**Dictionary using Text Recognition for Mobile App**

by

Muhammad Hanif B Faisal

Dissertation submitted in partial fulfilment of

the requirement for the

Bachelor of Technology (Hons)

Business Information System

JANUARY 2014

Universiti Teknologi PETRONAS

Bandar Seri Iskandar

31750 Tronoh

Perak Darul Ridzuan

CERTIFICATION OF APPROVAL

**Dictionary using Text Recognition for Mobile App**

by

Muhammad Hanif B Faisal

A project dissertation submitted to the

Business Information System Programme

Universiti Teknologi PETRONAS

in partial fulfillment of the requirement for the

BACHELOR OF TECHNOLOGY (Hons)
(BUSINESS INFORMATION SYSTEM)

Approved by,

_____

Dr. Jafreezal B Jaafar

UNIVERSITI TEKNOLOGI

PETRONAS TRONOH, PERAK

January 2014

CERTIFICATION OF
ORIGINALITY

This is to certify that I am responsible for the work submitted in this project, that the original work is my own except as specified in the references and acknowledgements, and that the original work contained herein have not been undertaken or done by unspecified sources or persons.

_____

MUHAMMAD HANIF B FAISAL

# ABSTRACT

Dictionary is reference wellspring of words in a dialect or order, organized sequentially. Notwithstanding characterizing the words, bigger lexicons additionally give data on the spellings, elocution, word sources (historical background), capacities, and diverse types of the word. Through modern years, technology has rapidly grown to serve humanity with better goods, hence dictionary should also abide to the growth of the technology. E-dictionary or online dictionary has served well for past few years since the introduction of smartphone which led to another possibilities and remove limitations of old age dictionary. With an Artificial Intelligence (AI) making debut as a new trend in passing years, integrating the AI technology into dictionary is the purpose of this study. Although, e-dictionary has solved most of readers problem which is searching the words manually rather consuming times, the search function in most advance e-dictionary prompt to human error such as mistyping and consuming time and limitation of input of language (i.e. Japanese, Mandarin, Thai, Arabic etc.) hence, the development of new kind of dictionary which integrate the AI technology using Android platform to solve these problems. Android is preferred since almost half of the world is using smartphones powered by operating system called Android. Android SDK will be used to develop this application integrating with the Optical Character Recognition (OCR) technology, in other word is the AI technology in more specific term. Tesseract which is being maintained by Apache is the most accurate OCR software available. Integrating these two platforms is the challenge for this application. For translation note, Google Translate API will be used hence another integration will be done. Google Translate API is chosen since Google has updated several of languages at the time of this study.

# ACKNOWLEDGEMENTS

First and foremost, I would like to extend my gratitude to my dedicated supervisor, A.P. Dr. Baharom b Baharudin and Dr Jafreezal B Jaafar who has given me his full support and guidance throughout my completion of Final Year Project. It has been a great privilege for me to be supervised by a helpful and experienced lecturer.

Special thanks are given to Universiti Teknologi PETRONAS and also to Computer and Information Sciences Department (CIS) who have put a lot of effort towards contributing in making the final year project course a great success and also for their generosity in providing me all the facilities and technical expertise in making my project become successful.

All of these also will not happen without my family members and friends who have been supporting me throughout completion of my Final Year Project. All of their efforts in ensuring the success and completion of this project are really appreciated.

# Table of Contents

# LIST OF FIGURES

# LIST OF ABBREVIATIONS

| Abby | Definitions |
|------|-------------|
| ADT  | Android Development Tool |
| AI   | Artificial Intelligence |
| API  | Application Programming Interface |
| IDE  | Integrated Development Environment |
| NDK  | Native Development Kit |
| NLP  | Natural Language Processing |
| OCR  | Optical Character Recognition |
| SDK  | System Development Kit |

# CHAPTER 1: INTRODUCTION

## 1.1 Background

A dictionary likewise called a word-stock, word reference, wordbook, lexicon, or vocabulary is an accumulation of expressions in one or more particular languages, frequently recorded alphabetically, with use informative content, definitions, historical backgrounds, phonetics, articulations, and other informative data or a book of expressions in one language with their equivalents in a different one, otherwise called a dictionary. As per Nielsen (2008) a dictionary may be viewed as a lexicographical item that is characterised by three noteworthy characteristics, it has been ready for one or more capacities, it holds information that have been chosen with the end goal of satisfying those capacities and its lexicographic structures join and create relationships between the information with the intention that they can help clients and fulfil the capacities of the dictionary.

Dictionary can be utilized in several manners such as finding the spelling and significance of expressions, pronunciation/audio, curved manifestations of expressions, underwriting, historical background, word division, relevant uses and significantly all the more, contingent on the lexicon. As well as expanding user vocabulary knowledge, composing articles and increasing proficiency for the particular language.

In this advanced era, books still managed to hold their throne in the heart of book enthusiasts as a tool to gain knowledge, expanding creativity, and playing with imagination while creating their own scenery according to the writer's plot. In order to achieve this imaginative desire to the fullest, a complete comprehension on the literature is a must, hence barrier such as different language and lack of vocabulary will set a limitation on one's imagination.

