites in one image (Toha & Ngah 2007). However, the parasite count was not clearly reported and validation of the result was not discussed.

Another study of malaria identification in thick blood film images was conducted by Elter et al. They focused on high detection sensitivity while accepting potentially low specificity per image and then reduced the number of false-positive detections to an acceptable level while maintaining the high detection sensitivity. The authors computed the arctan of the ratio of green and blue colour intensity to transform colour images into grayscale ones. Next, the black-top-hat morphological operator followed by threshold segmentation (the proportion of the green and the blue colour intensity as the threshold value) was used to separate potential malaria parasites from the background. A set of features consisting of statistical moment features, texture features, and colour features were extracted from each malaria parasite candidate. Then, a genetic algorithm was applied to the set of features to classify malaria parasite candidates. Of the 266 parasites and the 612 false-positive detections, at a reasonable high sensitivity of 0.97, this system achieved 3.2 false-positive detections per image without false-positive reduction and 0.8 false-positive detections per image with false-positive reduction. This demonstrated that the algorithm was suitable for the development of an automatic malaria diagnosis based on microscopic images (Elter, Haßlmeyer & Zerfaß 2011). However, parasitemia score was not presented.

In the same year, Hanif et al. focused on enhancing the quality of malaria-infected thick blood film images. Dark stretching was used to make images clearer and threshold segmentation was applied to separate parasites from the background. The threshold value was controlled manually. The method was tested on three images and three different threshold values (200, 220, and 230) were applied to each stretched image. Visually, the results showed that the dark stretching method can improve the image quality and the three threshold values, especially 220, were suitable for segmenting the three different stretched images to obtain malaria parasites (Hanif, Mashor & Mohamed 2011). Nevertheless, the study did not discuss parasitemia scores.

## 2.13.2 Estimation of Parasitemia

The estimation of infection level of malaria parasites may be determined based on the percentage of infected erythrocytes in thin blood films or the parasite quantitation in thick blood films.

A study by Diaz et al., also described above, counted the percentage of infected erythrocytes to determine parasitemia. Of 12,557 erythrocytes, 12,513 erythrocytes were classified automatically as infected or uninfected and 44 erythrocytes were manually evaluated by a user. The total assessment yielded the parasitemia of 5.6% for infected erythrocytes. The study went further to determine the life form of parasites using a multi-level perceptron classifier. The scheme was tested on the same data set consisting of 12,557 erythrocytes with 11,844 healthy erythrocytes, 521 erythrocytes infected by parasites in ring form, 109 infected by parasites in trophozoite form, and 83 infected by parasites in schizont form. Evaluated by experts, an average sensitivity of 78.8% and average specificity of 91.2% were reported for life stage identification (Díaz, González & Romero 2009). Meanwhile, estimation of parasitemia scores was not reported.

Frean estimated parasitemia scores based on thick blood films using open access software. Based on a range of 497 thick blood film images with varying levels of parasitemia (12 slides), strong correlation was achieved between the parasitemia estimation resulting from the software and that from microscopists (R = 0.99). Furthermore, on a database of 197 images from eight slides, the parasitemia scores resulting from this software correlated well (R = 0.97) with that from microscopists (Frean 2009). However, this software was run semi-automatically: leukocytes were counted manually and a parameter, the radius of parasite area, was also adjusted manually to remove outliers.

# Chapter 3: Classifying Erythrocytes for Estimating Parasitemia in Thin Blood Film Images

## 3.1 Overview

This chapter presents a study on estimating parasitemia based on thin blood films. Estimating parasitemia from thin blood films requires counting individual erythrocytes and determining the proportion that are infected with the malaria parasite. To do this automatically requires several linked steps: detecting erythrocytes, segmenting erythrocytes, determining which erythrocytes are infected with the malaria parasite and then translating the ratio of infected erythrocytes to parasitemia scores that may be compared with human expert scores.

These steps are presented in the chronological order in which they were addressed during this study instead of the order in which the steps are implemented in practice because the findings from the preliminary studies impacted choices made in subsequent steps. Accordingly, thin blood film images used in this study are described in Section 3.2. Section 3.3 presents a study on comparing two methods for segmenting erythrocytes given that an erythrocyte is known to be present at a certain location. This study was taken on initially to verify that segmentation was possible in principle before addressing the other tasks needed for a full system. Section 3.4 presents a method for full automatic segmentation of erythrocytes without prior knowledge of the location of erythrocytes in the image. Section 3.5 demonstrates the method used to determine if an erythrocyte is infected with a malaria parasite or not and Section 3.6 provides a method for relating the proportion of infected erythrocytes to parasitemia scores determined by human experts.

## 3.2 Thin Blood Film Images

For this part of the study, 610 images were randomly selected from the thin blood films of the seven positive slides (Section 1.12), approximately 87 images from each slide. Of these, twenty images were used for the erythrocyte segmentation study (Section 3.3), thirty images were used for training of erythrocyte classification (Section 3.5), 280 images were used for training the process for estimating parasitemia, and 280 were used for testing the process for estimating parasitemia. Typical thin blood film images contain two basic components (erythrocyte and thrombocytes) and some additional objects (parasites and occasionally some leukocytes).

44

## 3.3 Segmentation of known erythrocytes

Recognizing erythrocytes is an important topic in automatic malaria diagnosis in thin blood films. A preliminary study to the main project of automatic malaria diagnosis from a thin blood film was conducted to select an appropriate erythrocyte segmentation method. The work was conducted on a limited data set of images. This preliminary study described in this section and the main study described in the following sections were presented in Engineering International Conference and published in the proceedings (Sunarko et al. 2013) and (Sunarko et al. 2017), respectively.

For training and testing erythrocyte segmentation and classification, experimental results in this section were compared to ground truth established manually by a microscopist. This was because most normal and infected erythrocytes are clearly distinguishable. For parasitemia estimation, experiment results were compared to expert readers.

Two different approaches to segmenting erythrocytes were compared. One method was based on edge-detection using Canny's method followed by the circular Hough transform. The other method was based on setting a threshold using Otsu's method followed by the distance transform. The study demonstrated that the use of a threshold followed by the distance transform to segment erythrocytes from the background is superior to edge detection followed by the Hough transform. This preliminary result provided the green light for the remaining investigation forming the bulk of this thesis.

The reason for including this preliminary study here is that an important discovery made impacted both the processing steps and the focus of the main project. A naive implementation of Otsu's method and distance transform leads to reasonably good segmentation and identification of erythrocytes. The experiment leading to these conclusions is described in the following sections.

### 3.3.1 Images for Erythrocyte Segmentation Study

Twenty images from seven positive slides of thin blood films (Section 1.11) were used in this experimental study. Of these, two images were from slide 1 and three images were from each of the other slides. An example of an original colour image is shown in Figure 3.1.

Figure 3.1: An example of original thin blood film image. The image shows malaria free erythrocytes and thrombocytes.



(a)

(b)

(c)

(d)

Figure 3.2: Edge detection method for identifying erythrocytes. (a) An output of Canny's operator applied to the image in Figure 3.1. Most detected edges were not continuous and some erythrocytes were converted into two or more different circles. Some edges not related to erythrocytes were also found. (b) An output of the circle Hough transform (CHT) applied to Figure 3.2(a). Red plus signs indicate circle centers found by CHT. Each circle center represents a candidate erythrocyte. Most erythrocytes were detected as more than one erythrocyte. Several erythrocytes at the border of the image were not detected. (c) An output of a clustering applied to the circle centres on Figure 3.2(b). Some circle centers in the same circle were grouped and represented by a circle center. (d) An output of the CHT and clustering applied to Figure 3.1. A false detection (indicated by the index 28) can be seen lying between four true erythrocytes, in the upper-right quadrant.

### 3.3.2 Experimental Details

For the processing steps described here, the original colour images were converted to grayscale images based on luminance. MATLAB was used to implement the Otsu and Canny-based segmentation. The procedure for recognising erythrocytes is demonstrated by example images as shown in Figure 3.2 and Figure 3.3.



Figure 3.3: Threshold method for identifying erythrocytes. (a) Output of the Otsu's method applied to Figure 3.1(a). Some neighbouring erythrocytes in the original image were joined. Erythrocyte areas in binary images are wider than that of in the original images. Some erythrocytes, in the middle, were below threshold. In the upper-right corner, there is a cell with center below threshold. (b) Figure 3.3(a) after applying the median filter and flood-fill operation. (c) The distance transform was applied to Figure 3.3(b). Some erythrocytes have many local maxima (+). All erythrocytes were detected, including many at the edge, some non-erythrocytes were also detected. (d) Clustering was applied to Figure 3.3(c). Each erythrocyte was represented by a single +. (e) Identified erythrocytes of Figure 3.1.

### 3.3.2.1 The Edge Detection Method

Canny's algorithm (Section 2.4.1) was applied to the grayscale images to produce binary boundary images. To reduce noise from the binary images, median filtering (Section 2.7) using a 3x3-structure element was implemented. This resulted in many complete or nearly complete curves nominally corresponding to erythrocyte boundaries. In addition, some cell boundaries appeared as double boundaries and many spurious lines were found which were not associated with erythrocytes (Figure 3.2(a)). The CHT (Section 2.5.1) was applied to the boundary images to identify individual erythrocytes.

Applying the CHT required setting a minimum radius, $R_{01}$, and a maximum radius, $R_{02}$, for the circles to be detected as the boundaries of erythrocytes. For very small values of $R_{01}$ and $R_{02}$, many small line segments were falsely detected as candidate erythrocyte boundaries. On the other hand, for very big values of $R_{01}$ and $R_{02}$, small erythrocytes were not detected and leukocytes were detected as candidate erythrocyte. From some preliminary runs with images containing only erythrocytes, the value $R_{01} = 20$ and $R_{02} = 40$ were determined empirically to provide reasonable radii to retain erythrocytes and ignore leukocytes and thrombocytes.

The outcome of this was a set of three values location centre (x-y coordinate) and radius for each candidate of erythrocyte. Figure 3.2(b) shows all detected circles including erythrocyte-like objects, in which not all the objects represent true erythrocytes. To determine which candidate erythrocytes represented true erythrocytes, the average radius, $\mu_R$, and standard deviation, $\sigma_R$, of radii obtained from the CHT were calculated and the clustering method (Sunarko et al. 2013) described in the following section was used.

### 3.3.2.2 Simple Clustering Method

If several erythrocytes form a local cluster of cells, then their centres will be at least two cell radii apart. If candidate centres are closer than this number, then at least one of them is not the centre of a true erythrocyte and should be removed. In addition, if an erythrocyte is not perfectly round, then the Hough transform may result in several centres representing the same erythrocyte. These observations motivate the following algorithm for removing candidate centres that do not represent true erythrocytes.

Let $P_j$ denote the centre of an erythrocyte candidate.

For $P_j$, $j = 1, 2, 3, \ldots, n$

1. For every $i \neq j$, compute the distance $D_{ij}$ between $P_i$ and $P_j$

2.      If $D_{ij} > \mu_R - 3\sigma_R$, no change,

3.      Otherwise, remove $P_i$ or $P_j$ depending on which has the shorter radius.

### 3.3.2.3   The Threshold Method

Otsu's method (Section 2.4.2) was applied to the grayscale images resulting in binary images showing the footprints of erythrocyte as white spots (Figure 3.3(a)). Some erythrocytes were well segmented as isolated single erythrocytes, for example erythrocyte number 5 (Figure 3.3(d)); however, if erythrocytes are occluding or if they are just touching or closing together, for example erythrocyte A and B (Figure 3.1), only a single footprint Y is obtained for these cells (Figure 3.3(b)). This requires separate processing. For simplicity, the phrase "occluded erythrocytes" will be used to refer to cells that are just touching each other and cells that are actually occluding.

A median filter (Section 2.7) was implemented to remove noise from the binary images and a flood-fill operation was applied to the filtered image to fill the holes in erythrocyte footprints due to the lower intensity (Figure 3.3(b)). Next, the distance transform (Section 2.5.2) was applied to the filled images and the regional maxima values (Section 2.5.3) of the distance transform image were taken as the centres of erythrocytes. Subsequently, the radius of the erythrocyte was taken to be the distance between the centre and the minimum distance to the boundary. However, since erythrocytes are not perfect circles, one candidate can have many local maxima all having the same distance to the boundary. In addition, local maxima outside erythrocyte were found due to background noise or thrombocytes. As a result, there are many possible centres for each erythrocyte (Figure 3.3(c)). The most appropriate centres for erythrocytes were determined by selecting erythrocyte footprints, distinguishing single erythrocytes from occluded erythrocytes, and employing the simple clustering algorithm described in the previous section on occluded erythrocytes. Details of these steps follow.

First, bounds were found for the areas of erythrocyte footprints. To do this, 150 manually established regions of interest (ROI) of infected and normal erythrocyte were considered. These were used to determine a minimal area as $_{min}A = 1300$ pixels, a maximum area as $_{max}A = 3400$ pixels, and minimum radius of $_{min}R = 20$ for erythrocyte footprints. Footprints with area less than $_{min}A$ were assumed to be noise or thrombocytes and were not considered further. Footprints with area greater than $_{max}A$ were assumed to be leukocytes or occluded erythrocytes.

Leukocytes are very rare in thin blood film images since there are many fewer leukocytes than erythrocytes (Section 1.6); therefore, misidentifications of erythrocytes due to leukocytes do not significantly affect the counting of the number of erythrocytes. Thus, footprints due to leukocytes can be regarded as occluded erythrocytes without contributing much error. Subsequently, a simple clustering algorithm described in Section 3.3.2.2 (removing $P_i$ or $P_j$ depending on which has the lower image intensity) was applied to the assumed single and occluded erythrocytes containing many erythrocyte centres to determine the true erythrocyte centres.

### 3.3.3   Analysis Methods

The performances of the two methods for segmenting erythrocytes were compared in terms of accuracy and sensitivity. Subsequently, a pairwise t-test was used to compare each of the methods with human observers. This resulted in a probability (p-value) that the observed difference was due to chance alone.

### 3.3.4   Results and Discussion

#### 3.3.4.1   The Edge Detection Method

Figure 3.2(d) shows outputs of the image segmentation based on the edge detection method. The results of comparing a human reader to Canny based segmentation are summarized in a confusion matrix (Table 3.1). Of the 654 erythrocyte samples, 575 were correctly identified, providing an accuracy of 86.73% and sensitivity of 88.06%.

Table 3.1: Confusion matrix for erythrocyte classification performance of Canny's algorithm and the CHT. The accuracy was 86.60%.

|  |  | Predicted | |
|---|---|---|---|
|  |  | Erythrocyte | Non-erythrocyte |
| Actual class | Erythrocyte | 575 | 79 |
|  | Non-erythrocyte | 10 | 0 |

Although many erythrocyte were detected, the number of true detections reported by the edge-based segmentation was significantly lower than the human reader at the $\alpha = 0.05$ level (n = 20, p < 0.001). The method did not identify enough whole erythrocytes, especially erythrocytes at the edge of the image (Figure 3.2(b)). In addition, 10 blank areas lying between two or more erythrocytes, were incorrectly detected as erythrocytes (FP). The CHT used to recognize the erythrocytes is primarily to blame. The method is fundamentally driven

by recognition of circle shapes. The boundaries of some true erythrocytes did not form sufficiently consistent circles to be recognised and many false detections resulted from isolated line segments being recognised as circular boundaries. In addition, the method did not take into account the intensity patterns within the candidate circular boundaries. Therefore, this method is prone to inaccuracy.

### 3.3.4.2 The Threshold Method

Figure 3.3(e) shows the result of the image segmentation based on the threshold method. Potential erythrocytes were identified from the images with an accuracy of 99.37% and a sensitivity of 96.63% (Table 3.2). Furthermore, a paired t-test showed that there was no significant difference between human readers and the algorithm in terms of the number of true detections at the $\alpha = 0.05$ level ($p = 0.431$, $n = 20$). However, uncommon erythrocyte shapes lead to some false detections (*FP*), detected as more than one erythrocyte, and some missed erythrocyte (*FN*).

Table 3.2: Confusion matrix for erythrocyte classification performance of the Otsu's method and the distance transform. The accuracy was 95.90%

|  |  | Predicted | |
|---|---|---|---|
|  |  | Erythrocyte | Non-erythrocyte |
| Actual class | Erythrocyte | 631 | 23 |
|  | Non-erythrocyte | 4 | 0 |

In terms of *FP*, the Otsu-based segmentation identified a biconcave erythrocyte (Figure 1.5) as two erythrocytes. Also, the Otsu-based segmentation may not properly detect erythrocytes with unusual shapes (Figure 3.3(b, d, and e)). Unusual shapes can result from proximity to the boundary of the image or from the thresholding step if the intensity is not sufficiently different from the background.

### 3.3.5 Conclusion

These experimental results show that the Otsu-based segmentation gave similar results to a human reader. The Canny-based segmentation provided significantly poorer outcomes than a human reader ($n = 20$, $p < 0.001$). Hence, Otsu-based threshold segmentation is a better method for this task. Accordingly, Otsu's method and distance transform were used for erythrocyte segmentation and erythrocyte identification in the remaining experiments.

## 3.4 Fully Automatic Segmentation and Classification of Erythrocyte

Comparing the percentage of infected erythrocyte in thin blood films to parasitemia found by expert readers of thick blood films requires a process for automatically identifying large numbers of erythrocyte and classifying these as normal or infected.

In order to develop such a system, the approximate minimal erythrocyte area was determined (Section 3.3.2.3) to identify a single erythrocyte. Although the sizes of erythrocyte varied between slides, the value of $_{min}A$ was small enough to not eliminate true erythrocytes but large enough to eliminate large numbers of thrombocytes and artefacts in all images.

Steps for isolating footprints of individual erythrocyte has been described previously (Sunarko et al. 2013) but are included here for completeness. Connected objects in the binary images were extracted by a labelling algorithm (Haralick, Robert M. & Shapiro 1992). At this stage, connected components include artefacts, remaining thrombocytes, clusters of erythrocytes, regions touching the border of the image as well as proper erythrocytes (Figure 3.3(d)). The distance transform was used on each footprint to determine the distance of each pixel within the footprint to the boundary and a regional maxima algorithm was used to find the local maxima of this distance function. If the footprint is a circular disk, there is a unique local maximum of the distance function at the centre but for other shapes, there may be several local maxima. Accordingly, footprints were separated into three groups; those with area less than $_{min}A$ which were discarded as representing thrombocytes or artefacts, those with area greater than $_{min}A$ and containing many local maxima of the distance function viewed as possible clusters of erythrocyte or single erythrocyte with shape significantly different than circular, and those with area greater than $_{min}A$ but with only one local maximum. Footprints in the latter group were identified as representing a single isolated erythrocyte. With this criterion, it became possible to automatically identify isolated erythrocyte in images for use in further processing steps as described below.

The process for separating clusters of erythrocytes required knowing the average area, $\mu A$ and the standard deviation $\sigma A$ of the areas of the footprints of erythrocyte. These values were found not to be consistent from one image to the next due to differences in focus and possibly other image acquisition parameters. Accordingly, the method required that $\mu A$ and $\sigma A$ be computed separately for each image, unlike in Section 3.3 where a minimum size was found for several images from all slides. Thus, for each image processed, the isolated

erythrocytes identified by having area greater than $_{min}A$ and a single local maximum of the distance function were used to calculate $\mu A$ and $\sigma A$ for that image.

For footprints with area greater than $_{min}A$ and multiple local maxima of the distance function, the following steps were used to identify isolated erythrocyte and separate clusters into individual erythrocyte. The average radius $\mu R$ and standard deviation $\sigma R$ for isolated erythrocyte (assuming round erythrocyte footprints) were computed from $\mu A$ and $\sigma A$. Let $C_i$, $i = 1, 2, \ldots, n$ denote the local maxima of the distance function over a single footprint. Each $C_i$ is viewed as a candidate for representing an isolated erythrocyte. Let $D_{ij}$ denote the distance between $C_i$ and $C_j$. If $D_{ij} < \mu R - 3\sigma R$, either $C_i$ or $C_j$ was removed from the list of candidate erythrocyte depending on which had the lower image intensity value. This process was repeated until $D_{ij} \geq \mu R - 3\sigma R$ for all $i, j$ for which $C_i$ and $C_j$ remained on the list of local maxima of the distance function over the footprint. The motivation for this process is that erythrocytes comprising a cluster are often not round and local maxima of the distance function tend to appear in small clusters near the centres of the erythrocytes comprising the cluster. Completion of this step allowed the number of erythrocytes comprising a cluster to be determined by using the algorithm described in Section 3.3.2.2.

## 3.5   Classifying ROIs as Infected or Normal Erythrocyte

Having established a method for identifying erythrocytes, the next objective is to determine which erythrocytes are infected with the malaria parasite and which are not.

This study for training of erythrocyte classification is presented in the three sections below. The data involved in this experiment is described in Section 3.5.1 and the experiment details are explained in Section 3.5.2. Finally, Section 3.5.3 presents results, discussion and conclusion.

### 3.5.1   ROIs for Training of Erythrocyte Classification

Thirty thin blood film images with known malaria parasites were proportionally selected from the seven positive slides described in Section 1.11 based on parasitemia level. Each selected image consists of normal and infected erythrocytes. A total of sixty ROIs of individual well isolated erythrocytes were visually identified from these images. Of the sixty ROIs, thirty were infected erythrocytes and thirty were normal erythrocytes.

### 3.5.2 Experimental Details

After image pre-processing (Section 3.3.2) and segmentation using the threshold method (Section 3.3.2.2), each erythrocyte was represented by the histogram of intensity values in the original grayscale image within the footprint of the erythrocyte. The rationale for referring to the grayscale image was that, while the image is in colour, the various components of the image are distinguishable due to differences in how well they take up the stain. In other words, salience is due to the level of staining and so the image in essentially monochromatic (Figure 3.4). Examples of histograms of infected erythrocytes and normal erythrocytes are presented in Figure 3.5. For each footprint, the skewness and kurtosis of the distribution of grayscale intensity values for the erythrocyte were extracted.



Figure 3.4 Colour distribution of a colour image. Ten-thousand pixels were sampled from one of the thin blood film images (the image at coordinates X = 0, Y = 40 in a positive slide). For each pixel, R, G and B values were plotted as point in 3-dimensional RGB space.

Generally, infected erythrocytes can be divided into three parts with different colours and intensities (Section 1.8). The host erythrocyte is light red (moderate intensity), the chromatin and cytoplasm are dark red/blue (low intensity), and the vacuole is transparent, its colour intensity depending on the effect of acidity (pH) during the staining process. The intensity of the vacuole is higher or equal to that of the host erythrocyte. Accordingly, a Gaussian mixture model with three Gaussians was used to separate the histogram of the grayscale values of the erythrocyte into three groups nominally representing the parasite, the vacuole and the remaining components of the erythrocyte. The Gaussian with the lowest

54

mean value was assumed to be the one corresponding to the parasites as this component is reliably the darkest within the cell if parasites are present. Thus, for a parasite-infected erythrocyte, the lowest mean of the three Gaussians was expected to be noticeably lower than the lowest mean of the three Gaussians for erythrocyte with no parasites. Accordingly, the lowest of the three Gaussian means was used along with the skewness and kurtosis to represent the erythrocyte as a feature vector of length three.



Figure 3.5: The histogram of grayscale values of two individual erythrocytes: (a) an infected erythrocyte has low intensity value, negative skewness, and stronger peak; otherwise, (b) a normal erythrocyte has high intensity value, positive skewness, and wide peak.

Linear discriminant analysis (LDA) was used to classify the erythrocyte as infected or not based on the three features described in the previous paragraph. LDA is an optimal classifier if the underlying distributions are normal, but not otherwise. Visual inspection of the distributions indicated that simple thresholds for the features might provide better classification. Accordingly, an exhaustive search over possible thresholds for the three parameters was also conducted. More specifically, let $m$, $s$ and $k$ denote the smallest mean, skewness, and kurtosis of the intensity distribution over the erythrocyte footprint and let $M$, $S$, and $K$ denote threshold values for these parameters. Then the erythrocyte was assigned as infected if ($s < S$ or $k > K$) and ($m < M$) and was assigned as normal otherwise. The accuracy of classification was tested for 21 equally spaced values of S in the interval (-0.1, 0.1), 10 equally spaced value of $K$ in the interval (1, 10) and 51 equally spaced values of $M$ in the interval (150, 200). These ranges were determined by inspection and by some preliminary runs.

For ease of exposition, the term "threshold classifier" will be used to refer to exhaustive search over threshold values described in the previous paragraph and the term

"LDA classifier" will be used to refer to the classifier based on LDA. The training error of LDA classifier was 0.1883 and the lowest training error of the threshold classifier was 0.1667. Although the threshold classifier provided better classification, it was not successful in identifying a unique set of values for the thresholds $M$, $S$ and $K$ because the same lowest training error was obtained for many of the parameter combinations tested. The points for which the minimum training error was attained formed a connected set, $\Omega$, in the $M$; $S$; $K$ feature space. The set $\Omega$ comprised 343 of the 10,710 points tested in the feature space. Optimal thresholds were eventually found as part of the process for translating percentages of infected erythrocytes to parasitemia scores (described in Section 3.6).

### 3.5.3 Results, Discussion, and Conclusion

The values of statistical features for classifying ROIs as infected or normal erythrocytes are presented in Table 3.3 The infected erythrocyte tended to have skewness and mean of the lowest group of Gaussian higher than that of a normal erythrocyte. On the other hand, the mean of kurtosis of the normal erythrocyte was slightly lower than that of infected erythrocyte. Figure 3.6 shows a scatter plot of statistical features from these ROIs and the classification performance is presented in a confusion matrix (Table 3.4).

Table 3.3: Features of infected erythrocyte and normal erythrocyte. The $s$, $k$, and $m$ of infected and normal erythrocytes are significantly different ($p < 0.001$, n = 10, $\alpha = 0.05$).

|  | Infected Erythrocyte | | | Normal Erythrocyte | | |
|---|---|---|---|---|---|---|
|  | $s$ | $k$ | $m$ | $s$ | $k$ | $m$ |
| Minimum | -1.0927 | 2.5680 | 108.8287 | -0.0137 | 2.0047 | 141.8255 |
| Mean | -0.1870 | 4.1217 | 150.0926 | 0.8116 | 3.3044 | 164.8371 |
| Maximum | 0.6517 | 5.9702 | 169.2149 | 1.3157 | 4.2995 | 187.5006 |

Table 3.4: Confusion matrix obtained from erythrocyte classification using the LDA classifier. The classification accuracy was 83.33%

|  |  | Predicted | |
|---|---|---|---|
|  |  | Infected Erythrocyte | Normal Erythrocyte |
| Actual Class | Infected Erythrocyte | 20 | 10 |
|  | Normal Erythrocyte | 0 | 30 |

According to Table 3.4 and Figure 3.6, it seems that the algorithm distinguishes infected and normal erythrocytes reasonably well. These results were based on a relatively

small number of examples of infected and normal erythrocytes. Furthermore, the erythrocytes used were manually selected to be isolated cells for which determining the true state was very reliable.



Figure 3.6. Scatter plots of the statistical feature values for ROIs. The red asterisk and blue squares denote infected and normal erythrocyte, respectively. X-axis, Y-axis, and Z-axis are skewness, kurtosis, and min-mean respectively. The classification error was 0.117.

A thorough study on a large number of examples to determine the performance of the algorithm in classifying erythrocytes in their full diversity of appearance as expected in practice was not conducted for two reasons. First, assigning the true state (infected or not infected) for a large and diverse sample is problematic and is likely to result in some incorrect assignments even if performed by an expert. Incorrect assignments lead to incorrect performance scores and may lead to systematic error in classification as the classifier will be trained to make some wrong decisions. Second, the final objective is not to count infected and normal erythrocytes but to estimate parasitemia. Accordingly, further refinement and further testing of the algorithm was conducted using expert parasitemia scores of the full sample as the objective rather than the correct assignment of individual erythrocytes. This work is described in the next section.

## 3.6 Estimating Parasitemia

Although automatically identifying erythrocytes and classifying these as infected or not was a necessary part of the study, the final objective was to determine if estimates of the level of

infection determined by computing the percentage of infected erythrocytes correlate to values of parasitemia reported by expert human readers. Thus, the linear discriminant analysis classifier trained in the erythrocyte classification step (Section 3.5) was used to determine the percentage of infected erythrocyte to total erythrocyte for all seven slides from patients with malaria. Forty images were randomly selected from each slide for a total of 280 images. This set of images is denoted by $T1$. The linear discriminant analysis classifier was applied to the images in $T1$ without further training. Since the true percentage of infected erythrocyte was not known for these images, it could not be used as a performance criterion. Instead, the correlation between the ratios of infected erythrocyte and the parasitemia values reported by human experts based on thick blood films from the seven slides was computed.

The analogous experiment could not be conducted using the threshold classifier since the study on segmenting erythrocyte (Section 3.4) failed to produce a unique set of thresholds $M$, $S$, and $K$ to implement the classifier. Instead, the set of images $T1$ was used to retrain the threshold classifier using the criterion of correlation with human expert values of parasitemia. Eight combinations of thresholds $M$, $S$ and $K$ representing vertices of $\Omega$ and midpoints between vertices of this set were used to compute correlation with parasitemia.

The optimal values of $M$, $S$, and $K$ determined by this method were then tested on a separate set of images. In this testing step, 40 new images were randomly selected from each slide for a total of 280 images. This set of images is denoted $T2$. The threshold classifier using the optimal values of $M$, $S$, and $K$ found in the training step using the set of images $T1$ was applied to the set of images $T2$ to estimate the correlation between the percentage of infected erythrocytes computed by the algorithm and parasitemia levels determined by human experts.

Expert parasitemia values from thick blood film ranged over several orders of magnitude from 90 to $3 \times 10^6$ parasites/µl (Table 1.1) while percentages of infected erythrocytes found by the algorithm ranged from about 1 to 28. Accordingly, a linear relationship between these quantities cannot be expected. Instead, the relationship $P = bB^a$ was tested, where $P$ represents parasitemia judged by expert readers on thick blood films, $B$ is the percentage of infected erythrocytes found by the algorithm from the 280 thin blood film images comprising testing set $T1$, and $a$ and $b$ are constants. Hence, a regression was performed on log $P$ and log $B$. The resulting regression formula was then used to predict parasitemia from the percentage of infected erythrocyte measured on the 280 testing images comprising set $T2$.

### 3.6.1 Results for Estimating Parasitemia

The set $\Omega$ represents combinations of threshold values for $M$, $S$, and $K$ that resulted in the same minimal classification error during the preliminary stage of classifying erythrocyte as infected or normal. In this case, performance was measured according to the percentage of correct classifications of individual erythrocyte as infected or not. In the second stage, the thresholds in $\Omega$ yielded highly varying results (Table 3.5). In this case, performance was measured according to correlation with expert estimates of parasitemia based on thick blood films. The best performing combination of threshold values $M = 170$, $S = -0.02$ and $K = 5$ was used to classify 280 images in the testing set $T2$. The correlation between the percentage of infected erythrocyte as judged by the algorithm with this set of thresholds was 0.875 (Table 3.6), only slightly lower than the correlation found using the same combination of thresholds during training on the set $T1$ (Table 3.5).

Table 3.5: Percentage of infected erythrocytes in thin blood films for selected vertices and mid-points of $\Omega$ for training data. r is the coefficient of linear correlation with experts parasitemia scores based on thick blood films.

| Threshold values | | | Percentage of infected erythrocytes | | | | | | | r |
|---|---|---|---|---|---|---|---|---|---|---|
| $M$ | $S$ | $K$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | |
| 170 | -0.02 | 5 | 0.929 | 1.788 | 2.071 | 2.900 | 3.726 | 6.654 | 21.530 | 0.883 |
| 170 | 0.05 | 3 | 8.388 | 1.845 | 72.945 | 3.196 | 18.530 | 54.005 | 88.790 | 0.324 |
| 173 | 0.01 | 6 | 0.880 | 2.019 | 1.545 | 2.989 | 3.974 | 6.298 | 18.683 | 0.869 |
| 170 | -0.01 | 10 | 0.734 | 1.845 | 1.216 | 2.752 | 3.626 | 5.980 | 16.904 | 0.858 |
| 173 | 0.02 | 6 | 1.076 | 2.076 | 1.611 | 3.078 | 4.272 | 6.428 | 18.361 | 0.863 |
| 170 | -0.02 | 4 | 2.959 | 1.845 | 19.855 | 2.900 | 5.763 | 17.539 | 47.865 | 0.500 |
| 180 | -0.02 | 5 | 1.736 | 3.460 | 2.104 | 5.090 | 4.968 | 6.977 | 21.530 | 0.803 |
| 175 | 0.01 | 7 | 0.954 | 2.711 | 1.249 | 3.403 | 4.471 | 6.428 | 17.972 | 0.799 |

Table 3.6: Correlation between percentage of infected erythrocytes in thin blood film and parasitemia values in thick blood film. $B$ is percentage of infected erythrocytes in thin blood film and $P$ is estimation of parasitemia values in thick blood film for the same slide as thin blood film. $C$ is estimation of parasitemia scores based on standard correlation (Equation 1.2). r is as in Table 3.5.

| Slides | 1 | 2 | 3 | 4 | 5 | 6 | 7 | r |
|---|---|---|---|---|---|---|---|---|
| $B$ | 0.978 | 1.663 | 2.169 | 2.716 | 4.565 | 6.594 | 28.430 | 0.875 |
| $P$ | 155 | 612 | 1220 | 2181 | 8371 | 21699 | 955269 | - |
| $C$ | 48900 | 83150 | 108450 | 135800 | 228250 | 329700 | 1421500 | - |

The regression formula resulted from training step, $P = 164B^{2.59}$, is shown in Figure 3.7(a) and the resulting parasitemia estimation, using the regression formula, in the testing step is shown Figure 3.7(b) and Table 3.6.

Figure 3.7: Predicting parasitemia. (a) Regression for training images. For each slide, the algorithm produces an estimate of the percentage of infected erythrocytes, the log of which provides the horizontal coordinate of the open circles. For each slide, there are several estimates of parasitemia from human experts (Table 1.1), the logs of which provide the vertical coordinates of the open circles. The regression line (solid line) is $\log(P) = \log(b) + a \log(B)$, where $B$ is the proportion of infected erythrocytes from the algorithm, $P$ is the mean of parasitemia estimated by experts, $a = 2.59$ and $b = 1.64 \times 10^2$. (b) The open circles are as in (a) but with the log of the percentage of infected erythrocytes taken from the testing images. Since the estimates for the testing images are not identical to the estimates from the training images, the horizontal coordinates are not exactly the same as for (a). The vertical coordinates of the open circles are the same as in (a) since the experts' estimates are fixed for each slide. Here $\times$ indicates the estimate of parasitemia found by applying the regression formula found in (a) to the percentage of infected erythrocytes in the testing images.

### 3.6.2 Discussion and Conclusion

The algorithms presented in this section may be used to determine the percentage of erythrocytes infected with malaria parasites and to estimate the parasitemia scores. From training, the formula $P = 164B^{2.59}$ was obtained. This formula gives an estimate of the parasitemia scores, the units used by human experts when reporting parasitemia based on viewing thick blood films. For the seven parasitemia scores available for this study, this formula predicts a score for every slide that fell within the range of values reported by the human experts (Figure 3.7(b)). Accordingly, the method developed here may be viewed as a plausible alternative to reading thick blood film slides in terms of accuracy.

In this study, forty images from each slide were used. In a fully developed system, the number of images could be left variable according to the number of infected erythrocytes encountered. If the percentage of infected erythrocytes in the first few images is high, then fewer images are needed to attain a robust estimate of parasitemia. If the percentage of infected erythrocytes is low or zero, more images would be needed to establish a reliable value of parasitemia. A device for automatic detection of malaria will have to work in close to real time. However, this study was conducted only to see if estimating parasitemia

automatically on thin blood film images is plausible. The algorithms presented here have not been implemented to minimise run time but care has been taken to select well established methods that are amenable to efficient implementation.

There are many possible sources of inaccuracy that could be addressed in future work. For example, carbon dioxide rich erythrocytes are darker in colour than other erythrocytes and hence detecting parasites in these erythrocytes may be less reliable. A separate algorithm for detecting carbon dioxide rich erythrocytes could reduce the rate of misclassification of infected and normal erythrocytes. In addition, the original colour images were converted to grayscale according to luminance. This may result in some loss of information and so there may be room for improvement.

The core of the algorithm is a quick estimate of the number of infected and normal erythrocytes. This part of the algorithm was developed on only sixty visually selected ROI comprising a single erythrocyte each. This part of the algorithm was not tested directly on an independent data set of ROIs. The reason is that the process of identifying and extracting ROIs manually is time consuming and impractical. In addition, selection of ROIs of this type does not result in a sample of erythrocytes that represents the wide range of appearance of erythrocytes in images. Also, regardless of the amount of refining of this step of the algorithm, some infected erythrocytes will be incorrectly classified as normal and some normal erythrocytes will be incorrectly classified as infected. Since the objective of the study is not to classify erythrocyte but to estimate parasitemia, the crucial question is not to determine the misclassification rate, but to understand the impact on estimating parasitemia. Accordingly, validation was performed on final estimates of parasitemia instead of classification rates of individual erythrocyte.