Dictionary Using Text Recognition for Mobile Application will set a new form of trend among the users since this innovation enhance the practical use of an ordinary dictionary. For years, "flipping through pages" has dominating the world which contributes a lot of knowledge through manual searching and few years back, the world is being hit with a technology so called gadget. These gadgets that can be in the form of tablets and smartphones have changed the practicality of dictionary. Being known as e-dictionary application, it can be run on these tablets and smartphones which can be use as any other ordinary dictionary. Technology is improving life quality, instead of heavy traditional dictionary, this application weight equally to the weight of the gadget that bearer it and can be used like any ordinary dictionary with new features such as searching feature. This current application has served its purpose for this era, Dictionary Using Text Recognition is the next future of dictionary. By scanning the word through a medium such as camera (most gadgets have built in camera as standard feature), this application can recognize the word and automatically perform searching feature resulting in definition of the word will be display almost instantaneously.

## 1.2 Problem Statement

### Heavy

Ordinary dictionary which physically can be seen as ordinary book has mass. This mass or in common word called weight is one of major problem of a dictionary. A dictionary may consist up to 400 pages which make this book thick enough to be big as well as heavy. A pocket dictionary may have solved this problem but it hinders the purpose of dictionary since most of it consisted of common words or in different term, it is not complete. For a user to bring a dictionary on their routine life will be a burden, hence this Dictionary Using Text Recognition will be the solution since most user used a smartphone that this application can be run on.

**Mistyping**

Humans are prompt to human error, in this case, it is mistyping. A lot of factors contribute to mistyping where the relevant lists can be found below:

- Poor eyesight
- Incompatibility of user's finger and hardware size
- Typing too fast
- A new word to user usually hard to spell

These factors are leading to mistyping where in today's world, e-dictionary still need typing as one of the requirement to use it. With mistyping, a different word might emerge resulting in different totally different words and meaning.

## 1.3 Objectives and Scopes of Study

Objectives of this project are listed below:

- To dispose the unnecessary weight of common dictionary
- To reduce mistyping by using OCR technology
- To collaborate dictionary with gadgets such as smartphones and tablets
- To improve conveniences of using an e-dictionary

Scopes of study of this project will include:

- To learn on development of Dictionary using Text Recognition for Mobile Apps
- To research on OCR (Optical Character Recognition)
- To understand the concept of OCR
- To research on integration on Android platform and OCR technology

# CHAPTER 2: LITERATURE REVIEW

Dictionary Using Text Recognition is using Android based - platform in the form of application that can be run by any gadget running this OS. Hence a walkthrough of Android, Artificial Intelligence, Natural Language Processing, Text Mining and OCR need to be understood first. Questions that need to be covered are listed below:

- What is Android?
- What is Artificial Intelligence?
- What is Natural Language Processing?
- What is Text Mining?
- What is Optical Character Recognition?
- What platform needed to program Android's application?
- What platform can support Optical Character Recognition and Android to work simultaneously?

## 2.1 What is Android?

Android is a Linux-based working framework composed basically for touchscreen portable apparatuses, for example cell phones and tablet machines. At first advanced by Android, Inc., which Google sponsored monetarily and later purchased in 2005, Android was disclosed in 2007 in addition to the establishing of the Open Handset Alliance, a consortium of equipment, programming, and telecommunication organizations dedicated to propelling open norms for portable mechanisms. The predominant Android-powered telephone was sold in October 2008.

Android is open source and Google discharges the code under the Apache License. This open source code and tolerant permitting permits the programming to be openly changed and circulated by mechanism producers, remote bearers and fan designers. Also, Android has a huge neighbourhood of planners composing provisions ("applications") that amplify the practicality of smartphone, composed basically in a redid form of the Java customizing

dialect. In October 2012, there were more or less 700,000 applications accessible for Android, and the assessed number of requisitions downloaded from Google Play, Android's essential application store, was 25 billion.

These elements have helped towards making Android the planet's generally substantially utilized cell phone stage, overwhelming Symbian in the final quarter of 2010, and the programming of decision for engineering organizations who require an ease, customizable, lightweight working framework for high tech mechanisms without advancing one starting with no outside help. Accordingly, notwithstanding being basically intended for telephones and tablets, it has seen extra requisitions on Tvs, amusements supports, computerized Polaroids and different hardware. Android's open nature has further empowered a substantial group of designers and devotees to utilize the open source code as an establishment for neighborhood driven ventures, which include new characteristics for progressed clients or carry Android to mechanisms which were authoritatively discharged running other working frameworks.

## 2.2 What is Artificial Intelligence (AI)?

Regarding making complex informed decisions, computer cannot trade individuals. Anyway with artificial intelligence, computer could be prepared to think like people do. Artificial intelligence permits computer to gain experience as a matter of fact, distinguish examples in a lot of complex information and settle on complex choices dependent upon human learning and thinking aptitudes. Artificial intelligence has turned into an imperative field of study with a boundless of requisitions in fields running from pharmaceutical to agriculture.

**Expert Systems**

Two of the most important and most used branches of AI are neural networks and expert systems.

An expert system can illuminate true issues utilizing human information and accompanying human thinking abilities. Information and thinking procedures of experts are gathered and encoded into a learning base. Starting there on, the expert system could swap or support the

human experts in settling on complex choices by coordinating all the learning it has in its information base.
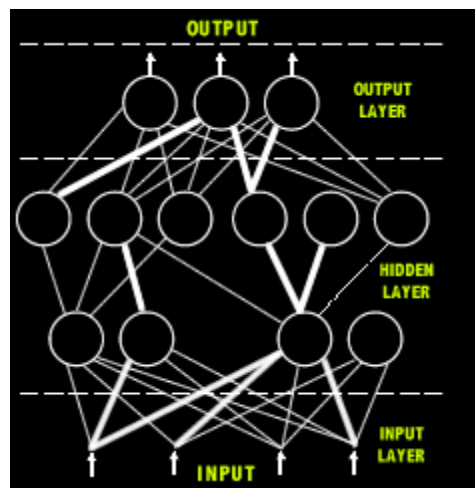
**Neural Networks**



Figure 1

Delineation of Neural Network. This chart speaks of an artificial neural network. A neural network is made of junctions masterminded in diverse examples speaking to the "intelligence" of the network. The line thickness shows the quality of the connections.