Although this study and that of Savkare and Narote (2011) similarly aim to estimate parasitemia, the respective results cannot be directly compared with each other due to significant differences in the validation process. While Savkare and Narote's validation yielded parasitemia per image, the study presented here went further by producing parasitemia levels per slide, which were validated using manual diagnosis by a number of experts based on thick blood films from the same subjects. This study clearly demonstrates that computer analysis of thin blood films is able to provide estimates of parasitemia that agree with human expert assessment of thick blood films. In this respect, the results of this study are more closely in line with that of Purwar et al. (2011), despite the latter's use of the same thin blood films in both automatic analysis and manual validation by a sole pathologist.

Meanwhile, the clinical method to correlate the percentage of infected erythrocytes in thin blood films to parasitemia scores (the number of parasite per µl) is found by multiplying the percentage of infected erythrocytes by the standard count of erythrocytes in a µl of blood (Equation 1.2). Based on the standard correlation (Equation 1.2), the parasitemia scores of this algorithm were shown in Table 3.6.

According to biological experiments (Dowling & Shute 1966; Trape 1985), parasitemia in thin blood films appears higher than that in thick blood films due to haemolysis (de-haemoglobinization) and the staining process. The loss of parasites during the processing of thick blood film varies greatly from 60% to 90% and 0% to 5% according to the experiment by Dowling and Shute and the experiment by Trape, respectively. For ease of exposition, the term "DS correlation" (Equation 3.1) and "TR correlation" (Equation 3.2) will be used to refer to the conversion of parsitaemia scores in thick blood films ($P_{thick}$) to actual parasitemia scores in thin blood films ($P_{thin}$) based on these two studies. Thus, the conversion using DS correlation is from

$$P_{thin\,DS} \in \left[ \frac{5}{4} P_{thick}, \ \frac{5}{3} P_{thick} \right]$$

(3.1)

and using TR correlation is from

$$P_{thin\,TR} \in \left[ P_{thick}, \ \frac{20}{19} P_{thick} \right]$$

(3.2)

Table 3.7: Parasitemia scores in thin blood film after conversion from experts' scores (Table 1.1) by DS correlation and TR correlation.

| Slides | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| DS | $4.80 \times 10^1$ - $3.70 \times 10^3$ | $9.80 \times 10^1$ - $2.99 \times 10^4$ | $8.00 \times 10^2$ - $2.48 \times 10^4$ | $1.20 \times 10^2$ - $1.35 \times 10^5$ | $7.63 \times 10^3$ - $3.60 \times 10^5$ | $3.93 \times 10^4$ - $7.83 \times 10^5$ | $7.55 \times 10^4$ - $9.77 \times 10^6$ |
| TR | $2.00 \times 10^1$ - $3.89 \times 10^2$ | $3.90 \times 10^1$ - $3.14 \times 10^3$ | $3.20 \times 10^2$ - $2.61 \times 10^3$ | $4.80 \times 10^1$ - $1.42 \times 10^4$ | $3.05 \times 10^3$ - $3.79 \times 10^4$ | $1.57 \times 10^4$ - $8.24 \times 10^4$ | $3.02 \times 10^4$ - $1.03 \times 10^6$ |

The algorithm presented here was moderately consistent with the work by Dowling and Shute (Dowling & Shute 1966). The system fitted the three high parasitemia slides well, slightly overestimated the middle parasitemia slide, and significantly overestimated the three low parasitemia slides (Table 3.7).

One reason for this discrepancy is that the algorithm is not applicable to crowded fields of thin blood films. One infected erythrocyte forming part of a normal cluster of occluded erythrocytes might cause all erythrocytes in the cluster to be detected as infected erythrocytes. This could increase the parasitemia level.

Another reason for over estimation of parasitemia is that some thrombocytes located near erythrocytes were captured by Otsu's method as parts of erythrocyte footprints. The number of thrombocytes is from around 5% of erythrocytes (Vander, Luciano & Sherman 2001) up to 8% of erythrocytes. Accordingly, this proportion is big enough to possibly play a role in over estimating parasitemia.

A dark normal erythrocyte touching, overlapping, or next to another normal erythrocyte may affect the histogram of the occluded normal erythrocyte due to the output of the Otsu's method. Thus, all normal erythrocytes forming part of a group of occluded normal erythrocytes might be identified as infected erythrocytes. If two or more erythrocytes are very close to each other, touching or overlapping then Otsu's method may result in a single footprint for these erythrocytes. By the process described in Section 3.3.4.2 the number of separate erythrocytes comprising this footprint will be determined for the purpose of counting erythrocytes. However, no processing steps have been included to actually separate these erythrocytes. Accordingly, a single disease state will be assigned to all these erythrocytes. Thus, a single infected erythrocyte or a single unusually dark erythrocyte may result in other nearby erythrocytes being incorrectly labelled as infected. For example, a dark normal erythrocyte, A, adjacent to another ordinary normal erythrocyte, B, (Figure 3.8(a) and (b)) was segmented by Otsu's method to be an occluded erythrocyte (Figure 3.8(c)). As a result, the skewness of the grayscale histogram of the occluded normal erythrocyte was negative (-0.21) and the lowest mean of the three Gaussians for the occluded erythrocyte was 128.21 (Figure 3.8(d)). Accordingly, both normal erythrocytes (A and B) were detected as infected erythrocytes. This also contributes to over estimation of parasitemia.

Despite the fact that plausible estimates of parasitemia were found using automatic image analysis on thin blood films, a final system based on this algorithm alone may not be practical for estimating very low levels of parasitemia or for reliably concluding that a subject is malaria free. Accordingly, the long-range role of the algorithm presented here is as part of a system that automatically analyses both thin and thick blood films. In combination, analysis of thin and thick blood films is more likely to provide accurate

estimates of levels of parasitemia for both very low and very high levels of infection and have the capacity to confidently declare subjects free of malaria.


(a)  (b)


(c)  (d)

Figure 3.8: Dark erythrocyte affecting feature of occluded erythrocyte. (a) Original image consisting of dark normal erythrocyte A and ordinary normal erythrocyte B. (b) the grayscale image of Figure 3.8(a). (c) An occluded erythrocyte footprint as the output of the Otsu's method applied to Figure 3.8(b) on erythrocyte A and B. (d) The histogram of grayscale value of the occluded erythrocyte matching with Figure 3.8(c).

To reduce the over estimation of parasitemia, parasites could be segmented independently of erythrocytes and then identified to determine infected erythrocytes. Such a process for identifying infected erythrocytes is described next in Chapter 4.

# Chapter 4: Parasite and Erythrocyte Identification for Estimating Parasitemia in Thin Blood Film Images

## 4.1 Overview

This chapter presents an alternative method for identifying infected erythrocytes and estimating parasitemia from thin blood films. The seven positive malaria thin blood films used for conducting experiments in this chapter were described in Section 1.11. This chapter, starting with Section 4.2, describes the thresholding method applied for parasite segmentation and establishing parasite threshold values. Leukocyte segmentation is discussed in Section 4.3. In Section 4.4, two candidate features for classifying parasite footprints are explored. This section also compares performance based on ROIs in Sections 4.4.1.1 and 4.4.2.1. This is followed by a method for identifying infected erythrocytes in Section 4.5. This section also evaluates performance of the method. Estimating parasitemia is presented in Section 4.6, which also includes the chapter discussion and conclusion.

## 4.2 Segmenting Parasites

In this study, the colour and brightness of the background was slightly different between images from the same blood films and significantly different between different blood films (Section 1.9.1 and Figure 1.6). An example of an original image containing erythrocytes, thrombocytes, and parasites is shown in Figure 1.6(a). To reduce the variation of backgrounds, the image normalization discussed in Section 6.3 was applied to the original images (Figure 1.6(a)). The image normalization described in Section 2.2.1 requires a reference image. In this case, a clear background image from one of the thin blood films was used as the image reference (Figure 4.1). Figure 4.2(a) shows the normalized image.



Figure 4.1: Reference image for image normalisation.

Figure 4.2: Pre-processing for parasite segmentation. (a) The normalized image of Figure 1.6(a) with Figure 4.1 as reference image, (b) the grayscale of the normalized image, (c) the inverted grayscale image of 4.1(b), (d) the histogram of the inverted image. (e) The binary image of the inverted grayscale image including footprints of erythrocytes, thrombocytes, and parasites.

The normalized image was converted to grayscale image (Figure 4.2(b)) based on luminance for the same reason as described in Section 3.5.2. There is no essential loss of information in converting to grayscale. At this stage, parasites and thrombocytes appear darker than normal erythrocytes, but erythrocytes are much darker than the background. For reasons explained in Section 6.4, the grayscale image was inverted (Figure 4.2(c)). The inverted grayscale image and its histogram (Figure 4.2(d)) clearly show that the foreground

and background have significantly different intensities. Thus, simple thresholding segmentation based on Otsu's method (Section 2.4.2) was applied to the inverted grayscale image to obtain a binary image (Figure 4.2(e)) comprising a foreground of isolated regions representing erythrocyte footprints and a small number of thrombocyte and parasite footprints.



Figure 4.3: Parasite segmentation. (a) The plot of the grayscale intensity profile along the line y =325 in Figure 4.2(c) passing a parasite, and (b) the plot of the grayscale intensity profile along the line y =367 in Figure 4.2(c) passing a thrombocyte.

From the discussion in Section 1.5, the area of the footprint of a thrombocyte is around 6.25% of the area of an erythrocyte and the number of thrombocytes is typically around 5% of the number of erythrocytes. Also considering that the maximum level of parasitemia is 10% (Garcia & Bruckner 1997), infected erythrocytes compared to total erythrocytes, and that the area of the ring-stage parasite is around 6.25% (Section 1.4 and Section 1.5), the maximum of total thrombocyte and parasite area is around 0.94% of the total erythrocyte area. Meanwhile, as discussed in Section 1.9.2, parasites and thrombocytes generally have higher intensities than erythrocytes and the intensities of parasites are normally higher than that of thrombocytes (Figure 4.3(a) and (b)). The parasites themselves do not account for sufficient area to influence the calculation above. This is because the number of parasites does not exceed 10% of the number of erythrocytes (above this level, the disease is fatal) and the area of the high intensity part of the parasite is roughly 6% of the area of an erythrocyte. Hence, the total area contributed by parasites is less than 1% of the area contributed by erythrocytes.

Figure 4.4: Candidate parasite footprints.

Thus, to segment parasites and thrombocytes from erythrocyte footprints, an adaptive threshold was set at 0.94% of the highest foreground intensity from the grayscale image. Figure 4.4 shows the parasite footprint as a result of the adaptive thresholding segmentation. This segmentation by adaptive threshold reliably retains parasites, but may also retain thrombocytes and some dark erythrocytes. Further steps to separate parasites from these other objects are described in subsequent sections.

In thin blood films, the number of leukocytes is much less than that of erythrocytes. Most thin blood film images do not contain a leukocyte. When a leukocyte is present in a thin blood film image, it is generally the only one. However, if a leukocyte exists in an image, the leukocyte will influence the results of parasite segmentation. To remove the existence of leukocytes in an image, a process of segmenting and removing is required. This process is explained in the next section.

## 4.3   Leukocyte Segmentation and Subtraction

In view of the previous section, removing leukocyte footprints from erythrocyte footprints before segmenting parasites is important for obtaining accurate parasitemia estimation. One way of doing this is by segmenting leukocytes and subtracting the leukocyte footprints from erythrocyte footprints. From the discussion on leukocytes in Section 1.9.2 and by inspection of seventy image intensity profiles consisting of leukocytes, Figure 4.5(b), as well as based on the relative intensities of leukocytes and other main blood components in thin blood film images, an adaptive leukocyte threshold value at 0.7 of the highest intensity pixels normally associated with leukocytes was determined empirically to segment leukocytes from other components. This segmentation by the adaptive threshold reliably retains leukocytes, but may also retain parasite and dark erythrocyte footprints in some images without leukocytes.

68

To reduce this error, $(\mu + 3\sigma)$ of maximum erythrocyte intensities was used as the threshold for leukocytes, instead of 0.7 of the highest image intensity.



(a)



(b)

(c)



(d)

(e)

Figure 4.5: Segmenting and removing leukocyte from erythrocyte footprints. (a) Inverted grayscale image of the Figure 1.6(c). A leukocyte at position x = 600 and y = 300. (b) A plot of the profile intensity along the line y = 300. (c) Leukocyte footprint of the image of 4.5(a), (d) erythrocyte footprints of the image of 4.5(a) including leukocyte footprint, and (e) erythrocyte footprints of the image of 4.5(a) without leukocyte footprint.

The leukocyte segmentation and subtraction procedure is illustrated through example images in Figure 4.5. The leukocyte threshold was used to identify the leukocyte footprints (Figure 4.5(b)) and Otsu's method was applied to the inverted grayscale image (Figure

4.5(a)) to obtain the footprints of the erythrocytes and leukocyte (Figure 4.5(c)). In this example, the leukocyte footprint was relatively free from noise. In other images, some noise or parasite footprints accompanied leukocyte footprints. To remove the noise and residual parasite footprints from leukocyte footprints, an opening filter (Section 2.3.2) with circular structure element radius $R = 6$ was applied. After that, the leukocyte footprint was subtracted from the erythrocyte footprint to remove the leukocyte footprint. The Figure 4.5(d) shows the erythrocyte footprints after removal of the leukocyte footprint. In this case, the leukocyte footprint was removed fully because the thresholding process resulted in a footprint of the leukocyte that closely matched the size of the leukocyte footprint in the erythrocyte binary images.

However, in other images, applying the leukocyte threshold resulted in leukocyte footprints that were smaller than those in erythrocyte footprints. In such cases, after the subtraction process, some leukocyte footprint borders remained. To avoid this, dilation with the same structuring element as the opening filter was applied to leukocyte footprints before the subtraction process.

According to the nature of a leukocyte (Section 1.5), leukocytes will engulf parasites. Thus, the probability of a parasite being located very close to a leukocyte is small and so the dilation process will not affect counting the number of parasites very much. Despite the dilation process, a few thin leukocyte footprint borders were present in some cases. Generally, these remaining borders were very small compared to thrombocyte footprints captured in erythrocyte footprints. Accordingly, these remaining borders did not significantly affect the estimate that the total thrombocyte and parasite (if present) area is around 0.94% of the total erythrocyte area (Section 4.2). However, these remaining leukocyte borders might be detected as parasites. This will mislead the parasitemia count. To distinguish the remaining borders from parasites, the parasite classification discussed in Section 4.4 was used.

Aside from normal leukocytes, infected thin blood film images may contain phagocytes or schizonts (Section 1.5 and Figure 1.6(d)). Therefore, leukocyte footprints might also contain schizonts footprints. To distinguish schizonts from leukocyte, leukocyte classification (Section 6.5) was applied to the leukocyte footprints.

## 4.4   Parasite Classification

To compute parasitemia in thin blood films, parasites must be identified and counted. The segmentation by adaptive threshold described above reliably retained parasites, but also retained some thrombocytes, dark erythrocytes, and leukocyte borders. To distinguish parasites from these other objects, the morphology of candidate parasites and the histograms of green values within the candidate parasite footprints were considered.

### 4.4.1   Morphological Feature Analysis to Recognise Parasite

This section describes data set, experimental details and results for training of parasite classification.

#### 4.4.1.1   Images for training of parasite classification

The data set involved in this part of the study comprised a total of 96 images from the seven thin blood films described in Section 1.11. Each image contained parasites, thrombocytes, dark erythrocytes, or leukocytes.

#### 4.4.1.2   Experimental Details

In the pre-processing step, the selected images were normalized by applying the colour normalization described in Section 6.3. Subsequently, the normalized images were converted to grayscale images and threshold segmentation (Section 4.2) was applied to get binary images containing parasite footprints, thrombocyte footprints, dark erythrocyte footprints, and leukocytes. Here, 140 ROIs were extracted from the training images. Of these, 35 were parasite footprints, 35 were thrombocyte footprints, 35 were dark erythrocyte footprints, and 35 were leukocyte border footprints. Since parasites, thrombocytes, dark erythrocytes, and leukocytes were not always present in the same images, some of these ROIs were from the same images and others were from different images.

In this experiment, the morphology values of the ROIs were extracted. The ellipticity and eccentricity (Section 2.6.2) for each selected ROI corresponding to parasite, thrombocyte, dark erythrocyte, or leukocyte border were calculated. Quadratic discriminant analysis, a variant of discriminant analysis (Section 2.9.1), was used to classify the candidate parasites as true parasites or not. Figure 4.6 shows an example of the process of parasite classification and the results.

Figure 4.6: Parasite identification. (a) Erythrocyte and candidate parasite footprints. Blue pluses indicate the centers of full single erythrocytes. Green pluses indicate the erythrocyte centers at occluded erythrocytes or border erythrocytes. Magenta dots indicate the centers of candidate parasites. (b) Parasite classification based on its location. Red and green stars indicate parasites and non-parasites, respectively.

### 4.4.1.3   Results, Discussion and Conclusion

Figure 4.7 shows a scatter plot of the morphology values of the ROIs and confusion matrices (Table 4.1 and 4.2) reports the classification results of ROIs obtained using the discriminant analysis. According to the confusion matrices and Figure 4.7, the ellipticity and eccentricity parameters of parasites were not significantly different from those of the other objects. Thus, the ellipticity and eccentricity were not reasonable as features for parasite classification. Therefore, these features were not included in classifying parasites in the remainder of the study.

Table 4.1: Confusion matrix obtained from parasite classification using quadratic discriminant analysis based on ellipticity and eccentricity of ROIs. The classification accuracy is 45%.

|  |  | Predicted | | | |
|---|---|---|---|---|---|
|  |  | Parasite | Thrombocyte | Dark Erythrocyte | Leukocyte Border |
| Actual class | Parasite | 7 | 16 | 2 | 10 |
|  | Thrombocyte | 3 | 31 | 1 | 0 |
|  | Dark Erythrocyte | 5 | 5 | 6 | 1 |
|  | Leukocyte Border | 10 | 5 | 19 | 19 |

Table 4.2: Confusion matrix of ROI (as in Table 4.1) classification based on parasite or non-parasite. The classification accuracy is 67.14%. The sensitivity and specificity are 20% and 82.86%, respectively.

|  |  | Predicted | |
|---|---|---|---|
|  |  | Parasite | Non-Parasite |
| Actual class | Parasite | 7 | 28 |
|  | Non-Parasite | 18 | 87 |

72

Figure 4.7: Ellipticity and eccentricity values of ROIs. Red asterisks are parasites, green diamonds are thrombocytes, magenta squares are dark erythrocytes, and blue circles are leukocyte borders.