The most significant requisition of neural networks is in example distinguishment. People, through neurons in their brains, study how to read human composing, distinguish a terrible fruit from a great one or recognize their youngsters from a set of children. Neural networks permit workstations to utilize the same standards that neurons in the brains use to distinguish and group distinctive examples. So in a manner, neural networks are a

computerized representation (in spite of the fact that very streamlined) of our brains. They are made of artificial neurons, associated by weights, which are demonstrative of the qualities of the associations. The neurons are organized in layers, and hinging upon the multifaceted nature of the requisition, there could be a couple of them or an extremely vast number of them (hundreds or thousands). Iterative engendering of data from one layer of neurons to the following (preparing) is the thing that empowers the neural network to gain experience for a fact.

Unlike people, when a neural is completely prepared, it can order and recognize examples in enormous measures of complex information. It could do this at high speeds that cannot be doubled by people.

## 2.3 What is Natural Language Processing?

Natural language processing (NLP) is the capability of a computer program to grasp human discourse as it is spoken. NLP is a part of artificial intelligence (AI).

The improvement of NLP provisions is testing in light of the fact that computer customarily require people to "talk" to them in a customizing language that is exact, unambiguous and exceedingly organized or, maybe through a set number of obviously articulated voice summons. Human discourse, then again, is not dependably exact --it is frequently vague and the semantic structure can rely on upon numerous complex variables, incorporating slang, territorial tongues and social connection.

Current methodologies to NLP are dependent upon machine studying, a sort of artificial intelligence that analyses and uses examples in information to enhance a program's own particular comprehension. The vast majority of the exploration being carried out on natural language processing spins around hunt, particularly undertaking search.

Regular NLP undertakings in programming systems today include:

- Sentence segmentation, part-of-speech tagging and parsing.
- Deep analytics.
- Named entity extraction.
- Co-reference resolution.

The point of interest of natural language processing could be seen when acknowledging the accompanying two comments: "Cloud computing protection ought to be part of each service level agreement" and "A great SLA guarantees a less demanding night's slumber --even in the cloud." If you utilize national language processing for hunt, the project will recognise that cloud computing is an entity, that cloud is a contracted manifestation of cloud computing and that SLA is an industry acronym for service level agreement.

## 2.4 What is Text Mining?

Text mining is the dissection of information held in natural language text. The requisition of text mining strategies to tackle business issues is called text analytics.

Text mining can help a conglomeration infer possibly important business experiences from text-based substance, for example word reports, message and postings on social media streams like Facebook, Twitter and Linkedin. Mining unstructured information with natural language processing (NLP), factual demonstrating and machine studying procedures might be challenging, in any case, since natural language text is regularly conflicting. It holds ambiguities brought about by conflicting sentence structure and semantics, incorporating slang, language particular to vertical businesses and age assemblies, two sided sayings and mockery.

Text analytics programming can help by transposing words and phrases in unstructured information into numerical qualities which can then be connected with organized

information in a database and dissected with universal information mining systems. With an iterative approach, a conglomeration can solidly utilize text analytics to increase knowledge into substance particular qualities, for example assessment, feeling, power and importance. On the grounds that text analytics innovation is still acknowledged to be a rising technology, be that as it may, comes about and profundity of examination can differ uncontrollably from vendor to vendor.

## 2.5 What is Optical Character Recognition?

Assume the necessity to digitize a magazine article or a printed contract. Spending hours retyping and afterward rectifying misprints is not efficient. Alternately, this can change over all the needed materials into advanced configuration in some minutes utilizing a scanner (or a digital camera) and Optical Character Recognition software.

## What exactly is meant by OCR?

OCR (optical character recognition) is the recognition of printed or composed content characters by a computer. This includes photo scanning of the content character-by-character, examination of the checked in picture, and afterward interpretation of the character picture into character codes, for example ASCII, usually utilized within information preparing.

In OCR transforming, the examined in picture or bitmap is investigated for light and dull regions keeping in mind the end goal to recognize every alphabetic letter or numeric digit. The point when a character is distinguished, it is changed over into an ASCII code. Exceptional circuit sheets and computer chips planned explicitly for OCR are utilized to accelerate the recognition procedure.

OCR tends to be utilized by libraries to digitize and save their property. OCR is additionally used to process checks and charge card slips and sort the mail. Billions of magazines and letters are sorted each day by OCR machines, respectably accelerating mail conveyance.

Optical Character Recognition, or OCR, is an engineering that empowers to change over distinctive sorts of records, for example checked paper reports, PDF documents or pictures caught by a digital camera into editable and searchable information.

Envision a paper report -for instance, magazine article, hand-out, or PDF contract an accomplice sent via message. Clearly, a scanner is insufficient to make this qualified information accessible for altering, say in Microsoft Word. Every one of the scanner can do is make a picture or a preview of the report that is nothing more than a gathering of dark and white or colour specks, regarded as a raster picture. To concentrate and repurpose information from examined reports, camera pictures or image-only PDFs, you need OCR software that might single out letters on the picture, put them into expressions and afterward -statements into sentences, therefore empowering to enter and alter the substance of the original document.

**What Technology lies behind OCR?**

The careful instruments that permit people to distinguish protests are yet to be grasped, however the three essential standards are now well known by researchers – integrity, purposefulness and adaptability (IPA). These standards constitute the centre of OCR permitting it to repeat common or human-like recognition.

Would not it be great if how OCR distinguishes content can be explained. First and foremost, the project breaks down the structure of record picture. It partitions the page into components, for example squares of writings, tables, pictures, and so forth. The lines are separated into expressions and afterward -into characters. When the characters have been singled out, the system contrasts them and a set of example pictures. It developments

various theories about what this character is. Basing on these speculations the system dissects diverse variants of breaking of lines into statements and expressions into characters. In the wake of transforming colossal number of such probabilistic theories, the project at long last takes the choice, displaying you the distinguished text.

## 2.6 What platform needed to program Android's application?

There are few platforms that can be used to code Android's application such as AppsInventer, Android SDK (Software Development Kit) usually integrated with Eclipse for it IDE and Android Studio which is based on IntelliJ IDEA.

### Android SDK

Android software development is the methodology by which new applications are made for the Android working framework. Applications are generally advanced in the Java modifying dialect utilizing the Android Software Development Kit, however other development tools are accessible. As of October 2012, more than 700,000 applications have been produced for Android, with in excess of 25 billion downloads.

The Android software development kit (SDK) incorporates a thorough set of development tools. These incorporate a debugger, libraries, a handset emulator dependent upon QEMU, documentation, sample code, and tutorials. As of now underpinned development platforms incorporate workstations running Linux (any present day desktop Linux circulation), Mac OS X 10.5.8 or later, Windows XP or later; for the minute one can advance Android software on Android itself by utilizing [aide -Android IDE -Java, C++] application and [android java editor] application. The authoritatively underpinned reconciled nature's turf (IDE) is Eclipse utilizing the Android Development Tools (ADT) Plugin, however Intellij IDEA IDE (all releases) completely underpins Android development out of the container, and Netbeans IDE additionally backs Android development through a plugin. Moreover,

20

developers might utilize any content manager to alter Java and XML indexes, then use charge line tools (Java Development Kit and Apache Ant are needed) to make, fabricate and debug Android applications and in addition control appended Android units (e.g., triggering a reboot, instituting software package(s) remotely).

Improvements to Android's SDK run as an inseparable unit with the generally speaking Android platform development. The SDK additionally backs more seasoned forms of the Android platform in the event that visionaries wish to focus on their applications at more seasoned units. Development tools are downloadable segments, so after one has downloaded the most recent variant and platform, more seasoned platforms and tools can additionally be downloaded for similarity testing.

Android applications are bundled in .apk group and archived under /data/app organizer on the Android OS (the envelope is receptive just to the root client for security explanations). APK bundle holds .dex records (incorporated byte code documents called Dalvik executables) and resource files.

**Android NDK**

The NDK is a toolset that allows user to implement parts of app using native-code languages such as C and C++. For certain types of apps, this can be helpful so user can reuse existing code libraries written in these languages.

**Eclipse**

In computer programming, Eclipse is a multi-language Integrated development environment (IDE) involving a base workspace and an extensible fitting in framework for altering the environment. It is composed basically in Java. It could be utilized to improve requisitions in Java and, by method of different plug-ins, other programming language incorporating Ada,

C, C++, COBOL, Fortran, Haskell, Javascript, Perl, PHP, Python, R, Ruby (counting Ruby on Rails framework), Scala, Clojure, Groovy, Scheme, and Erlang. It can additionally be utilized to advance bundles for the programming Mathematica. Development environments incorporate the Eclipse Java development instruments (JDT) for Java and Scala, Eclipse CDT for C/c++ and Eclipse PDT for PHP, around others.

The beginning codebase started from IBM Visualage. The Eclipse programming development unit (SDK), which incorporates the Java development apparatuses, is implied for Java designers. Clients can augment its capabilities by instating fitting ins composed for the Eclipse Platform, for example development tool compartments for other modifying dialects, and can compose and help their own particular attachment in modules.

Discharged under the terms of the Eclipse Public License, Eclipse SDK is free and open source programming (in spite of the fact that it is incongruent with the GNU General Public License). It was one of the first Ides to run under GNU Classpath and it runs without issues under Icedtea.

**Android Studio**

Android Studio is another Android development environment dependent upon Intellij IDEA. Comparable to Eclipse with the ADT Plugin, Android Studio gives mixed Android developer tools for development and debugging. On top of the abilities you need from Intellij, Android Studio offers:

• gradle-based form help.

• android-particular refactoring and fast fixes.

• lint instruments to get execution, ease of use, adaptation similarity and different issues.

• Proguard and application marking proficiencies.

• template-based wizards to make regular Android plans and parts.

**Google Translator API**

Google Translate is a tool that automatically translates text from one language to another language (e.g. French to English). Programmer can use the Google Translate API to programmatically translate text in their webpages or apps.

**2.7 What platform can support Optical Character Recognition and Android to work simultaneously?**

**Tesseract**

Tesseract is an optical character distinguishment motor for different working frameworks. It is free programming, discharged under the Apache License, Version 2.0, and advancement has been supported by Google since 2006. Tesseract is viewed as a standout amongst the most precise open source OCR motors presently accessible.

Tesseract was in the main three OCR motors regarding character precision in 1995. It is accessible for Linux, Windows and Mac OS X, be that as it may, because of restricted assets just Windows and Ubuntu are thoroughly tried by designers.