### 4.4.2 Colour Intensity Feature Analysis to Recognise Parasite

An alternative method for parasite classification was carried out based on colour intensity features as parameters. According to Figure 4.8(b), (c), and (d), the intensity of the green channel dominates the colour image intensity profile, especially for parasites, thrombocytes, and erythrocytes. Visually, the profiles of the green channels of parasites are symmetric and spiky. Meanwhile, the profiles of thrombocytes are also symmetric but the curves are wider than that of parasites. In addition, the green channels of dark erythrocytes show characteristic plateaus. This means that the green channels of parasites, thrombocytes, and dark erythrocytes are distinguishable. Thus, the statistical descriptors of the green channels were used to classify parasites, thrombocytes, dark erythrocytes, or leukocytes. In particular, variance, skewness, and kurtosis were used.

#### 4.4.2.1 Data and Experimental Details for Colour Feature Analysis

A set of 140 training ROIs were selected from binary images computed as part of the pre-processing step explained in Section 4.4.1.2 was used in this colour feature analysis. For each ROI, the green channel was extracted from the corresponding normalized colour image (Section 6.3) and the ROI was represented by the variance, skewness and kurtosis of its green channel values. Discriminant analysis (Section 2.11.1) was used to classify the ROI as representing parasites or not.

73

Figure 4.8: Colour profiles: (a) a plot of the colour intensity profile along the line y = 325 in Figure 4.2(a). There is a parasite at approximately x = 160 and erythrocytes at around x = 75-160, 170-225, 450-500, and 550-600. (b) A plot of the colour intensity profile along with the line y = 367 in Figure 4.2(a). A thrombocyte exists at around x = 350-375. (c) A plot of the colour intensity profile along the line y = 225 in Figure 4.2(a). There is a dark erythrocyte at around x = 225-275.

### 4.4.2.2    Results of Training, Discussion and Conclusion

The green channel values of the ROIs are shown in Figure 4.9 and classification performance based on green colour feature for distinguishing parasites, thrombocytes, dark erythrocytes, or leukocytes are displayed as confusion matrices (Table 4.3). According to Table 4.3, the classification accuracy was 78.57%.

Table 4.3: Confusion matrix of classification based on green colour channel feature for parasite, thrombocyte, dark erythrocyte, or leukocyte border.

| | | Predicted | | | |
|---|---|---|---|---|---|
| | | Parasite | Thrombocyte | Dark Erythrocyte | Leukocyte Border |
| Actual class | Parasite | 22 | 0 | 0 | 8 |
| | Thrombocyte | 8 | 33 | 1 | 1 |
| | Dark Erythrocyte | 0 | 1 | 33 | 0 |
| | Leukocyte Border | 5 | 5 | 1 | 22 |

74

However, the main purpose of this part of this study was to determine a pattern for classifying ROIs as parasites or non-parasites. When viewed as a classification into two classes, parasites and non-parasites (thrombocyte, dark erythrocyte, and leukocyte border), the classification accuracy is 85% (Table 4.4). The sensitivity and specificity are 73.33% and 88.18%, respectively.



Figure 4.9: Green channel values of ROIs. Red asterisks are parasites, green diamonds are thrombocytes, magenta squares are dark erythrocytes, and blue circles are leukocyte borders. X-axis is standard deviation, Y-axis is skewness, and Z-axis is kurtosis. Scatter plots for the red and blue channels are very similar to the green channel and are not shown separately.

Table 4.4: Confusion matrix of classification as parasite or non-parasite.

| | | Predicted | |
|---|---|---|---|
| | | Parasite | Non-Parasite |
| Actual class | Parasite | 22 | 8 |
| | Non-Parasite | 13 | 97 |

This classification experiment confirmed that colour intensity features, in particular, statistical descriptors of the green channel, are able to distinguish ROIs with parasites from ROIs without parasites. Subsequently, this classification experiment will be referred to the parasite recognition.

## 4.5   Identifying Infected Erythrocytes

The purpose of this section is to identify infected erythrocytes by distinguishing parasites infecting erythrocytes from independent parasites and thrombocytes. In Section 4.2, parasites were segmented by the threshold method. In some cases, the parasite segmentation also produced footprints of thrombocytes. There was also the possibility that a single erythrocyte contained multiple footprints, which may imply the presence of either several parasites in one cell or a single parasite detected as more than one footprint. However, since the use of thin blood films in this study is not to calculate the number of parasites but to identify infected erythrocytes, only one parasite footprint in each erythrocyte is recognised as a representation of the cluster of footprints in the cell. The representative footprint was selected through the clustering method explained in Section 3.3.2.2.

Two consecutive procedures were used to distinguish parasites from thrombocytes. The first procedure used the fact that thrombocytes appear outside of erythrocytes, as explained in Section 1.9.2, as the feature. If (Xe,Ye) represents the centroid of the erythrocyte footprint, (Xp,Yp) represents the centroid of the candidate parasite footprint and $R$ is the radius of the erythrocyte footprint, then if $|Xp_i-Xe_i|>R$ AND $|Yp_i-Ye_i|> R$, the footprints are located outside erythrocyte and thus can be confirmed as thrombocytes or independent parasites (ignored in calculation of infected erythrocytes); otherwise, the footprints may be either thrombocytes touching an erythrocyte or true parasites. This indicates that the location of thrombocytes by itself, despite being an essential feature, is not sufficient to differentiate parasites from thrombocytes.

To resolve these two possibilities, these first steps were followed up with a second procedure utilising colour features to distinguish between parasites and thrombocytes. In this procedure, the green channel of each candidate parasite footprint was extracted from the associated normalized colour image and the variance, skewness and kurtosis of the green channel were extracted. These three green feature values were then classified by means of the pattern discovered in Section 4.4.2 as true parasites or non-parasites. In this case, the parasite candidates classified as true parasites represent infected erythrocytes. The combination of erythrocyte detection, location feature, and parasite identification method described above will be referred as the erythrocyte identification method.

### 4.5.1    Images for Testing the Erythrocyte Identification

For validation of the performance of the erythrocyte identification method described in the previous section, a data subset of around 1000 images consisting of 75075 erythrocytes from the testing data set (Section 4.6.1) were selected randomly. Around 100 images were collected from each thin blood film.

### 4.5.2    Results of Testing, Discussion and Conclusion

The experiment of the erythrocyte identification yielded a confusion matrix (Table 4.5). The accuracy was 99.84%, the sensitivity was 73.88%, and the specificity was 99.87%.

Table 4.5. Confusion matrix of testing classification for infected or normal erythrocytes.

|  |  | Predicted | |
| --- | --- | --- | --- |
|  |  | Infected | Normal |
| Actual class | Infected | 266 | 94 |
|  | Normal | 32 | 74683 |

In the context of classification performance, the specificity of the erythrocyte identification was similar to that of the parasite recognition (Section 4.4.2). However, the accuracy and specificity of the erythrocyte identification were significantly different from those of the parasite recognition. According to Table 4.3, the big contributor of misclassification in the parasite recognition is *FP* due to thrombocytes and leukocyte borders detected as parasites. The location feature used to filter thrombocytes and leukocyte borders from parasites in erythrocytes had a role in reducing the misclassification. Thrombocytes and leukocyte borders, which are far away from erythrocytes, are not involved in the classification process of the erythrocyte identification. This means that the location feature is combined with colour intensity pattern to classify erythrocytes in the erythrocyte identification. This significantly reduces the effect of thrombocytes and leukocyte borders. Thus, the accuracy and the specificity of the erythrocyte identification method are much higher than those of the parasite recognition method using only the colour intensity feature to distinguish between parasites and non-parasites.

In addition, in the erythrocyte identification, the number of normal erythrocytes was much higher than the number of infected erythrocytes. The maximum proportion of infected erythrocytes is 10% (Garcia & Bruckner 1997). Statistically, the probability of *TN* in the erythrocyte identification is much higher than that in the parasite recognition. Furthermore,

most normal erythrocytes have a lower intensity than the parasite adaptive threshold. This means that the normal erythrocytes were automatically assigned as true normal erythrocytes (*TN*) without being involved in the parasite classification process. Accordingly, the specificity and accuracy of the erythrocyte identification are significantly higher than those of the parasite recognition.

For classification, Table 4.5 indicates that the classification performance of the proposed algorithm performs at least as well as that of Tek's study (Tek, F. B., Dempster & Kale 2010) in three categories: accuracy 93.3%, sensitivity 72.4%, and specificity 97.6%. For specificity, the proposed algorithm slightly outperforms or is relatively similar to Tek's study. However, for accuracy and specificity, the proposed algorithm performs substantially better. This might be because the proportion of normal erythrocytes in Tek's data set was much less than in this study.

The erythrocyte identification method was applied to testing data (Section 4.6.1) for parasitemia estimation, as described in the following section. The outcome of the diagnosis was then compared with that of parasitemia scores by expert readers.

## 4.6   Estimating Parasitemia

The parasitemia may be estimated by noting the number of infected erythrocytes as a percentage of the total number of erythrocytes and using Equation 1.2 to give the parasite count per µl of blood.

### 4.6.1   Images for Testing Parasitemia Estimation

From the seven positive thin blood films and three negative thin blood films described in Section 1.11, a data set of 1590 images from the seven positive thin blood films (the same as data set used in Chapter 3) and 900 images from the three negative thin blood films were randomly selected (between 200 and 300 images from each slide) for testing parasitemia estimation.

### 4.6.2   Testing for Parasitemia Estimation

Every candidate parasite footprint was identified and the location feature (Section 4.5) was applied to the testing data set (Section 4.6.1). If the candidate parasite was detected inside an erythrocyte, then the colour pattern method (Section 4.5) was applied to classify the candidate parasite as parasite or non-parasite.

### 4.6.3 Results of Testing, Discussion, and Conclusion

In this chapter, three steps for segmentation and identification were used to estimate parasitemia in thin blood film: the first step was erythrocyte segmentation followed by identification (Section 3.4); the second step was parasite segmentation followed by classification; and, the final step was parasitemia estimation using the standard correlation (Equation (1.2)). From the first and second step, the percentage of infected erythrocytes was obtained. The percentage of infected erythrocytes resulting from the testing step is described in Table 4.6 and Table 4.7. In addition, the resulting parasitemia estimation from the testing step is shown in Figure 4.10 and Figure 4.11.

Regression analysis was used to determine if there was a correlation between the percentages of infected erythrocytes resulting from this algorithm in thin blood films and parasitemia scores from expert readers in the corresponding thick blood films. The correlation result is shown in Figure 4.10. The result indicates that there is good correlation, $r = 0.91$. Furthermore, the estimates of parasitemia and parasitemia scores from expert readers have a strong correlation, $r = 0.98$.



(a)                              (b)

Figure 4.10: (a) The regression for testing images. For each slide, there are green open circles representing estimates of parasiatemia from human experts in thick blood films. The correlation between percentage of infected erythrocytes for testing images in positive thin blood films and parasitemia scores from human experts in thick blood film was 0.91. (b) The estimation of parasitemia in thin blood films and parasitemia scores from expert readers in thick blood films. The open green circles are as in Figure 4.10(a) and the blue stars are the estimates of parasitemia from the percentage of infected erythrocytes for testing images resulted from this algorithm after correlation with parasitemia scores by standard correlation (Equation 1.2).

Figure 4.11: Estimates of parasitemia scores. (a) The green vertical lines are ranges of parasitemia scores from expert readers in thick blood films after converting to thin blood film scores using DS correlation (Equation 3.1). (b) The magenta vertical lines are ranges of parasitemia scores from human experts in thick blood films after converting to thin blood film scores using TR correlation (Equation 3.2).

Table 4.6: Percentage of infected erythrocytes in negative thin blood films. The mean was 0.004% and standard deviation was 0.0001%. *NS1, NS2,* and *NS3* are negative slide 1, 2, and 3.

| Slides | *NS1* | *NS2* | *NS3* |
|--------|--------|--------|--------|
| *B* | 0.0040 | 0.0039 | 0.0041 |
| *P* | 200 | 193 | 205 |

Table 4.7: Percentage of infected erythrocytes in positive thin blood films (*B*) and their estimation of parasitemia scores (*P*) based on the standard correlation. *PS1, PS2, PS3, ..., PS7* are positive slides 1, 2, 3, …, 7.

| Slides | *PS1* | *PS2* | *PS3* | *PS4* | *PS5* | *PS6* | *PS7* |
|--------|--------|--------|--------|--------|--------|--------|--------|
| *B* | 0.0047 | 0.0128 | 0.0375 | 0.1262 | 0.3310 | 0.6951 | 2.2291 |
| *P* | 234 | 639 | 1,875 | 6,309 | 15,372 | 34,754 | 111,455 |

Subsequently, a t-test was used to determine if the difference in estimation performance of parasitemia estimates due to this algorithm compared with parasitemia scores from expert readers. The results indicate that there is no significant difference between the performance ($p = 0.37$, $n = 7$, $\alpha = 0.05$) and the mean of experts' parasitemia scores. Similarly, the parasitemia resulting from this algorithm was found to be not significantly different ($p = 0.40$, $n = 7$, $\alpha = 0.05$) from the median of experts' parasitemia scores.

For the negative thin blood films, zero parasitemia scores were not obtained (Table 4.6). This may be caused by the presence of inaccuracy in the process of parasite

classification (Table 4.3 and Table 4.4). This could be addressed in future work, given that thrombocytes touching normal erythrocytes were not considered in the training process. Classification inaccuracy is due to false positives (*FP*) in distinguishing between thrombocytes and parasites. Some thrombocytes located near (less than the average radius of erythrocytes) or on erythrocytes were analysed as candidate parasite footprints. The number of thrombocytes is from around 5% of erythrocytes (Vander, Luciano & Sherman 2001) up to 8% of erythrocytes. Accordingly, this proportion was big enough to play role in the over estimation of parasitemia and caused less sensitivity, 73.88% (Table 4.5). On the other hand, some bright parasites, probably young parasites, were mostly not detected because their intensities on inverted grayscale images are lower or not significantly different from normal erythrocytes. This produced false negative misclassification and reduced sensitivity. However, if the estimates of parasitemia scores from the negative thin blood films were joined with the scores from the positive thin blood films, the performance of the parasitemia estimation was not still significantly different (p = 0.54, n = 10, α = 0.05) from the parasitemia scores from expert readers.

The algorithms presented in this chapter did provide better estimation of parasitemia than that in (Sunarko et al. 2017). Generally, results from this algorithm indicate that estimates of parasitemia were in line with parasitemia scores from expert readers in thick blood films both directly (Figure 4.10(b)) and in terms of DS correlation (Figure 4.11(a)) or TR correlation (Figure 4.11(b)). However, the parasitemia value resulting from this algorithm in the lowest parasitemia slide was among the outliers.

Comparisons between the study in (Sunarko et al. 2017) and this study must be made judiciously. The study in (Sunarko et al. 2017) was focused on distinguishing infected erythrocytes from uninfected erythrocytes based on erythrocyte analysis. Here, infected erythrocytes were determined based on not only erythrocyte identification but also parasite analysis. Parasites were segmented and analysed independently from their erythrocytes. Therefore, false positives due to the effects of infected erythrocytes or thrombocytes on the other normal occluded erythrocytes could be avoided and thus reduce the error.

In (Sunarko et al. 2017), erythrocyte classification was based on erythrocyte features. In some cases of occluded erythrocytes, a histogram computed over all occluded erythrocytes represented not only individual erythrocytes but also the other erythrocytes or thrombocytes joined in the occluded cells. This leads to miss-interpretation in erythrocyte classification. Meanwhile, parasite features were strictly based on a local parasite candidate.

81

Hence, a histogram represented an individual parasite candidate. In addition, here, the erythrocyte identification based on parasite features was coupled with location feature and clustering analysis. Before being classified as a parasite infecting an erythrocyte, a parasite candidate was determined by location feature whether the candidate was inside or outside of an erythrocyte. If the candidate was outside, it was ignored because this step consists of identifying infected erythrocytes, not noting the parasites themselves. Subsequently, the parasite candidate was classified based on parasite features as parasite or non-parasite. In cases where there are many parasites infecting an erythrocyte, a simple clustering algorithm (Section 3.3.2.2) was applied to select one of them and ignore the others.

The method can estimate parasitemia via automatic analysis of thin blood film images with error no greater than the variation between expert malaria readers. The important observation in this chapter is that the combination of erythrocyte identification and parasite analysis outperforms features from erythrocyte analysis alone. This significantly reduced the overestimation of parasitemia obtained in (Sunarko et al. 2017).

# Chapter 5:   Morphological Features

As discussed in Section 1.6, both thick and thin blood films are valuable and microscopists are highly recommended to look at both when conducting a malaria examination (Moody 2002). In this chapter and the next, computational methods are developed to automatically diagnose and estimate parasitemia based on thick blood film images.

Thick blood film images (Section 1.9.3) contain three main components: leukocytes, thrombocytes, and cytoplasm, and possibly parasites if the images are from positive slides. Parasites contain one or two dot chromatins. A source of error in malaria diagnosis based on thick blood film images is in distinguishing between parasites and thrombocytes. The probable reason is that thrombocytes also contain chromatin (Elter, Haßlmeyer & Zerfaß 2011) or chromatin-like material (Howard & Hamilton 2013; Thein 2001). Human readers use physical appearance of parasites and thrombocytes as well as colour and size in malaria diagnosis.

This preliminary study to the main project of automatic parasitemia estimation based on thick blood film images was conducted to ascertain if physical appearance and size can be used to distinguish parasites from thrombocytes. The most commonly seen stage of parasites is the trophozoite stage which, in thick blood films, is characterised by the presence of a red chromatin dot together with blue cytoplasm of lower intensity in the shape of a ring, a partial ring or more irregular shapes called "amoeboid" (WHO 2010). In *P. falciparum*, a second chromatin dot is often seen (WHO 2010).

Two studies were undertaken to determine the extent to which these features could be used to distinguish between parasites and thrombocytes in thick blood films. In the first study, the size of secondary structure in the neighbourhood of each candidate parasite (if present) was recorded and the size (total intensity) of the secondary feature was recorded as a feature for classification. This method will be referred to as the "second structure method". In the second method, the presence of a second dot was used to distinguish between parasites and thrombocytes. The motivation for this approach was that the first dot identified could represent either a chromatin dot associated with a parasite or a thrombocyte. The presence of a second dot would suggest a second chromatin dot and hence a *P. falciparum* parasite. All seven malaria positive slides in the data set were known to be *P. falciparum* infections. This method will be referred to as the "chromatin dot method".