Tesseract up to and including form 2 could just acknowledge TIFF pictures of straightforward one section message as inputs. These early forms finished not incorporate format dissection thus inputting multi-ordered content, pictures, or mathematical statements generated a jumbled yield. Since adaptation 3.00 Tesseract has upheld yield content organizing, hocr positional data and page format investigation. Help for various new picture arrangements was included utilizing the Leptonica library. Tesseract can catch whether content is monospaced or relative.

The beginning forms of Tesseract could just distinguish English dialect content. Beginning with form 2 Tesseract could handle English, French, Italian, German, Spanish, Brazilian Portuguese and Dutch. Beginning with adaptation 3 it can distinguish Arabic, English, Bulgarian, Catalan, Czech, Chinese (Simplified and Traditional), Danish, German (standard and Fraktur script), Greek, Finnish, French, Hebrew, Croatian, Hungarian, Indonesian, Italian, Japanese, Korean, Latvian, Lithuanian, Dutch, Norwegian, Polish, Portuguese, Romanian, Russian, Slovak (standard and Fraktur script), Slovenian, Spanish, Serbian, Swedish, Tagalog, Thai, Turkish, Ukrainian and Vietnamese. Tesseract could be prepared to work in different dialects as well.

On the off chance that Tesseract is utilized to process right-to-left content such Arabic or Hebrew the outcomes are requested just as it is left-to-right content.

Tesseract is suitable for utilization as a backend, and might be utilized for additional muddled OCR errands including design examination by utilizing a frontend, for example, Ocropus. Tesseract deals with Linux, Windows (with Vc++ Express or Cygwin) and Mac OSX. It can additionally be assembled for different stages, including Android and the iphone, however these are not too tried stages.

# CHAPTER 3: METHODOLOGY

**3.1 Project Flow Chart**

**Literature Review**
- Research on journal related to this project.
- Understanding the scope of study of project.

**Learning**
- Understand the development process of project
- Learning various platforms that can developed this project

**Data Collection**
- Data mining on the internet regarding OCR
- Forums, to get more information based on user experiences.

**Data Analysis**
- Collected data are used and analyze to create the development of this project.

**Improvement**
- Need more information regarding OCR and Android integration.

Figure 2

## 3.2 Tools & Equipment

- Asus S550C – device used to complete report and project
- Microsoft Word – report writing
- Android SDK with Eclipse IDE – Testing development kit
- Android NDK
- HTC One Max
- Command Prompt
- Tesseract Library

## 3.3 System Architecture



Figure 3: Architecture

The application will run and access the camera, using the camera, the text will be scanned and Tesseract will perform the algorithmic calculation to recognize the characters before the data is being processed and transfer to search function which is integrate with the Google

Translator API. Result will be shown after the search is done in the form of translation preferred by the user.

## 3.4 Model Framework

| Application |
|---|
| Dictionary Text Recognition · OCR Platform |

| Application Framework (Smartphone) |
|---|
| Task Manager · Telephony Manager · View System |

| Libraries | Android Runtime |
|---|---|
| SQLite · WebKit · SGL · SSL | Core Libraries · Dalvik Virtual Machine |

| Linux Kernel |
|---|
| Display Driver · Camera Driver · Flash Memory Driver · Keypad Driver · Power Management · WiFi Driver |

Figure 4: Model Framework

**3.5 Requirement Analysis & Specification**

**User**

This application will be used by any users who can comprehend English as well as Malay and if we taking it to a higher level, it will be used by bookworms whose using dictionary as their medium of translating new words. This application might as well being used by students who newly try to learn English. It will become an interesting application among children applicable with the current trend that this generation is a generation Information Technology (IT).

**Functions**

- Input can be done using manual input and OCR.
- Scan text in any document which automatically being detect by the application.
- Search function of the dictionary (manually search still applicable).
- Display the result of input in the form of explanation and definition (both input).
- Automatic detection of the language input (OCR functionality).

**Limitation**

- It will not be a real-time basis OCR, hence reducing user friendly practicality.
- Dictionary provided will be English – Malay and vice versa.
- Accuracy of OCR depends on the OCR platform being used.
- Font's recognition limitation.

## 3.6 Gant Chart

| No. | Project Activities (FYP1) | Week | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| 1. | **Selection of Project Title** | ▓ | ■ | | | | | | | | | | | | |
| | Search for Project Title | ▓ | ▓ | | | | | | | | | | | | |
| 2. | **Planning & Research Analysis** | | | ▓ | ▓ | ▓ | ▓ | ▓ | ▓ | ▓ | ■ | | | | |
| | Literature review research | | | ▓ | ▓ | ▓ | | | | | | | | | |
| | Define application scope | | | | | ▓ | ▓ | ▓ | | | | | | | |
| | Determine application outline | | | | | | | ▓ | ▓ | ▓ | | | | | |
| | Testing various platform | | | | | | | ▓ | ▓ | ▓ | ▓ | | | | |
| 3. | **User Design** | | | | | | | | | | | ▓ | ▓ | ▓ | ■ |
| | Design user interface | | | | | | | | | | | ▓ | ▓ | | |
| | Preliminary application design | | | | | | | | | | | | ▓ | ▓ | ▓ |

▓ Process
■ Suggested Milestone

| No. | Project Activities (FYP2) | Week | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| 4. | **System Construction** | ▓ | ▓ | ▓ | ▓ | ▓ | ▓ | ▓ | ▓ | ▓ | ▓ | ■ | | | |
| | • Build | ▓ | ▓ | ▓ | ▓ | ▓ | | | | | | | | | |
| | Develop User Interface | ▓ | ▓ | | | | | | | | | | | | |
| | Write coding | | | ▓ | ▓ | | | | | | | | | | |
| | Integrate OCR into application | | | | ▓ | ▓ | | | | | | | | | |
| | • Demonstrate | | | | | | ▓ | ▓ | ▓ | | | | | | |
| | Run simple test to show the workability | | | | | | ▓ | ▓ | | | | | | | |
| | Ensure all components interrelated and working | | | | | | | ▓ | ▓ | | | | | | |
| | • Refine | | | | | | | | | ▓ | ▓ | ▓ | | | |
| | Debug | | | | | | | | | ▓ | ▓ | | | | |
| | Reconstruct the system | | | | | | | | | | ▓ | ▓ | | | |
| 5. | **System Cutover** | | | | | | | | | | | | ▓ | ▓ | ■ |
| | Testing system functionality and usability | | | | | | | | | | | | ▓ | ▓ | |
| | Check system specification aligned with requirements | | | | | | | | | | | | | ▓ | |
| | System implementation | | | | | | | | | | | | | | ▓ |

▓ Process
■ Suggested Milestone

Figure 5: Project's Gant Chart

## 3.7 Experimental Setup

In order to build this application, various ways are possible which include using the official Android Development platform called Android Studio or AppsInventer. Instead of these two platforms, Eclipse IDE and Android SDK are chosen since Android Studio is still new hence integrating it with other platform might be impossible at this time of study. AppsInventer has limitation since it is basically to develop a simple Android application hence making it out of equation.

## Setup the Android SDK and Eclipse IDE

There are two possible ways of setup the Eclipse IDE with Android SDK. First, is by manually downloading the Eclipse IDE from https://www.eclipse.org/downloads/ .



Figure 6: Eclipse Download

After that, go to http://developer.android.com/sdk/index.html to download the Android SDK. This Android SDK is provided by Google for developers.



Figure 7: Android SDK Download

Install both software after they have finish downloading. Follow below step to setup the Android SDK in Eclipse IDE:

1. Start Eclipse, then select **Help** > **Install New Software**.

2. Click **Add**, in the top-right corner.

3. In the Add Repository dialog that appears, enter "ADT Plugin" for the *Name* and the following URL for the *Location*:

   https://dl-ssl.google.com/android/eclipse/

   **Note:** The Android Developer Tools update site requires a secure connection. Make sure the update site URL you enter starts with HTTPS.

4. Click **OK**.

5. In the Available Software dialog, select the checkbox next to Developer Tools and click **Next**.

6. In the next window, you'll see a list of the tools to be downloaded. Click **Next**.

7. Read and accept the license agreements, then click **Finish**.

   If you get a security warning saying that the authenticity or validity of the software can't be established, click **OK**.

8. When the installation completes, restart Eclipse.

The second way is much simpler which is to download the ADT (Android Development Tools) Bundle which is available at http://developer.android.com/sdk/index.html .


Figure 8: ADT Bundle Download

Extract the ADT Bundle in your C:\ and run your Eclipse IDE. This ADT Bundle has already integrated both Android SDK and Eclipse.

**Setup the Tesseract Library/Platform**

Download the Tesseract Library at https://code.google.com/p/tesseract-ocr/downloads/list .



Figure 9: Tesseract Download

In order to use the Tesseract as a Library in Eclipse IDE, the Tesseract must be build first using Android NDK. Since Tesseract is built using C/C++ as its native language, Linux is the most preferable Operating System to build the Tesseract Library. There are three ways to set up the Tesseract Library:

a. Using Command Prompt

Since the Tesseract must be built on a Linux operating system, *ndk-build* command will not be recognize by the command prompt. Hence a few modifications to the Windows Enviroment Variables must be done. Go to *Control Panel > System and Security > System > Advanced system settings > Environment Variables* under *System variables* tab choose *Path* and click *Edit*. Add a path of the Android NDK folder (whichever path user extract it to) i.e. *C:\Android\android-ndk*.

Figure 10: Environment Variables

Now to build the Tesseract Library, open the *cmd.exe* using *Administrator* privilege.
Changed the path of the directory to path of folder Tesseract i.e. *cd C:\Tesseract.*



Figure 11.1: cmd.exe

Run the *ndk-build* command by pressing *Enter*.

Figure 11.2: cmd.exe

Wait till it finishes building the Tesseract library like shown below.



Figure 11.3: cmd.exe

Import the Tesseract library into Eclipse workplace. It is now ready to be code.

b. Using Ubuntu operating system

Open up the terminal in the Ubuntu system.

Figure 12: Terminal

Type the following command which will build the Tesseract Library:

*cd tess*

*cd tess-two*

*ndk-build*

*android update project --path .*

*ant release*

Note: The directory in the command above is depends on user preference of Tesseract Library directory.

Wait until the build is finish and import it into Eclipse workplace. Tesseract is ready to use.

c. Using Eclipse IDE

Run Eclipse and *import* Tesseract as *General Project* into *Workplace*. Right-click the Tesseract project and click *Properties*. Click *Builder* tab and create *New* to create a NDK Builder using Eclipse.

Figure 13.1: NDK Builder

Click *Program > Ok* and a new window will pop up. At the *Name* put *NDK Builder*. Under *Location*, insert the directory of Android NDK and search for *ndk-build.exe*. For *Working Directory*, insert the directory of Tesseract.



Figure 13.2: NDK Builder

Press *OK* and *Refresh* the Tesseract library which will lead to building of the library. The library is ready to use.