For this study, 80 images from each slide were used to establish feasibility. The following pre-processing steps were common to both methods.

Colour images were converted to grayscale images based on luminance and templates for leukocytes were constructed for each image as in Section 6.5. The number of leukocytes present in each image determined by the number of connected components in the leukocyte template was recorded.

## 5.1   Stripe Artefact Removal

When small patches of images were observed (tantamount to zooming), images from some slides were found to show significant horizontal striping patterns and, in some cases vertical striping to a lesser degree. When present, horizontal stripes consisted of alternating rows of light and dark stripes of single pixel width but extending over the full width image. Vertical stripes, when present, tended to appear in blocks of size $8 \times 8$ but did not necessarily follow a pattern of alternating light and dark stripes (Figure 5.1(a)).

A program was written to mitigate these stripe artefacts. To remove the horizontal striping pattern, the intensity sum over all odd numbered horizontal rows of pixels was computed (odd row sum) and the intensity sum over all the even numbered rows of pixels was computed (even row sum). A constant, $C$ was added to every pixel in the odd rows to insure that the new odd and even row sums would be equal. Thus the constant $C$ was given by

$$C = \frac{\text{even row sum - odd row sum}}{N/2}$$

Here $N$ is the total number of pixels in the image. This strategy resulted in reducing (albeit not totally eliminating) the apparent horizontal striping (Figure 5.1(b)).

After adjusting for the horizontal striping artefacts, the vertical stripes were addressed for each $8 \times 8$ block of pixels at a time according to the following rule. If, within and $8 \times 8$ block, the intensity sum of a column was less than the column sums of each of the columns on either side, then a constant was added to all the pixels in the column so that the new column sum would be equal to the mean of the column sums of the two columns on either side. This strategy resulted in reducing the vertical striping (Figure 5.1(c)).

Figure 5.1: Removing stripe pattern – local patch. A zoomed view of an image patch consisting of 48 rows and columns. (a) is the original patch. (b) is the same patch after removing horizontal stripes. (c) is the same patch after removing horizontal and vertical stripes and (d) is the patch after subsequent local smoothing.

Next, a local smoothing filter was used to remove the remnant striping patterns. This filter consisted of a 3×3 patch of pixels with intensity value 1/9. The striping patterns in the resulting image were noticeably reduced (Figure 5.1(d)). Applying a smoothing filter alone did not reduce the striping patterns nearly as well.

## 5.2   Finding Candidate Parasites

A filter, h, was used to search for high intensity dots associated either with parasites (a chromatin dot) or thrombocytes (Figure 5.2(a) and (b)). The filter consisted of an image patch with intensity values given by

$$h(x,y) = \begin{cases} -A, & if \quad R^2 < x^2 + y^2 \leq (R+1)^2 \\ 1, & if \quad x = 0, y = 0 \\ 0, & if \quad otherwise \end{cases}$$

where the number $A$ was chosen so that $\sum_{x,y} h(x,y) = 0$. Here, R is the radius of the circle on which the filter is positive. The motivation for this filter is that, being zero sum, the filter

effectively sets the background to have mean zero. Since $h(0, 0) = 1$ and $h(x, y) = 0$ for points $(x, y)$ near $(0, 0)$, the filter locally acts as the identity. Thus, this filter preserves local intensity structure well but flattens the image to have mean zero background.



Figure 5.2: Candidate parasites. (a) An example grayscale image. (b) The same image after flattening by filter h. (c) After applying the intensity threshold $T = 20$. (d) After applying the leukocyte template to remove remaining edge effects from leukocytes.

The filtered image was inverted (so that foreground features were represented by positive intensity values) and thresholded at zero to remove background fluctuations. Several intensity profiles were examined to determine that an intensity threshold of 20 was suitable for identifying parasites and thrombocytes (Figure 5.2(c)). The resulting binary image will be referred to as the parasite-thrombocyte footprint image.

Since the spatial extent of the flattening filter was very small compared to the area of a leukocyte, the interior of regions corresponding to leukocytes were automatically set to background values. However, the edges of leukocytes gave rise to high responses to the filter as would be the case for any linear zero-sum filter. Thus, edges of leukocytes were present in the parasite-thrombocyte footprint image. To remove these edge responses, the parasite-thrombocyte footprint image was multiplied element-wise by the leukocyte template. Finally, the parasite-thrombocyte footprint images also contained a number of small specks,

too small to be considered candidate parasites. Accordingly, all connected components of area 5 pixels or less in the parasite-thrombocyte footprint image were removed (Figure 5.2(d)).

The remaining connected components in the parasite-thrombocyte footprint image were viewed as associated either with parasites or non-parasite objects, including thrombocytes and remaining artefacts. Since only the number of true parasites was of interest, attention focused on classifying these components as representing parasites or not.



Figure 5.3: Second structure method. (a) The bright spot at the centre is the intensity image associated with a connected component *C* in the parasite-thrombocyte footprint image and the image shown includes the disk *D* extracted from the flattened image. The lower bright spot is not part of *C* but is part of *D*. (b) After removing *C*. (c) After identifying the largest component in *D* after removing *C*. (d) The intensity image associated with the second largest structure near *C*.

## 5.3 The Second Structure Method

For each connected component in the parasite-thrombocyte footprint image, several processing steps were taken. Let *C* denote a connected component in the parasite-thrombocyte footprint image and let (*x0, y0*) denote the centre of this set. A disk shaped region, *D*, of radius 15 pixels centred at (*x0, y0*) was extracted from the flattened image (the grayscale image after flattening but prior to applying the threshold) (Figure 5.3(a)). The intensity function in *D* corresponding to the connected component *C* was set to zero so that

the only pixels with positive intensity values in $D$ corresponded to objects in the image other than the candidate parasite associated with $C$ (Figure 5.3(b)). A new intensity threshold ($T = 8$) was applied to $D$ resulting in a binary patch, $P$, the size of $D$ and comprising footprints of structures near the candidate parasite associated with $C$, but excluding this candidate itself (Figure 5.3(c)). The total intensity (integral of the intensity function) in the flattened image over each of the connected components in $P$ was computed and the component with the largest such total intensity was taken to represent the most prominent secondary structure near $C$ (Figure 5.3(d)). The total intensity over this component was recorded as the size, s, of the most prominent secondary structure near $C$. Since the original (primary) structure $C$ was removed from the parasite-thrombocyte footprint image, this structure was not eligible to play the role of a second structure for a different candidate parasite and hence a single pair of dots was only considered once as one potential parasite.

For a given threshold $t$, the component $C$ was designated as a parasite if $s \geq T$ and as a non-parasite if $s < T$. For each slide, a parasitemia score was obtained by multiplying the number of parasites found in the 80 sample images from this slide by $8000/L$, where $L$ is the number of leukocytes found within the 80 images. The number 8000 comes from the standard value for the number of leukocytes per µl of blood. The parasitemia scores produced in this manor were compared to the average parasitemia scores reported by human experts (Figure 5.5).

## 5.4   The Chromatin Dot Method

Starting with the parasite-thrombocyte footprint image described in Section 1.2, all connected components of area less than 65 pixels were considered candidate chromatin dots. If two candidate chromatin dots were found with centres less than 25 pixels apart, the two were considered to be pairs of chromatin dots and counted as a single parasite. Connected components without a neighbour within 25 pixels distance were considered isolated dots not associated with parasites (Figure 6.4). Parasitemia was computed by the same method as described in Section 1.3 using the leukocyte count derived from the leukocyte template.

Figure 5.4: Parasite detection based on pairs of chromatin dots. The image is the same example as in Fig. 5.2(a). Red circles indicate detections of chromatin dot pairs associated with parasites and green circles indicate isolated dots rejected as being associated with parasites. The mean human expert parasitemia score for the slide from which this image came was more than 300,000.

## 5.5 Results

Neither method for estimating parasitemia agreed well with the parasitemia scores reported by expert readers (Figures 5.5 and 5.6).



Figure 5.5: Expert and second structure parasitemia. For each slide, the log of the parasitemia score estimated using the second structure method is plotted according the horizontal axis and the log of the human expert score is plotted according to the vertical axis. Each slide is represented by several blue circles, one for each human expert. The red line is the regression line. The results are shown for the threshold on the size of the second structure yielding the lowest error. Logarithms were used because the parasitemia values ranged over five orders of magnitude.

For the second structure method, the estimate of parasitemia depends on the threshold used for the size (total intensity) of the second structure. A search over threshold values

89

within the range of second structure sizes was conducted to find the best fit according to the sum mean square error between estimates of parasitemia for each slide and the mean of the human expert estimates of parasitemia. The fit associated with the threshold yielding the least error by this method indicates that the estimates based on the second structure method does not agree well with human expert estimates (Figs. 5.5). For the chromatin dot method, there was no parameter over which to optimize the fit (Figure 5.6).



Figure 5.6: Expert and chromatin dot parasitemia. For each slide, the log of the parasitemia score estimated using the chromatin dot method is plotted according the horizontal axis and the log of the human expert score is plotted according to the vertical axis. Each slide is represented by several blue circles, one for each human expert. The red line is the regression line. Logarithms were used because the parasitemia values ranged over five orders of magnitude.

## 5.6 Discussion and Conclusion

The two methods presented in this chapter were designed to exploit visual clues used by human experts to distinguish between parasites and non-parasite objects of similar size and image intensity, including thrombocytes. The approaches were both aimed specifically at detecting *P. falciparum* in the trophozoite stage and were simple in that shapes of structures were not carefully analysed. It was hoped that such an approach would be flexible and allow the detection of all the variations of the ring shapes and the presence of a second chromatin dot. Only results on a training set of 80 images per slide were presented, because results were so poor that there was little point in measuring performance on a testing set.

The poor results indicate that much more care is needed to characterise the shapes of structures rather than measuring just their total image intensity. Probably, modern image

analysis methods could be used to extract shape features from candidate parasites that would result in accurate detection of *P. falciparum* in the trophozoite stage. However, this study indicates that such feature extraction methods would likely be somewhat complex and that a large suite of such algorithms would be needed to identify the various species of malaria parasites and their various stages. In all, the processing time required per image would increase substantially.

Long processing times for images are prohibitive in a study that aims at the long-term goal of developing a simple, cost effective and accurate device that could be used in under developed regions of the world to diagnose malaria. In order to separate malaria negative subjects from low parasitemia malaria positive subjects, hundreds of images from a thick blood film are needed. Thus, a modest increase in processing time per image easily results in unrealistic total processing time for the slide.

Accordingly, this study supports the notion that an algorithm that imitates the human expert's process for detection of parasites may not be optimal for an automatic system. Instead, an alternative study involving intensity and colour features was conducted (Chapter 6) to improve on these results.

# Chapter 6: Estimating Parasitemia from Thick Blood Film Images

## 6.1 Overview

This chapter presents the identification of leukocytes and malaria parasites for the purpose of estimating parasitemia scores based on thick blood films. This chapter starts with describing the data used in this study in Section 6.2. Methods for image normalisation are demonstrated in Section 6.3. Section 6.4 and 6.6 explain the process of segmenting leukocytes and parasites, respectively. The classification of leukocytes and parasite footprints are explored in Section 6.5 and 6.7, respectively. Section 6.8 discusses counting parasites and computing parasitemia scores. Results for estimating parasitemia are discussed in Section 6.9. Finally, Section 6.10 reviews this chapter.

## 6.2 Thick Blood Film Images

Images for this part of the study were taken from thick blood films. A total of 171 images were selected from seven positive slides (Section 1.11) for the training steps (Section 6.3.1, 6.5.1, and 6.7.1) and a total of 2424 images were selected from ten slides consisting of three negative slides and seven positive slides (Section 1.11) for testing parasitemia estimation (Section 6.9). Between 200 and 300 images were randomly selected from each slide.

Typical thick blood film images contain leukocytes, thrombocytes, and, in infected thick blood films, parasites. The original images were RGB colour images having different intensity and colour backgrounds. These differences might occur due to variation in staining or image acquisition. To compensate for these differences, the original images were normalized as described in the following section.

## 6.3 Image Normalisation

Two normalization methods were tested to explore image normalization effects on consistency of image intensity: grey-world normalization and median normalization. These two normalization algorithms were used for the following reasons. Qualitative experiments using the first method revealed that the method was very effective and robust against the different input image colour characteristics (Tek, F B, Dempster & Kale 2006). On the other hand, the second normalization method was simpler but applies only to grayscale images.

### 6.3.1 Images for Normalization Experiments

Sixty three images from seven slides of thick blood films (Section 6.2) were involved in this experimental study. The primary objective of this normalization study was to get consistent images that fit all slides. As such, this data set comprised nine images from each slide.

### 6.3.2 Experimental Details

The first normalization method was applied to colour images. In this colour normalization study, the grey-world algorithm based on a reference image (Figure 4.1) described in Section 2.2.1 was implemented separately on each colour channel to establish colour consistency over the various slides. The algorithm described in Section 2.2.1 used the mean value of image intensity to generate illumination factors. However, in this study, the median was used instead of the mean because the median is more robust than the mean value in representing the population (Moore, McCabe & Craig 2012). For this reason, the illumination factors were sought by dividing the medians of each channel of the reference $med[i]^c$ by those of unknown $med[i]^u$, yielding the following equation:

$$m[i] = \frac{med[i]^c}{med[i]^u} \qquad i = [r,\ g,\ b] \tag{6.1}$$

Otherwise, the grey-world normalisation followed the method described in Section 2.2.1. The grey-world normalisation procedure is demonstrated by example images shown in Figure 6.1. The histograms of the normalised images are shifted (see Figure 6.1(g) and (h)) and aligned to the reference histogram (see Figure 6.2(b)). In addition, although the original colour RGB images were normalised, in this case, small differences in background intensity between images still remained. This may due to the effect of staining resulting in different tissues appear as having roughly the same colour but at different intensities. Even though the blood film images produced by modern microscopes coupled with dedicated cameras are colour images, the distribution of colour forms essentially a one-dimensional subset of the colour space (Figure 3.4). The colours are not fully distributed in all directions of the colour space. Thus, the colour space is essentially one-dimensional.

Figure 6.1: Image normalisation: (a) and (b) are original images, (c) and (d) are normalised images of (a) and (b), (e) and (f) are histogram of original images of (a) and (b) with median 148 and 195, respectively, (g) and (h) are histogram of normalised images of (c) and (d) with median 169 and 183, respectively.

Accordingly, image normalization can be carried out on grayscale images. After converting the input colour images to grayscale images, the median normalisation (Section 2.2.2) was used to replace the grayscale image by Equation (2.1). In this case, α was set to the *Nth* least value, where *N* is 10% of total number of pixels.

### 6.3.3    Results, Discussion, and Conclusion

To confirm the usefulness of the image normalization, qualitative and quantitative experiments were conducted. A set of images from different slides with different intensities (Section 6.3.1) was used to compare outputs.

Table 6.1: Median of image intensity. *G1* is the grayscale of the original images. *G2* is the grayscale of the normalized RGB images. *G3* is the grayscale of the normalized grayscale images. SD is standard deviation.

|      | Min of Median | Max of Median | Mean of Median | SD of Median |
|------|---------------|---------------|----------------|--------------|
| *G1* | 93            | 244           | 170.80         | 40.66        |
| *G2* | 183           | 194           | 186.67         | 3.20         |
| *G3* | 0.875         | 1.161         | 1.010          | 0.063        |

From a qualitative perspective, Figure 6.1 shows some examples of original and normalized images using the modified grey-world algorithm described above. The original images have a wide variety of colour background intensities (see Figure 6.1(a) and (b)) and the medians of their histograms (Figure 6.1(e) and (f)) are significantly different, with values 148 and 169, respectively. After the normalization process, the normalized images are relatively consistent as indicated by the histograms of the normalized images. They have median values at the same indexes, at around 190 (see Figure 6.1(g) and (h)) and, visually, the colour backgrounds of the normalized images are similar (see Figure 6.1(c) and (d)).

Subsequently, the quantitative experiment demonstrated the consistency of images as shown in Table 6.1. Sixty-three images from seven positive slides (Section 1.11) were involved. The original images consist of a great range of median intensities from 93 to 244 with mean and standard deviation at 170.80 and 40.66, respectively. In contrast, the normalized images have a much narrower range of intensity from 183 to 194 with median and standard deviation at 186.67 and 3.20, respectively. Also, using the median normalization (Section 2.2.2) yields a narrow range for median intensity (from 0.875 to 1.161) with the average of the median and standard deviation of median being at 1.010 and 0.063 respectively.

Table 6.2: Maximum intensity of leukocytes.

| Slides | Colour Normalized Images | | Grayscale Normalized Images | |
|---|---|---|---|---|
| | Minimum | Maximum | Minimum | Maximum |
| 1 | 162 | 209 | 0.6219 | 0.9563 |
| 2 | 142 | 208 | 0.4122 | 0.7266 |
| 3 | 166 | 233 | 0.7234 | 1.1959 |
| 4 | 169 | 210 | 0.7212 | 1.0118 |
| 5 | 154 | 237 | 0.4895 | 0.9779 |
| 6 | 159 | 243 | 1.0152 | 1.3002 |
| 7 | 215 | 255 | 1.1302 | 1.4681 |



(a)                                    (b)

Figure 6.2: Reference image: (a) grayscale image, (b) A plot of the profile intensity along the line y = 200 in Figure 6.2(a). The leukocyte at about x = 200 and the parasites at around x = 150, 730, 820, and 950 have lower intensities in the profile.

In addition to the experiment described in Section 6.3.2, a brief study involving 420 leukocyte footprints from seven positive slides was also undertaken to ascertain if retaining full colour information provides a better final detection of malaria compared to reducing each image to grayscale based on luminance. The results indicated that there was no difference in final performance whether images were processed as full colour images or converted to grayscale. Table 6.2 shows that patterns of normalized RGB images and normalized grayscale images were the same, namely there were still differences in the ranges of leukocyte intensity in different slides. Accordingly, images were converted to grayscale and then the median normalization was applied to get more consistent images for subsequent processing steps.

Figure 6.2(a) shows the grayscale image converted from Figure 6.1(a). In the grayscale images, darker pixels, such as leukocytes and parasites, have lower intensity (Figure 6.2(b)).

## 6.4 Segmenting Leukocytes

Automatic malaria diagnosis requires a process for automatically identifying and classifying leukocytes as normal or phagocytes. In this study, the term "phagocytes' is used to refer to leukocytes or erythrocytes containing one or more malaria parasites. The phagocytes might be leukocytes engulfing malaria parasites, gametocytes, or groups of parasites resulting from haemolysed schizonts (a form of the malaria parasites, as described in Section 1.4 and 1.7). Many image analysis methods are designed with the assumption that target objects of interest are bright (foreground) compared to dark background. For this reason, the grayscale images were inverted prior to further processing steps (Figure 6.3(a)), meaning that leukocytes (as a target) have higher intensity values than other objects in images. However, leukocytes, in fact, have a variable intensity range and are superimposed on varying background (Figure 6.3(b)). The effects of various backgrounds within images will be discussed in the parasite segmentation step (Section 6.6). Furthermore, because of the remaining small differences in image background intensity between images from different slides, an adaptive threshold was determined to segment leukocytes from background and other objects. To do this, seventy ROIs containing leukocytes from seven slides were used to determine an adaptive leukocyte intensity threshold, minimum leukocyte area, and the minimum intensity of leukocytes.



(a)                                     (b)

Figure 6.3: Pre-processing steps: (a) the inverted image from Figure 6.2(a), (b) A plot of the profile intensity along the line y = 200 in Figure 6.3(a). The leukocyte at about position x = 200 has a wide tall spike profile. The parasites at approximately x = 150, 730, 820, and 950 result in thin narrow spikes in the profile. The other components have a varying profile.

Based on Equation (1.3), the number of leukocytes must be greater than zero in order to obtain a definite parasitemia score. Also, considering the nature of thick blood films explained in Section 1.9.2, the highest intensity value is always located at leukocyte pixels and the leukocyte intensities vary between images, especially between slides (Table 6.2). Then, by inspection of seventy image intensity profiles, e.g., Figure 6.3(b), an adaptive

leukocyte threshold value of 0.7 of the highest intensity pixels generally associated with leukocytes was decided empirically to segment leukocytes from other components.

The reason for using the maximum leukocyte intensity as reference for the threshold of leukocytes and parasites was that leukocytes must be present to calculate parasitemia score (Equation (1.3)). The leukocyte threshold value was low enough to keep leukocytes but high enough to eliminate large numbers of parasites in all images. By using this adaptive leukocyte threshold, segmented images consisted of leukocyte footprints and a few parasite footprints in some cases.

The minimum leukocyte intensity was found to be 0.412 (see Table 6.2). In cases where the highest image intensity pixel was less than the minimum leukocyte intensity (for example, there is no leukocyte in the image), the leukocyte intensity threshold was the minimum leukocyte intensity (0.412). The adaptive threshold was applied to each inverted image to obtain binary images consisting of a foreground of isolated regions, referred to as footprints of leukocytes, and some noise (Figure 6.4(a)). The noise is removed by using an opening filter.



Figure 6.4: Segmented leukocytes. (a) the binary image results from the grayscale image in Figure 6.3(a) after applying the optimum adaptive threshold, (b) the footprints of leukocytes results from the binary image in Figure 6.4(a) after applying an opening filter.

The opening filter, described in Section 2.3.2, requires a circular disk structure element constructed by a size parameter, $R$. For very small values of $R$, small noise is removed. On the other hand, very big values of $R$ results in eliminating leukocyte footprints. The value R = 5 was determined empirically to provide a reasonable size to retain leukocytes and remove noise. The opening filter was applied to remove noise. The resulting image

(Figure 6.4(b)) contained only footprints corresponding to leukocytes and a phagocyte, in the bottom-left quadrant.

## 6.5   Leukocytes Classification

To compute parasitemia, both the number of leukocytes and the number of parasites must be counted. Since parasites also lie within phagocytes, the phagocyte footprints, which might present among leukocyte footprints, must be investigated. To classify candidate leukocyte footprints (e.g., Figure 6.4(b)) as leukocyte or phagocyte footprints, the histogram of the grayscale intensity values of the leukocytes and phagocytes matching with their leukocyte footprints were considered.

### 6.5.1   ROIs for training of leukocyte classification

One hundred ROI of individual isolated leukocytes were identified visually from 33 images. Of these, fifty were normal leukocytes and fifty were phagocytes. Of the fifty normal leukocytes, 25 were randomly selected for training to select significant features and 25 were reserved for testing to determine the best objective function for classification. Likewise, fifty phagocytes were randomly separated into two groups of equal size for training and testing.

### 6.5.2   Experimental Details

The intensity values for leukocytes were extracted from the grayscale image using the footprint as a template. For each footprint, the mean, the standard deviation, skewness, and kurtosis of the distribution of grayscale intensity values for the leukocytes were extracted. Accordingly, an exhaustive search over possible features for the two parameters was also conducted. In the following, $m$, $v$, $s$, and $k$ denote the mean, standard deviation, skewness, and kurtosis of the intensity distribution over the leukocyte footprint.

In the image pre-processing step, images containing leukocytes or phagocytes were converted to grayscale images, and pre-processed by applying image normalization described in Section 6.3. Thresholding segmentation (Section 6.4) was applied to the normalized images and 100 ROIs (Section 6.2) were selected. For each selected ROI corresponding to a leukocyte, $m$, $v$, $s$ and $k$ were extracted from the normalized grayscale images.

Feature selection used the sequential forward feature selection method (Section 2.10.2) using a discriminant analysis classifier (Section 2.9.1), and 10-fold cross validation

(Section 2.10.3) was applied to compute classification error of every feature or combination of features. The set of training ROIs was used to determine optimal feature combinations and the set of testing ROIs was used to determine classifier performance on unseen data.

### 6.5.3   Results

The values of the statistical features for classifying ROIs as normal leukocyte or phagocytes are presented in Table 6.3. Table 6.4 shows the results of feature selection for the set of training ROIs and a confusion matrix for selected features (v, k, s) is presented in Table 6.5.

Table 6.3: Features from leukocytes and phagocytes.

|  | Phagocytes | | | | Leukocytes | | | |
|---|---|---|---|---|---|---|---|---|
|  | $m$ | $v$ | $s$ | $k$ | $m$ | $v$ | $s$ | $k$ |
| min | 2.4209 | 0.0166 | -1.0700 | 1.8823 | 2.6127 | 0.3763 | -1.1139 | 1.5898 |
| mean | 2.9450 | 0.1399 | -0.5766 | 2.6032 | 3.2272 | 0.6931 | -0.6650 | 2.2087 |
| max | 3.6234 | 0.2846 | 0.0128 | 4.0891 | 3.8120 | 1.3252 | -0.0987 | 3.3250 |
| SD | 0.3085 | 0.0793 | 0.2765 | 0.6523 | 0.2980 | 0.2655 | 0.2360 | 0.4351 |

Table 6.4: Classificfation performance on the training set using sequential forward feature selection and 10-fold cross validation.

| 50 observation, 10-Fold Cross validation | | | |
|---|---|---|---|
| Features | Error rate | Features | Error rate |
| $m$ | 0.41 | $v, s$ | 0.11 |
| $v$ | 0.10 | $v, k$ | 0.10 |
| $s$ | 0.33 | $v, k, m$ | 0.11 |
| $k$ | 0.47 | $v, k, s$ | 0.05 |
| $v, m$ | 0.11 | $v, k, s, m$ | 0.05 |

### 6.5.4   Discussion and Conclusion

According to Table 6.3 and Table 6.4, the feature subset of $v$, $k$, and $s$ and the feature subset of $v$, $k$, $s$ and $m$ have the least misclassification error. In other words, these feature subsets have the greatest potential for classifying leukocyte footprints as representing leukocyte or phagocytes.

Since the chance of over fitting increases with the number of features, smaller feature sets are generally preferred. In this case, the error rate did not improve by including $m$ and the feature subset $v, k, s$ is preferred over the subset $v, k, s, m$. Also, the pairs of the features (Table 6.4 and Figure 6.5) indicate that the feature $m$ is not specific for distinguishing

leukocyte or phagocyte. The classification error of the combination of *v* and *m* in Table 6.4 may appear satisfactory just because the selected leukocyte ROIs are mostly associated with high variance intensity while phagocytes are mostly associated with low variance intensity (Figure 6.5). This means that variance of intensity might be playing the main role in classifying leukocyte and non-leukocyte footprint. Otherwise, the role of mean alone is unclear. Therefore, this feature was not used in classifying leukocytes in the remainder of the study. On the other hand, the subset *v*, *k* and *s* was reasonable for adoption as the set of features for classification.



Figure 6.5: Classification results using quadratic discriminant analysis for leukocyte classification. The classification error was 0.11.

## 6.6 Segmenting Parasites

Aside from the leukocyte segmentation and classification described in Sections 6.4 and 6.5, automatic malaria diagnosis requires parasite detection and enumeration. To segment parasites, a difference of Gaussian filters (Section 2.8) was applied to remove varying background and retain the spikes in the image associated with parasites. In order to implement the filter, the size of the filter and the standard deviation must be set. The key to choosing the size of the filter (the size of the non-zero patch on which the function is defined) is to make it large enough that near the edge, the filter has value zero. If not, then the jump from non-zero values at the edge and the implied zero values beyond the edge of the filter will create artefacts. Otherwise, the filter should be as small as possible since the size determines the run time. The size of 31 x 31 pixels and standard deviations $\sigma = 3.0$ and $\sigma = 3.9$ were determined empirically (Figure 6.6(a) and (b)). The standard deviation for the first

Gaussian was found by trying a few values to get a good balance between enhancing parasites and blurring of the image. The effect of applying the filter in Figure 6.3(a) results in "flat" image though blurred (Figure 6.7(a)).



(a)                    (b)

Figure 6.6: Filter Gaussian: (a) the difference of Gaussian filter for background removal and enhancing parasites, (b) profiles of a parasite and the filter superimposed. The blue line is a part of the profile of Figure 6.3(b) focused on a spike due to a parasite. The red line is the profile of the first Gaussian filter.

By inspection of some profiles, e.g., Figure 6.7(d), the parasite threshold value of 0.03 of the highest intensity pixels, was determined. This threshold was applied to the filtered images from which the leukocytes had been removed to segment parasites from other objects. The parasite threshold value was low enough for parasites to emerge but high enough to eliminate thrombocytes and plasma. The resulting images, Figure 6.8(a), contain only footprints associating with parasites. However, an effect of applying the filter of size $(2n+1)$ x $(2n+1)$ is the presence of white lines near the edge (within $n$ pixels of the edge) of the images. Compared to the whole area, the edge is not significant. Hence, the border of width n pixels was removed (Figure 6.8(b)). The next step was to count the number of parasites.

Figure 6.7: Processing for parasite segmentation. (a) The flat image of Figure 6.3(a), (b) the flat image minus non-dilated leukocytes, (c) the flat image minus dilated leukocytes, (d) parasite intensity profile.



Figure 6.8: Segmented parasites. (a) Parasite image with white lines near the edge of image, (b) parasite image.

## 6.7 Parasite Classification

In some cases, thrombocytes have higher intensity than the parasite threshold. Also, some small sections of leukocyte border occasionally remain even though leukocyte dilation and subtraction (Section 6.6) were applied. As a result, some parasite footprint images (Section

6.6) contain thrombocyte footprints and small leukocyte borders. A method for distinguishing parasites from thrombocytes and leukocyte borders suitable for thick blood films was developed in the following experiment.

### 6.7.1 Data and Experimental Details

Seventy-five ROIs were identified visually from 75 images. Of these, 25 were parasite footprints, 25 ROIs were individual isolated thrombocyte footprints, and 25 ROIs were leukocyte border footprints.

In this experiment, the intensity values derived from the ROIs were extracted. For each selected ROI corresponding to parasite, thrombocyte, or leukocyte border, a variance of luminance feature ($Y$), heterogeneity ($h$) (Section 2.6.3), and the four central moments: mean ($m$), variance ($v$), skewness ($s$) and kurtosis ($k$) (Section 2.6.1) were calculated. Quadratic discriminant analysis (Section 2.9.1) was used to classify the candidate parasites as true parasites or not.

### 6.7.2 Results, Discussion, and Conclusion

The values of features mentioned in the previous section for classifying ROIs as parasites or not are presented in Table 6.5. The classification error resulting from the combination of $m$, $v$, and $h$ (Table 6.6) appears to be the most satisfactory. This means that the mean and variance of intensity and the heterogeneity are suitable for classifying parasite and non-parasite footprints in thick blood films. Therefore, the other features, $s$, $k$, and $Y$, were not used in classifying parasites in the remaining of the study.

Table 6.5: Features from parasites, thrombocytes, and leukocytes. $P$, $T$, and $L$ denote parasites, thrombocytes, and leukocyte borders. $SD$ is standard deviation. In this case, $vY$ is variance of $Y$ value.

|   |   | $m$ | $v$ | $s$ | $k$ | $h$ | $vY$ |
|---|---|---|---|---|---|---|---|
| $P$ | min | 1.7462 | 0.0475 | 0.0012 | 1.5041 | 0.1933 | 2.0539 |
|  | mean | 2.6091 | 0.2192 | 0.2101 | 2.2931 | 0.3952 | 5.2791 |
|  | max | 3.8606 | 0.5437 | 0.9407 | 3.1318 | 0.6500 | 12.0602 |
|  | $SD$ | 0.5266 | 0.1295 | 0.2111 | 0.3860 | 0.1066 | 2.0914 |
| $T$ | min | 2.9631 | 0.0477 | 0.0744 | 1.7942 | 0.3500 | 1.6348 |
|  | mean | 3.6596 | 0.2304 | 0.4099 | 2.4017 | 0.4688 | 6.3246 |
|  | max | 5.0565 | 0.3952 | 1.1081 | 3.8485 | 0.5667 | 8.6746 |
|  | $SD$ | 0.6111 | 0.0846 | 0.3019 | 0.5454 | 0.0622 | 1.9498 |
| $L$ | min | 2.8149 | 0.0772 | 0.0149 | 1.8260 | 0.2788 | 2.6953 |
|  | mean | 4.9813 | 0.4130 | 0.4806 | 2.4532 | 0.4508 | 9.2472 |
|  | max | 6.8808 | 0.8481 | 0.8196 | 3.4544 | 0.5833 | 17.4172 |
|  | $SD$ | 1.0980 | 0.1676 | 0.2039 | 0.3852 | 0.0702 | 2.9292 |

Table 6.6: Feature selection with quadratic discriminant analysis classifier to classify parasites, leukocyte borders, or thrombocytes. *FC* is feature combination and *ER* is error rate.

| 75 observations, 5-Folds Cross-validation | | | | | | | |
|------|-------|-------|-------|-----------|-------|----------------|-------|
| *FC* | *ER* | *FC* | *ER* | *FC* | *ER* | *FC* | *ER* |
| *m* | 0.307 | *m, v* | 0.280 | *m, h, s* | 0.253 | *m, h, vY, s, v* | 0.227 |
| *v* | 0.560 | *m, s* | 0.360 | *m, h, k* | 0.307 | *m, h, vY, s, k* | 0.307 |
| *s* | 0.533 | *m, k* | 0.32 | *m, h, vY* | 0.227 | *m, h, vY, s, v, k* | 0.227 |
| *k* | 0.747 | *m, h* | 0.253 | *m, h, vY, v* | 0.253 | | |
| *h* | 0.520 | *m, vY* | 0.280 | *m, h, vY, s* | 0.240 | | |
| *vY* | 0.440 | *m, h,* | 0.267 | *m, h, vY, k* | 0.253 | | |

The classification performance based on mean, variance, and heterogeneity for distinguishing parasites, thrombocytes, or leukocyte borders are displayed as confusion matrices (Table 6.7). According to Table 6.7, the classification accuracy was 0.85. However, the main purpose of this part of this study was to determine a pattern for classifying ROIs as parasites or non-parasites. When viewed as a classification into two class, parasites and non-parasites (thrombocytes and leukocyte borders), the accuracy was 88% (Table 6.8).

Table 6.7: Confusion matrix for classification of parasites, leukocyte borders, or thrombocytes. The classification accuracy is 85%.

| | | Predicted | | |
|---|---|---|---|---|
| | | Parasites | Leukocyte borders | Thrombocytes |
| Actual Class | Parasites | 18 | 2 | 0 |
| | Leukocyte borders | 6 | 17 | 6 |
| | Thrombocytes | 1 | 5 | 20 |

Table 6.8: Confusion matrix for classification of parasite or non-parasites in the training stage. The classification accuracy is 88%. The sensitivity and specificity are 90% and 87.27%, respectively.

| | | Predicted | |
|---|---|---|---|
| | | Parasites | Non-parasites |
| Actual Class | Parasites | 18 | 2 |
| | Non-parasites | 7 | 48 |

## 6.8 Estimating Parasitemia

As a preliminary study, samples from several images were used to identify individual parasites and determine an average size. A simple way to do this is to visually examine

images, such as in Figure 6.8(b), and make sure that only bright spots associated with parasites in the original image are used for such a calculation. Once an average parasite area is known, then counting the number of parasites in the image is just a matter of summing over the area of white pixels and dividing by the standard average parasite size determined previously. If there are some noise pixels in the image, these do not contribute much to the total area and so do not have much effect on the number of parasites counted.

For validation of parasite classification, a data subset of around 1000 images from the testing data set (Section 6.2) was selected randomly and proportionally. This subset was used to test the parasite classification method described above. The performance of the method was also validated by comparing these automatically identified parasites to those identified by an informal reader.

The final step, the calculation of parasitemia was done using the conventional method (Equation (1.3)) described in Section 1.8. The number of leukocytes and parasites were calculated over all images in the testing data set described in Section 6.2. The performance of this algorithms for estimating parasitemia was validated by comparing these parasitemia estimations to the parasitemia scores given by expert readers. Beside correlation, due to the small sample involved in this study, a nonparametric test, t-test was used to validate this algorithm.

## 6.9   Results for Estimating Parasitemia

During initial experiments to obtain a standard threshold of leukocyte and parasites, iteration with various thresholds applied to several image samples from all slides showed that an optimum threshold for producing precise areas of leukocytes and parasites compared with visual investigation were 0.7 and 0.03 of the maximum intensity of leukocytes, respectively.

Table 6.9. Confusion matrix of parasite classification as a part of the estimating parasitemia experiment.

|  |  | Predicted | |
| --- | --- | --- | --- |
|  |  | Parasites | Non-parasites |
| Actual Class | Parasites | 7492 | 2774 |
|  | Non-parasites | 520 | 1335 |

Validation scores for parasite classification were recorded in a confusion matrix (Table 6.9). The accuracy of classification was 72.82%, the sensitivity was 72.98%, and the specificity was 71.97%.