**Setup the Google Translator API**

Google Translator API is a charged service which certain fees need to be paid in order to use it. In order to use in the application it first need to be identified to Google. To find application's API key, do the following:

a. Go to the Google Developers Console.
b. Select a project.
c. In the sidebar on the left, select APIs & auth. In the list of APIs, make sure the status is ON for the Google Translate API.
d. In the sidebar on the left, select Credentials.
   Note: Translate API does not have any methods that require OAuth 2.0 authorization. Instead, click on the "Create new key" button and generate a server, browser, Android, or iOS key depending on your application's needs.

The Console enables user to create server, browser, Android and iOS API Keys. Once user have created a key, it is extremely important that user secure their key by not disclosing or making it available in a public forum since anyone who has access to their key can use it to incur charges on their bill.

# Chapter 4: Result & Discussion

**4.1 Prototype**

Using Eclipse IDE which requires Java Programming, Dictionary using Text Recognition for Mobile App managed to be developed. This application which using Tesseract as its OCR AI to recognize text and Google Translator API to translate the data received from OCR engine managed to serve its purpose as AI Dictionary. Since it is using Google Translator API, this application will used internet hence making it an online useable application.



Figure 14.1: DictionaryTextRecognition

In Figure 14.1 the first interface is accessed when the application first start up. It automatically accessed the camera with a few function added. The first thing can be seen is the box. It can be adjust according to the text size. It purposes is to let user focus on the text that need to be recognize by the OCR engine. The camera button let the user to capture the text and the OCR engine will perform the text recognition as shown in Figure 14.2. It also will show the translation according to user preference if available.
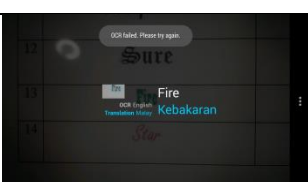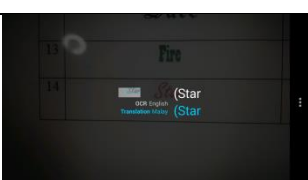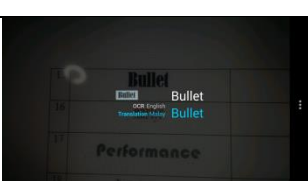
Figure 14.2: DictionaryTextRecognition

A continuous preview is also available which allow user to preview the text and making sure the text is being recognize by the OCR engine accurately. This function is useful since most OCR engine other than Tesseract is not accurately. It can be seen in at top left corner Figure 14.3 below.



Figure 14.3: DictionaryTextRecognition

A menu can be access by clicking "…" at the below of the application which will access the setting. In the setting, the continuous preview can be "off" or "on" based on user preference. It also will let user to choose the recognize text language and translation language. Figure 14.4 below will show the detail.



Figure 14.4: DictionaryTextRecognition

After the user has select the recognize language, OCR engine will update the data into the phone by downloading it through the internet. Figure 14.5 show the downloading of English

database into the phone and Figure 14.6 show the downloading of Malay database into the phone.



Figure 14.5: DictionaryTextRecognition



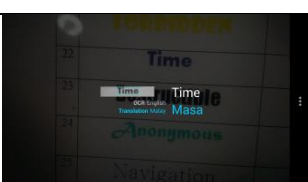Figure 14.6: DictionaryTextRecognition

Most of function in this application can be used since it is a prototype, the functions of Translator engine restricted to Google Translator API, it may use Bing Translator in near future. OCR engine available only Tesseract which in near future it can integrate with Cube engine.
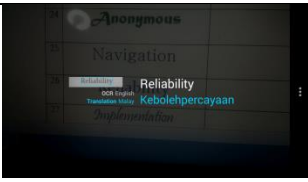
## 4.2 System Evaluation

In this part, a test case will be used to perform user testing for Dictionaries with Text Recognition for Mobile Application. There are two objectives of evolution need to be done. The first one is sensibility of OCR engine and secondly the accuracy of translation. In the table below, both criteria will be discussed in term of accepted and rejected.

| # | Text to Test | Result | Comment |
|---|---|---|---|
| 1 | Application |  | OCR : Accepted<br>Translation : Accepted |
| 2 | BODY |  | OCR : Accepted<br>Translation : Accepted |
| 3 | Searching |  | OCR : Rejected<br>Translation : Rejected |
| 4 | Document |  | OCR : Rejected<br>Translation : Rejected |
| 5 | Create |  | OCR : Accepted<br>Translation : Accepted |
| 6 | Believe |  | OCR : Rejected<br>Translation : Rejected |

| 7 | From |  | OCR : Accepted<br>Translation : Accepted |
|---|---|---|---|
| 8 | *Introduction* |  | OCR : Rejected<br>Translation : Rejected |
| 9 | Reference |  | OCR : Accepted<br>Translation : Accepted |
| 10 | Discover |  | OCR : Rejected<br>Translation : Rejected |
| 11 | *Request* |  | OCR : Accepted<br>Translation : Accepted |
| 12 | Sure |  | OCR : Accepted<br>Translation : Accepted |
| 13 | Fire |  | OCR : Accepted<br>Translation : Accepted |
| 14 | Star |  | OCR : Accepted<br>Translation : Accepted |
| 15 | **Bullet** |  | OCR : Accepted<br>Translation : Accepted |