Figure 6.9: Estimating parasitemia scores of positive slides. For each slide, the algorithm produces an estimate of the parasitemia scores indicated with blue stars. For each slide, there are several estimates of parasitemia scores from human experts (Table 1.1), the logs of which provide the vertical coordinates of the green open circles and the blue stars. The log of median provides the horizontal coordinate of the green open circles and the blue stars.

Parasitemia scores per slide generated by the automatic diagnosis process fitted well with expert parasitemia scores (r = 0.79 and p-value = 0.40 at α = 0.05; Table 6.10 and Figure 6.9).

Table 6.10: Comparison between estimated parasitemia and median of expert's scores. P-value obtained from t-test is 0.38. *EP* is estimate of parasitemia scores and *EM* is median of experts' parasitemia scores. *PS1*, *PS2*, *PS3*, …, *PS7* are positive slide 1, 2, 3, …, 7. *NS1*, *NS2*, and *NS3* are negative slide 1, 2, and 3.

| Slides | NS1 | NS2 | NS3 | PS1 | PS2 | PS3 | PS4 | PS5 | PS6 | PS7 |
|--------|-----|-----|------|-----|-----|------|-----|------|-------|--------|
| EP | 39 | 51 | 1631 | 140 | 129 | 1769 | 898 | 8934 | 54677 | 63440 |
| EM | 0 | 0 | 0 | 103 | 634 | 1340 | 5583 | 11442 | 32730 | 319393 |

## 6.10 Discussion and Conclusion

In the course of this study, choices were made regarding the inclusion of colour intensity features and details of extracting and computing these features, multiple image segmentation, and estimating parasitemia. At the lowest level, the choices were made

according to the experience of microscopists and known manifestations of malaria parasites in thick blood films. Examples include focusing on the ROIs of leukocytes and parasites and investigating colour and features associated with intensity and area. However, the understanding in biology and physics of the appearance of leukocytes and malaria parasites in constrained thick blood films explained in Section 1.9.3 and Section 6.2, is not sufficient to accurately predict colour features that contribute to good automatic identification of parasites. Accordingly, some choices regarding the details of extracting features were made without specific guidance. In these cases, choices were made based on small experiments and observations.

Although this study and that of Frean (Frean 2009) similarly aim to estimate parasitemia on thick blood films, the respective ways and methods are different. Frean diagnosed malaria parasites semi-automatically, leukocytes were identified manually. In this study, parasitemia was estimated in a fully automatic way. Both leukocytes and parasites were enumerated automatically. Furthermore, this study was able to detect phagocytised parasites.

In the context of parasite classification, the performance of testing is relatively less than that of parasite classification (Section 6.7). According to Table 6.10, the big contributor of misclassification was *FN*, which may due to phagocyte (gametocyte) borders detected as non-parasites. Phagocyte borders were not considered in the training. Another contributor was a lower *TN*. Parasite segmentation resulted in parasite candidates, in which the number of parasites are much higher than the number of non-parasites. Statistically, the probability of getting non-parasites (*TN*) in testing is much lower than that of getting parasites (*TP*). This was significantly different from training, where the number of ROIs of parasites and the number of ROIs of non-parasites were the same. Accordingly, the specificity in testing was lower than that in training.

Another big contributor of misclassification is negative images coming from a thin patch of thick blood films (Figure 6.10(a)). A thin patch is a patch from the edges of thick blood films. In thin patch areas, the staining of the residual haemolysed erythrocytes is relatively transparent (thin) and clear. In this area, the thrombocytes strongly respond to the stain used (Figure 6.10(b)); therefore, these thrombocytes have a darker colour (low intensity) almost similar to parasites (Figure 6.10(f)) from a positive slide. Thrombocytes in thin areas of thick blood films and parasites have almost the same response to the stain. This produces many *FP,* lowering the performance of parasite classification. Conversely, the

residual haemolysed erythrocytes in thick patches of thick blood films (e.g., Figure 6.10(c)) is relatively thicker than that in the thin patches, is evenly distributed, and have absorbed more dye. Therefore, the colour response of thrombocytes is similar to that of the residual haemolysed erythrocytes (Figure 6.10(d)). As a result, the colour intensities of inverse images from this area are lower than that of the parasite adaptive threshold. Thus, these thrombocytes were not assigned as parasite candidates. This means they were not involved in the parasite classification process.



| (a) | (b) |
| (c) | (d) |
| (e) | (f) |

Figure 6.10. Patch (Pa) areas of thick blood films. (a) A patch of images from thin area of a negative slide. (c) A full thick patch area of a negative slide. (e) A semi full thick patch area of positive slide. (b), (d), and (f) are image examples from (a), (c), and (e), respectively.

Furthermore, in some cases, patches of imperfect staining in thick blood films meant that the colour response of gametocytes were brighter (lower intensities in inverse images) than the parasite adaptive threshold. As a result, these gametocytes were not assigned as parasite candidates. This means they were not involved in the parasite classification process

and thus, the probability of *TP* decreased. Accordingly, the accuracy and sensitivity in testing are lower than those in training.

Despite that fact that sources of misclassification can be identified, overall, the methods presented in this chapter are useful for estimating parasitemia in the sense that estimates fell within the ranges of expert parasitemia scores.

# Chapter 7: Concluding Remarks and Future Work

The main objective of this study was to develop an image analysis method for automatic malaria diagnosis and parasitemia estimation. Previous computer-aided malaria diagnosis schemes that were based mainly on thin blood film images, only considered erythrocytes and parasites, or emphasized reporting just the accuracy of identifying parasites (Section 1.10), which is not the same as estimating parasitemia. Unlike the previous studies, this study considers all the main components in blood as must be done in fully automatic systems. This study goes farther than many previous studies in that blood films with different levels of parasitemia are considered, not just one level. In addition, this study is based not only on thin blood film images, but also on thick blood film images. In malaria diagnosis, the thick blood film is important because one thick blood film consists of many layers of erythrocytes, so that large amounts of blood can be examined quickly and easily. Usually, experts diagnose malaria parasites and estimate the parasitemia based on thick blood films (WHO 2010). So, one may expect that automatic malaria diagnosis based on thick blood film images may also be important. However, this study found that automatic malaria diagnosis based on thick blood film images was limited.

The objective of developing a method for accurate diagnosis was not successful. A number of false positive parasites were detected in the negative blood films; therefore, zero parasitemia scores were not obtained. However, the objective of developing a method for estimating parasitemia was successful. This study reveals that the method is able to estimate parasitemia based on automatic analysis of thin and thick blood film images with error no greater than the variation between expert readers. Relevant discussions and conclusions were included at the ends of Chapters 3, 4, 5, and 6.

Chapter 3 investigated infected erythrocytes based on erythrocyte features to estimate parasitemia scores. This study perceived that the estimation approach was moderately consistent with Dowling and Shute's (1966). The system is well suited for the three high parasitemia slides and slightly overestimated for the middle parasitemia slide, and overestimated for the three low parasitemia slides. Moreover, the results were not in line with the Trape's (1985) study. This study found that the overestimations were from three reasons described in Section 3.6.2.

Chapter 4 introduced a new feature for parasite segmentation of thin blood film images and investigated erythrocytes and parasites to identify infected erythrocytes and

estimate parasitemia scores. The results show that the parasitemia estimation fitted well with parasitemia scores from expert readers (coefficient correlation r = 0.97 and p-value = 0.54 at α = 0.05. In addition, the outcomes were in line with both Dowling and Shute's (1996) study and Trape's (1985) study. Therefore, the combination of erythrocyte detection, location feature, and parasite identification to investigate infected erythrocytes resulted in better parasitemia estimation than that in Chapter 3.

Chapter 5 investigated parasite detection in thick blood film images based on physical appearance and size of parasite footprints to estimate parasitemia scores. The morphological approach was not successful for estimating parasitemia. In addition to size, this study indicated that intensity feature is needed to extract parasite profiles for parasite classification.

Chapter 6 developed a parasite detection method in thick blood film images to estimate parasitemia scores. Following the suggestion by Kumar et al. (Kumar et al., 2012) and that of Chapter 5 in parasite identification, this study evaluated grayscale and green colour channel to extract parasite profile for estimating parasitemia. The results demonstrated that the parasitemia estimation was not significantly different from experts' parasitemia scores (Table 6.9 and Figure 6.9), with coefficient correlation r = 0.79 and p-value = 0.40 at α = 0.05.

It should be noted that parasitemia estimation based on thin blood film images described in Chapter 3 systematically overestimates parasitemia. Thus, this method may be impractical for estimating low parasitemia slides or for credible assignment of negative slides. On the other hand, identifying erythrocytes combined with analysing parasites separately from the erythrocyte host has significantly better performance in terms of reducing the overestimation of parasitemia.

Meanwhile, parasitemia estimation based on thick blood film images indicated that there was no significant difference between this automatic malaria diagnosis system and experts' parasitemia scores. The parasitemia scores in all positive slides fitted well within manual expert readers. Nevertheless, there were some slight overestimation of parasitemia.

Each of these observations have possibly important ramifications and are seeds for future work. Although these results do not establish that a particular method for analysing blood film images is clearly best for estimating parasitemia scores, an important set of

features has been identified. This will be of benefit to future work on combining automatic malaria diagnosis based on thin and thick blood film images.

The parasite segmentation presented in this thesis uses a multi-threshold method. This method has provided an unplanned, but possibly important, new direction of research. This illustrates that the field of segmentation for malaria diagnosis based on thin and thick blood film images has great potential.

In this work, parasitemia estimation was tested on a limited number of slides (blood films). Large numbers of high quality images of blood films showing the full range of clinical presentation in terms of parasitemia levels (including many negative cases), different species of malaria parasites and at different life-stages of the parasite life cycle are needed in order to conduct a full study on automatic detection of malaria. No such data set currently exists, or is difficult to obtain since blood films are not normally imaged during examination and special equipment is needed to do so.

In summary, the main contribution of this study is that automatic methods for estimating malaria parasitemia were developed which predict parasitemia scores that are not significantly different from expert readers. The first contribution is the introduction of a new feature to segment malaria parasites in thin blood film images. Supported by a literature review and empirical evidence, the study contributes the following:

1. This study refines the benefit of colour image normalization from Tek, et al. (Tek, F B, Dempster & Kale 2006) using the diagonal model (Barnard, Cardei & Funt 2002) and the database grey-world algorithm (Hordley & Finlayson 2004). To obtain the grey values, this study uses a real background image of a slide as a reference image (Section 4.2).
2. This study introduces a new feature derived from the natural characteristics of components of blood film to determine an adaptive threshold for image segmentation.

The further contribution is the estimation of parasitemia scores based on the percentage of infected erythrocytes in thin blood film images. As a validation, these results are compared to parasitemia scores from expert readers based on thick blood films and manual erythrocyte identification from informal readers in thin blood films.

# Bibliography

Aikawa, M, Iseki, M, Barnwell, JW, Taylor, D, Oo, MM & Howard, RJ 1990, 'The pathology of human cerebral malaria', *Am J Trop Med Hyg*, vol. 43, no. 2, Aug, pp. 30-7.

Al-khatib, HB & Al-Horani, A 2012, 'Predicting financial distress of public companies listed in Amman Stock Exchange', *European Scientific Journal*, vol. 8, no. 15, July, pp. 1-17.

Altman, EI, Marco, G & Varetto, F 1994, 'Corporate distress diagnosis: Comparisons using linear discriminant analysis and neural networks (the Italian experience)', *Journal of banking & finance*, vol. 18, no. 3, January, pp. 505-29.

Arce, GR 2005, *Nonlinear signal processing: a statistical approach*, John Wiley & Sons, New Jersey.

Arias-Castro, E & Donoho, DL 2009, 'Does median filtering truly preserve edges better than linear filtering?', *The Annals of Statistics*, vol. 37, no. 3, June, pp. 1172-206.

Avci, D & Varol, A 2009, 'An expert diagnosis system for classification of human parasite eggs based on multi-class SVM', *Expert Systems with Applications*, vol. 36, no. 1, pp. 43-8.

Bain, BJ 2014, *Blood cells: a practical guide*, 4th edn, Blackwell, Victoria, Australia.

Baird, J, Richie, TL, Marwoto, H & Gunawan, S 1995, 'Epidemic malaria among transmigrants in Irian Jaya', *Buletin Penelitian Kesehatan*, vol. 23, no. 3, September, pp. 18-34.

Ballard, DH & Brown, CM 1982, *Computer Vision*, Prentice Hall, Inc., New Jersey, USA.

Bangs, MJ & Subianto, DB 1999, 'El Niño and associated outbreaks of severe malaria in highland populations in Irian Jaya, Indonesia: a review and epidemiological perspective', *Southeast Asian Journal of Tropical Medicine and Public Health*, vol. 30, no. 4, December, pp. 608-19.

Barnard, K 1999, 'Practical Colour Constancy', PhD thesis, Simon Fraser University, British Colombia, Canada.

Barnard, K, Cardei, V & Funt, B 2002, 'A comparison of computational color constancy algorithms-Part I: Methodology and experiments with synthesized data', *IEEE Transactions on Image Processing*, vol. 11, no. 9, September, pp. 972-84.

Basalamah, S 2012, 'Histogram based circle detection', *International Journal of Computer Science and Network Security*, vol. 12, no. 8, August, pp. 40-3.

Binghan, L, Xiuduan, F, Weizhi, W & Zhiyong, Z 2002, 'Automatic separation of overlapping objects', *Proceedings of the 4th World Congress on Intelligent Control and Automation*, IEEE Robotics and Automation Society, Sanghai, China, pp. 2901-5.

Bishop, CM 2006, *Pattern Recognition and Machine Learning*, Springer Berlin Heidelberg, New York, USA.

Borgefors, G 1984, 'Distance transformations in arbitrary dimensions', *Computer vision, graphics, and image processing*, vol. 27, no. 3, February, pp. 321-45.

Burger, W & Burge, MJ 2016, *Digital image processing: an algorithmic introduction using Java*, 2nd edn, Springer-Verlag, London, UK.

Canny, J 1986, 'A computational approach to edge detection', *IEEE transactions on pattern analysis and machine intelligence*, no. 6, November, pp. 679-98.

CDC 2016, *Malaria*, Global Health – Division of Parasitic Diseases and Malaria, Atlanta, viewed 14 November 2016, <https://www.cdc.gov/malaria/>.

Chan, H-P, Petrick, N & Sahiner, B 2000, 'Computer-Aided Breast Cancer Diagnosis', in A Jain, A Jain, S Jain & L Jain (eds), *Artificial Intelligence Techniques in Breast Cancer Diagnosis and Prognosis*, World Scientific Publishing Co. Pte. Ltd., Singapore, vol. 39, pp. 179-264.

Charayaphan, C & Marble, AE 1992, 'Image processing system for interpreting motion in American Sign Language', *Journal of Biomedical Engineering*, vol. 14, no. 5, September, pp. 419-25.

Dakić, Z, Ivović, V, Pavlović, M, Lavadinović, L, Marković, M & Djurković-Djaković, O 2014, 'Clinical significance of molecular methods in the diagnosis of imported malaria in returning travelers in Serbia', *International Journal of Infectious Diseases*, vol. 29, August, pp. 24-30.

Danielson, P-E 1978, 'A new shape factor', *Computer Graphics and Image Processing*, vol. 7, no. 2, September, pp. 292-9.

Dash, M & Liu, H 1997, 'Feature selection for classification', *Intelligent data analysis*, vol. 1, no. 3, March, pp. 131-56.

Díaz, G, González, FA & Romero, E 2009, 'A semi-automatic method for quantification and classification of erythrocytes infected with malaria parasites in microscopic images', *Journal of Biomedical Informatics*, vol. 42, no. 2, 4//, pp. 296-307.

Dluzewski, A, Mitchell, G, Fryer, P, Griffiths, S, Wilson, R & Gratzer, W 1992, 'Origins of the parasitophorous vacuole membrane of the malaria parasite, Plasmodium falciparum, in human red blood cells', *J Cell Sci*, vol. 102, no. 3, pp. 527-32.

Dorini, LB, Minetto, R & Leite, NJ 2007, 'White blood cell segmentation using morphological operators and scale-space analysis', *XX Brazilian Symposium on Computer Graphics and Image Processing*, IEEE Computer Society, Los Alamitos, Brazil, pp. 294-304.

Dowling, M & Shute, G 1966, 'A comparative study of thick and thin blood films in the diagnosis of scanty malaria parasitaemia', *Bulletin of the World health Organization*, vol. 34, no. 2, p. 249.

Duda, RO & Hart, PE 1972, 'Use of the Hough transformation to detect lines and curves in pictures', *Communications of the ACM*, vol. 15, no. 1, January, pp. 11-5.

Dudani, SA, Breeding, KJ & McGhee, RB 1977, 'Aircraft identification by moment invariants', *IEEE transactions on computers*, vol. 100, no. 1, January, pp. 39-46.

Efford, N 2000, *Digital Image Processing: A Practical Introduction Using Java 2000*, Pearson Education, Essex, UK.

Efron, B 1983, 'Estimating the error rate of a prediction rule: improvement on cross-validation', *Journal of the American Statistical Association*, vol. 78, no. 382, June, pp. 316-31.

Elbehiery, H, Hefnawy, A & Elewa, M 2005, 'Surface Defects Detection for Ceramic Tiles Using Image Processing and Morphological Techniques', *International Journal of*

*Computer, Electrical, Automation, Control and Information Engineering*, vol. 1, no. 5, pp. 158-62.

Elter, M, Haßlmeyer, E & Zerfaß, T 2011, 'Detection of malaria parasites in thick blood films', *2011 33rd Annual International Conference of the IEEE Engineering in Medicine and Biology Society*, IEEE Engineering in Medicine and Biology Society, Massachutsetss, USA, pp. 5140-4.

Elyazar, IR, Hay, SI & Baird, JK 2011, 'Malaria distribution, prevalence, drug resistance and control in Indonesia', in *Advances in parasitology*, Elsevier Ltd., vol. 74, pp. 41-475.

Farnoosh, R & Zarpak, B 2008, 'Image segmentation using Gaussian mixture model', *IUST International Journal of Engineering Science*, vol. 19, no. 1-2, April, pp. 29-32.

Fisher, RA 1936, 'The use of multiple measurements in taxonomic problems', *Annals of eugenics*, vol. 7, no. 2, September, pp. 179-88.

Frean, J 2009, 'Reliable enumeration of malaria parasites in thick blood films using digital image analysis', *Malaria Journal*, vol. 8, no. 1, September, p. 218.

Freitas, RA 1999, *Nanomedicine, volume I: basic capabilities*, Landes Bioscience, Texas, USA.

Fu, Z & Wang, L 2012, 'Color Image Segmentation Using Gaussian Mixture Model and EM Algorithm', *Proceedings of the 2nd International Conference on Multimedia and Signal Processing*, Springer, Shanghai, China, pp. 61-6.

Fukunaga, K 1990, *Introduction to statistical pattern recognition*, 2nd edn, Academic press, California, USA.

Garcia, LS & Bruckner, DA 1997, *Diagnostic medical parasitology*, 5th edn, American Society for Microbiology Press, Washington, DC, USA.

Gonzalez, RC, Woods, RE & Eddins, SL 2001, *Digital Image Processing*, 3rd edn, Pearson Education, London, UK.

Hadjiiski, L, Sahiner, B, Chan, H-P, Petrick, N, Helvie, MA & Gurcan, M 2001, 'Analysis of temporal changes of mammographic features: computer-aided classification of

malignant and benign breast masses', *Medical Physics*, vol. 28, no. 11, November, pp. 2309-17.

Hall, A, Doberstyn, E, Mettaprakong, V & Sonkom, P 1975, 'Falciparum malaria cured by quinine followed by sulfadoxine-pyrimethamine', *Br Med J*, vol. 2, no. 5961, April, pp. 15-7.

Hanif, N, Mashor, M & Mohamed, Z 2011, 'Image enhancement and segmentation using dark stretching technique for Plasmodium Falciparum for thick blood smear', *2011 IEEE 7th International Colloquium on Signal Processing and its Applications*, IEEE, Penang, Malaysia, pp. 257-60.

Haralick, RM & Shapiro, LG 1992, *Computer and Robot Vision*, 1st edn, vol. 1, Addison-Wesley, New York, USA.

Haralick, RM, Sternberg, SR & Zhuang, X 1987, 'Image Analysis Using Mathematical Morphology', *IEEE transactions on pattern analysis and machine intelligence*, vol. PAMI-9, no. 4, July, pp. 532-50.

Harwood, VJ, Whitlock, J & Withington, V 2000, 'Classification of antibiotic resistance patterns of indicator bacteria by discriminant analysis: use in predicting the source of fecal contamination in subtropical waters', *Applied and Environmental Microbiology*, vol. 66, no. 9, September, pp. 3698-704.

Hastie, T, Tibshirani, R & Friedman, J 2001, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, Springer-Verlag, New York, USA.

Hordley, SD & Finlayson, GD 2004, 'Re-evaluating colour constancy algorithms', *Proceedings of the 17th International Conference on Pattern Recognition*, IEEE Computer Society, Washington, DC, USA, pp. 76-9.

Hough, PVC 1962, *Method and means for recognizing complex patterns*, Google Patents, patent, US 3069654 A.

Howard, MR & Hamilton, PJ 2013, *Haematology: an illustrated colour text*, 4th edn, Churchill Livingstone, London, UK.

Hu, M-K 1962, 'Visual pattern recognition by moment invariants', *IRE transactions on information theory*, vol. 8, no. 2, pp. 179-87.

Huang, CT & Mitchell, OR 1994, 'A Euclidean distance transform using grayscale morphology decomposition', *IEEE transactions on pattern analysis and machine intelligence*, vol. 16, no. 4, April, pp. 443-8.

Huang, TS, Yang, G & Tang, G 1979, 'A fast two-dimensional median filtering algorithm', *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 27, no. 1, February, pp. 13-8.

Huang, Y-L, Jao, Y-L, Hsieh, T-Y & Chung, C-W 2008, 'Adaptive automatic segmentation of HEp-2 cells in indirect immunofluorescence images', *2008 IEEE International Conference on Sensor Networks, Ubiquitous and Trustworthy Computing*, IEEE, Taichung, Taiwan, pp. 418-22.

Huang, Y, Englehart, KB, Hudgins, B & Chan, AD 2005, 'A Gaussian mixture model based classification scheme for myoelectric control of powered upper limb prostheses', *IEEE Transactions on Biomedical Engineering*, vol. 52, no. 11, November, pp. 1801-11.

Jagessar, RC & Rampersaud, E 2014, 'A survey of the status of malaria in Guyana and treatments: synthetic and herbal', *Journal of Physics: Conference Series*, vol. 8, March 2014, pp. 26-34.

Jähne, B 2005, *Digital Image Processing*, 6th edn, Springer Berlin Heidelberg, Berlin, Germany.

Jain, AK 1989, *Fundamentals of digital image processing*, Prentice-Hall, Inc. , New Jersey, USA.

Jayant, N 1976, 'Average-and median-based smoothing techniques for improving digital speech quality in the presence of transmission errors', *IEEE Transactions on Communications*, vol. 24, no. 9, September, pp. 1043-5.

Jayaraman, S, Esakkirajan, S & Veerakumar, T 2009, *Digital Image Processing*, Tata McGraw-Hill Education, New Delhi, India.

Karcheva, M, Atanasova, M & Rainova, I 2017, 'A Case of malaria transmitted in Bulgaria from abroad', in *Archives of the Balkan Medical Union*, Balkan Medical Union, Bulgaria, vol. 52, pp. 95-8.

Kemenkes, RI 2015, *Profil Kesehatan Indonesia Tahun 2014*, Kementerian Kesehatan RI, Jakarta, Indonesia.

Kim, K, Jeon, J, Choi, W, Kim, P & Ho, Y-S 2001, 'Automatic Cell Classification in Human's Peripheral Blood Images Based on Morphological Image Processing', in M Stumptner, D Corbett & M Brooks (eds), *AI 2001: Advances in Artificial Intelligence*, Springer Berlin Heidelberg, vol. 2256, pp. 225-36.

Kohavi, R & Provost, F 1998, 'Glossary of terms', *Machine Learning*, vol. 30, no. 2-3, pp. 271-4.

Kries, JV 1878, 'Beitrag zur physiologie der gesichtsempfinding', *Arch. Anat. Physiol*, vol. 2, pp. 5050-524.

Kruk, ZA, Bottema, MJ, Forder, R, Reyes-Veliz, L, Pitchford, WS & Bottema, CDK 2016, 'Image Analysis Method of Marbling Quality Parameters', Unpublished.

Kumar, A, Choudhary, A, Tembhare, PU & Pote, CR 2012, 'Enhanced Identification of Malarial Infected Objects using Otsu Algorithm from Thin Smear Digital', *International Journal of Latest Research in Science and Technology*, vol. 1, no. 2, pp. 159-63.

Larson, SC 1931, 'The shrinkage of the coefficient of multiple correlation', *Journal of Educational Psychology*, vol. 22, no. 1, p. 45.

Liangwongsan, S, Marungsri, B, Oonsivilai, R & Oonsivilai, A 2011, 'Extracted circle Hough Transform and circle defect detection algorithm', *World Academy of Science, Engineering and Technology*, vol. 5, pp. 432-6.

Löffler, H & Rastetter, J 2000, *Atlas of clinical hematology*, 5th edn, Springer Berlin Heidelberg, New York, USA.

Lu, J, Plataniotis, KN & Venetsanopoulos, AN 2003, 'Face recognition using kernel direct discriminant analysis algorithms', *IEEE Transactions on Neural Networks*, vol. 14, no. 1, January, pp. 117-26.

Lyons, MJ, Budynek, J & Akamatsu, S 1999, 'Automatic classification of single facial images', *IEEE transactions on pattern analysis and machine intelligence*, vol. 21, no. 12, December, pp. 1357-62.

Maichomo, M, McDermott, J, Arimi, S, Gathura, P, Mugambi, T & Muriuki, S 1999, 'Study of brucellosis in a pastoral community and evaluation of the usefulness of clinical signs and symptoms in differentiating it from other flu-like diseases', *African journal of health sciences*, vol. 7, no. 3-4, pp. 114-9.

Makkapati, VV & Rao, RM 2011, 'Ontology-based Malaria Parasite Stage and Species Identification from Peripheral Blood Smear Images', Boston, Massachusetts USA.

Martens, P, Kovats, RS, Nijhof, S, de Vries, P, Livermore, MTJ, Bradley, DJ, Cox, J & McMichael, AJ 1999, 'Climate change and future populations at risk of malaria', *Global Environmental Change*, vol. 9, pp. S89-S107.

Martens, W, Niessen, LW, Rotmans, J, Jetten, TH & McMichael, AJ 1995, 'Potential impact of global climate change on malaria risk', *Environmental health perspectives*, vol. 103, no. 5, May, p. 458.

Marwoto, HA & Sulaksono, STE 2003, 'Peningkatan Kasus Malaria di Pulau Jawa, Kepulauan Seribu dan Lampung', *Media Litbang Kesehatan*, vol. XIII, no. 3.

Mboera, L, Fanello, C, Malima, R, Talbert, A, Fogliati, P, Bobbio, F & Molteni, F 2013, 'Comparison of the Paracheck-Pf® test with microscopy, for the confirmation of Plasmodium falciparum malaria in Tanzania', *Annals of tropical medicine and parasitology*, vol. 100, no. 2, October, pp. 115-22.

McLachlan, G 2004, *Discriminant analysis and statistical pattern recognition*, John Wiley & Sons, New Jersey, USA.

McLachlan, G & Peel, D 2004, *Finite mixture models*, John Wiley & Sons, New York, USA.

Metz, CE 1978, 'Basic principles of ROC analysis', *Seminars in nuclear medicine*, vol. 8, no. 4, October, pp. 283-98.

Miller, MT, Jerebko, AK, Malley, JD & Summers, RM 2003, 'Feature selection for computer-aided polyp detection using genetic algorithms', *Proceedings of Medical Imaging 2003: Physiology and Function: Methods, Systems, and Applications*, SPIE-International Society for Optics and Photonics, San Diego, USA, pp. 102-10.

Moody, A 2002, 'Rapid diagnostic tests for malaria parasites', *Clinical microbiology reviews*, vol. 15, no. 1, January, pp. 66-78.

Moore, DS, McCabe, GP & Craig, BA 2012, *Introduction to the practice of statistics*, Seventh edn, Ruth Baruth, New York.

Mosteller, F & Wallace, DL 1963, 'Inference in an authorship problem: A comparative study of discrimination methods applied to the authorship of the disputed Federalist Papers', *Journal of the American Statistical Association*, vol. 58, no. 302, June, pp. 275-309.

Murray, PR, Rosenthal, KS & Pfaller, MA 2015, *Medical microbiology*, 4th edn, Elsevier, Philadelphia, USA.

NEQAS 2016, *Effects of Anticoagulant on Malarial Parasites*, Hospital for Tropical Diseases, London, UK, viewed 27 October 2016 2016, <http://www.ukneqasmicro.org.uk/parasitology/index.php/blood-parasitology/malaria-species>.

Oaks Jr, SC, Mitchell, VS, Pearson, GW & Carpenter, CC 1991, *Malaria: obstacles and opportunities*, National Academies Press, Washington, D.C., USA.

Otsu, N 1979, 'A threshold selection method from gray-level histograms', *Automatica*, vol. 11, no. 285-296, January, pp. 23-7.

Palakuru, SS 2016, *Malaria*, Texas, USA, viewed 8 October 2016, <http://www.austincc.edu/microbio/2704w/pf.htm>.

Permuter, H, Francos, J & Jermyn, IH 2003, 'Gaussian mixture models of texture and colour for image database retrieval', *Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing*, IEEE, pp. III-569-72.

Pratt, WK 1978, *Digital Image Processing*, Third edn, John Wiley & Sons, New York, USA.

Purwar, Y, Shah, SL, Clarke, G, Almugairi, A & Muehlenbachs, A 2011, 'Automated and unsupervised detection of malarial parasites in microscopic images', *Malaria Journal,* vol. 10, no. 1, December, pp. 364-373.

Rabiner, L, Sambur, M & Schmidt, C 1975, 'Applications of a nonlinear smoothing algorithm to speech processing', *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 23, no. 6, December, pp. 552-7.

Raudys, Š & Young, DM 2004, 'Results in statistical discriminant analysis: a review of the former Soviet Union literature', *Journal of Multivariate Analysis*, vol. 89, no. 1, February, pp. 1-35.

Raudys, SJ & Jain, AK 1991, 'Small sample size effects in statistical pattern recognition: recommendations for practitioners', *IEEE transactions on pattern analysis and machine intelligence*, vol. 13, no. 3, March, pp. 252-64.

Refaeilzadeh, P, Tang, L & Liu, H 2008, *Cross-Validation*, Arizona State University, Arizona.

Reynolds, DA, Quatieri, TF & Dunn, RB 2000, 'Speaker verification using adapted Gaussian mixture models', *Digital signal processing*, vol. 10, no. 3, July, pp. 19-41.

Rosenthal, PJ, Lee, GK & Smith, RE 1993, 'Inhibition of a Plasmodium vinckei cysteine proteinase cures murine malaria', *Journal of Clinical Investigation*, vol. 91, no. 3, March, pp. 1052-6.

Ross, NE, Pritchard, CJ, Rubin, DM & Duse, AG 2006, 'Automated image processing method for the diagnosis and classification of malaria on thin blood smears', *Medical and Biological Engineering and Computing*, vol. 44, no. 5, April, pp. 427-36.

Ruberto, C, Dempster, A, Khan, S & Jarra, B 2000, 'Automatic thresholding of infected blood images using granulometry and regional extrema', IEEE, pp. 441-4.

Ruberto, C, Dempster, A, Khan, S & Jarra, B 2000, 'Segmentation of blood images using morphological operators', *Proceedings 15th International Conference on Pattern Recognition*, IEEE, Barcelona, Spain, pp. 397-400.

Sachs, J & Malaney, P 2002, 'The Economic and Social Burden of Malaria', *Nature*, vol. 415, pp. 680-5.

Savkare, SS & Narote, SP 2011, 'Automatic Detection of Malaria Parasites for Estimating Parasitemia', *International Journal of Computer Science and Security*, vol. 5, no. 3, pp. 310-5.

Serra, J 1984, *Image Analysis and Mathematical Morphology*, Academic Press, Orlando, USA.

Soille, P 2013, *Morphological image analysis: principles and applications*, Springer Berlin Heidelberg, New York, USA.

Sunarko, B, Williams, S, Prescott, WR, Byker, SM & Bottema, MJ 2013, 'Comparative study of two methods for blood cell segmentation', *Proceeding 2nd Engineering International Conference*, Semarang State University, Semarang, Indonesia, pp. II-96-100.

—— 2017, 'Correlation between automatic detection of malaria on thin film and experts' parasitaemia scores', *AIP Conference Proceedings*, AIP Publishing, Semarang, Indonesia, pp. 020054-1-10.

Suradkar, PT 2013, 'Detection of Malarial Parasite in Blood Using Image Processing', *International Journal of Engineering and Innovative Technology*, vol. 2, no. 10, April 2013, pp. 124-6.

Tanizaki, R, Kato, Y, Iwagami, M, Kutsuna, S, Ujiie, M, Takeshita, N, Hayakawa, K, Kanagawa, S, Kano, S & Ohmagari, N 2014, 'Performance of rapid diagnostic tests for Plasmodium ovale malaria in Japanese travellers', *Tropical medicine and health*, vol. 42, no. 4, August, p. 149.

Tek, FB, Dempster, AG & Kale, I 2006, 'Malaria Parasite Detection in Peripheral Blood Images', *BMVC06 Proceedings*, Edinburg, UK, pp. 347-56.

Tek, FB, Dempster, AG & Kale, I 2009, 'Computer vision for microscopy diagnosis of malaria', *Malar J*, vol. 8, July, p. 153.

—— 2010, 'Parasite detection and identification for automated thin blood film malaria diagnosis', *Comput Vision and Image Understanding*, vol. 114, August, pp. 21-32.

Thein, S 2001, 'Hematology–Landmark Papers of the Twentieth Century', *British Journal of Haematology*, vol. 114, no. 3, September, pp. 736-7.

Toha, SF & Ngah, UK 2007, 'Computer aided medical diagnosis for the identification of malaria parasites', *2007 International Conference on Signal Processing, Communications and Networking*, IEEE, Chenmai, India, pp. 521-2.

Trape, JF 1985, 'Rapid evaluation of malaria parasite density and standardization of thick smear examination for epidemiological investigations', *Transactions of the Royal Society of Tropical Medicine and Hygiene*, vol. 79, no. 2, pp. 181-4.

Tukey, J 1974, 'Nonlinear (nonsuperposable) methods for smoothing data', *Congr. Rec. 1974 EASCON*, vol. 673.

Vander, AJ, Luciano, D & Sherman, J 2001, *Humanphysiology: the mechanisms of body function*, 8th edn, McGraw-Hill Higher Education, Boston, US.

Vincent, L 1993a, 'Grayscale area openings and closings, their efficient implementation and applications', *Proceedings EURASIP Workshop on Mathematical Morphology and its Applications to Signal Processing*, Barcelona, Spain, pp. 22-7.

—— 1993b, 'Morphological grayscale reconstruction in image analysis: applications and efficient algorithms', *IEEE Transactions on Image Processing*, vol. 2, no. 2, April, pp. 176-201.

Wheater, PR, Burkitt, HG & Daniels, VG 1979, *Functional histology. A text and colour atlas*, Churchill Livingstone, Edinburgh, Scotland.

WHO 2010, *Basic Malaria Microscopy Part I: Tutor's Guide*, 2nd edn, World Health Organization, Geneva, Switzerland.

—— 2014a, *Technical consultation to update the WHO Malaria microscopy quality assurance manual*, World Health Organization, Geneva, Switzerland.

—— 2014b, *World malaria report 2014*, World Health Organization, Geneva, Switzerland.

—— 2016, *World Malaria Report 2016*, World Health Organization, Geneva, switzerland.

Yadollahi, M & Procházka, A 2011, 'Image Segmentation for Object Detection', *Proceedings of the 19th International Conference Technical Computing*, Prague, Czech Republic.

Yang, M-H & Ahuja, N 1998, 'Gaussian mixture model for human skin color and its applications in image and video databases', *Storage and Retrieval for Image and Video Databases*, SPIE, San Jose, USA, pp. 458-66.

Young, IT, Verbeek, P & Mayall, BH 1986, 'Characterization of chromatin distribution in cell nuclei', *Cytometry*, vol. 7, no. 5, April, pp. 467-74.

Zitová, B & Flusser, J 2002, 'Estimation of camera planar motion from blurred images', *Proceedings IEEE 2002 International Conference on Imaging Processing*, IEEE, New York, USA, pp. II-329-32.

Zivkovic, Z 2004, 'Improved adaptive Gaussian mixture model for background subtraction', *Proceedings of the 17th International Conference on Pattern Recognition*, IEEE, Cambridge, UK, pp. 28-31.