| 16 | Martyr |  | OCR : Accepted<br><br>Translation : Rejected |
| --- | --- | --- | --- |
| 17 | Performance |  | OCR : Accepted<br><br>Translation : Accepted |
| 18 | Increase |  | OCR : Accepted<br><br>Translation : Accepted |
| 19 | Development |  | OCR : Rejected<br><br>Translation : Rejected |
| 20 | New |  | OCR : Accepted<br><br>Translation : Accepted |
| 21 | FORBIDDEN |  | OCR : Rejected<br><br>Translation : Rejected |
| 22 | Time |  | OCR : Accepted<br><br>Translation : Accepted |
| 23 | Destructible |  | OCR : Rejected<br><br>Translation : Rejected |
| 24 | Anonymous |  | OCR : Accepted<br><br>Translation : Accepted |

| 25 | Navigation |  | OCR : Accepted<br>Translation : Accepted |
|---|---|---|---|
| 26 | Reliability |  | OCR : Accepted<br>Translation : Accepted |
| 27 | *Implementation* |  | OCR : Rejected<br>Translation : Rejected |
| 28 | Apply |  | OCR : Accepted<br>Translation : Accepted |
| 29 | Success |  | OCR : Accepted<br>Translation : Accepted |
| 30 | *Graduate* |  | OCR : Accepted<br>Translation : Accepted |

## 4.3 Discussion on Findings/Results

Based on system evaluation, 30 samples of text consisting of different kind of font and colour reveal that the OCR engine managed to recognize 21 samples out of 30 samples. This showed that the sensibility of the OCR engine is high and perform well than expected. This is because most of the words are made up from rarely used fonts on daily basis. The font such as Sitka Small, Constantia and variety of font are rarely used in the daily basis making this application more applicable in real life situation.

The colours of the text took a role to affect the sensibility of the OCR engine. This is proved because colours with high brightness such as yellow, orange etc. tend to make the text unreadable by the OCR engine. In term of colours which are contrast made the text more readable and recognized by the OCR engine. This shows that this problem exist because of hardware limitation, in other term, the camera of a smartphone or the smartphone itself. Hence a better smartphone with an autofocus camera might detect such invisible colours which will helps the software, the OCR engine recognize the text.

In term of translation accuracy evaluation, most of the text are being well translated into Malay language since the preferences used is English to Malay conversion. It only fail to translate when it the OCR engine fail to recognize the word. Only one case, which is a higher level of word fail the translation test, Martyr is the one. This shows that, using Google Translator API managed to perform the basic function of a dictionary. Google always updating their resources hence a case like one sample failed out of 21 samples shows the application is running good to serve its purpose because it used Google Translator API real time translation.

This application also in term of weight it is basically based on user smartphone's weight. Hence, it will be handy to use rather than old fashioned dictionary which is thick and heavy.

# CHAPTER 5: CONCLUSION

## 5.1 Relevancy to the Objectives

This project will be developed to address the objectives stated. Generally, there are four (4) objectives to be achieved in this project:

- To dispose the unnecessary weight of common dictionary
- To reduce mistyping by using OCR technology
- To collaborate dictionary with gadgets such as smartphones and tablets
- To improve conveniences of using an e-dictionary

This project is relevant towards all objectives as being shown in Chapter 4 of this study.

## 5.2 Suggested Future Work for Expansion & Continuation

The future work for this project providing some suggestion for expansion and continuation in future work are:

- To integrate different OCR engine such as Cube to make text recognition more sensibility by combining two OCR engines.
- To integrate more translation API such as Bing Translator to make variety of words available.
- To combine the dictionary with different scanner such as QR Code and Bar Code to make it more practical in daily used.

## 5.3 Conclusion

As a conclusion, this application, Dictionary using Text Recognition for Mobile App meets all the four objectives stated earlier. This application now going through the evaluation phase based on Chapter 4. This application will go through a lot more correction and modification until it meet the requirement set earlier.

# REFERENCES

1. Dictionary - Wikipedia, the free encyclopedia. (n.d.). Retrieved June 30, 2013, from http://en.wikipedia.org/wiki/Dictionary

2. What is OCR and OCR Technology. (n.d.). Retrieved from http://finereader.abbyy.com/about_ocr/whatis_ocr/

3. Tesseract (software). (n.d.). In Wikipedia, the free encyclopedia. Retrieved July 10, 2013, from http://en.wikipedia.org/wiki/Tesseract_(software)

4. Android (operating system). (n.d.). In *Wikipedia, the free encyclopedia*. Retrieved July 10, 2013, from http://en.wikipedia.org/wiki/Android_(operating_system)

5. *What is Artificial Intelligence?* (n.d.). Retrieved from http://interests.caes.uga.edu/eai/ai.html

6. *What is natural language processing (NLP)? - Definition from WhatIs.com*. (n.d.). Retrieved from http://searchcontentmanagement.techtarget.com/definition/natural-language-processing-NLP

7. *What is text mining (text analytics)? - Definition from WhatIs.com*. (n.d.). Retrieved from http://searchbusinessanalytics.techtarget.com/definition/text-mining

8. Eclipse (software). (n.d.). In *Wikipedia, the free encyclopedia*. Retrieved July 10, 2013, from http://en.wikipedia.org/wiki/Eclipse_(software)

9. Android software development Wikipedia, the free encyclopedia. (n.d.). Retrieved July 10, 2013, from https://en.wikipedia.org/wiki/Android_software_development

10. Retrieved April 4, 2014, from https://console.developers.google.com/project/118826163607/apiui/api/translate/method/language.detections.list

11. *java - Using google translate in android application - Stack Overflow*. (n.d.). Retrieved March 24, 2014, from http://stackoverflow.com/questions/16809232/using-google-translate-in-android-application