



UNIVERSITY OF HUDDERSFIELD

MENG GROUP PROJECT

Cryptic Crossword Solver

Authors:

Mohammad RAHMAN
Leanne BUTCHER
Stuart LEADER
Luke HACKETT

Supervisor:

Dr. Gary ALLEN

Examiner:

Dr. Sotirios BATSAKIS

Friday, 18th April 2014

Abstract

Cryptic crosswords are a unique style of crosswords in which the answer to each given clue is a word puzzle. An answer can only be obtained if the cryptic clue is read in the correct way. Often when the clue is surface read, the clue makes no sense at all. The challenge is to find a way in which the reading of the clue leads to a solution.

Cryptic crosswords are a popular type of puzzle found in many parts of the world. Most of the national UK newspapers will print cryptic crosswords of varying difficulty on a daily basis. Many users can often become frustrated when a clue appears to be unsolvable.

Clues that appear to be unsolvable are often all part of the fun of cryptic crosswords. However each clue can take some time to solve due to the fact that there are many forms of wordplay.

Contents

1 Problem Analysis	7
1.1 Problem	7
1.2 Product	8
1.3 Client & Stakeholders	8
1.4 Users	9
1.5 Research Areas	13
2 Research	16
2.1 Crosswords	17
2.2 Natural Language Processing	28
2.3 Mobile Platforms and the Market	34
2.4 Web Services	43
3 Development Methodologies	50
3.1 Waterfall	51
3.2 Spiral	52
3.3 Agile	53
3.4 Rapid Application Development	54
3.5 Summary	55
4 Software Specification	56
4.1 Aims & Objectives	57
4.2 Constraints	58
4.3 Facts and Assumptions	59
4.4 Scope	60
4.5 Functional Requirements	63
4.6 Non-Functional Requirements	64
4.7 Risk Assessment	67
4.8 Feasibility Study	73
5 Design	74
5.1 Primary Path	75
5.2 Activity Diagrams	77

5.3	Use Case	84
5.4	Sequence Diagrams	86
5.5	Class Diagrams	89
5.6	Database Design	99
5.7	User Interface	100
6	Development	105
6.1	Web Service & Servlets	106
6.2	User Interface	109
7	Testing	112
	Glossary of Terms	112

List of Figures

1.1	Results for survey questions 1 - 6	11
1.2	Results for survey questions 7 - 10	12
2.1	Creating a parse tree for a given sentence	28
2.2	A single sentence is already in the desired state	29
2.3	Assigning <i>Parts of Speech</i> tokens to an input sentence	31
2.4	A demonstration of the OpenNLP DocumentCategorizer	32
2.5	Mobile phone market share September 2013	40
2.6	Screenshots of Puzzler Super Cryptic Crosswords on iPhone	41
2.7	Screenshots of Cryptic Crossword on iPhone	42
4.1	Hierarchical Structure of the team	60
4.2	A visualisation of the identified risks of the project	69
5.1	The complete Activity diagram	78
5.2	Activity diagram showing how a user would input a clue	80
5.3	Activity diagram showing how a solver is to be chosen	81
5.4	Activity diagram showing how results should be handled	83
5.5	Use case illustrating the actor's range of possibilities	84
5.6	System Use Case	85
5.7	Sequence diagram illustrating a user submitting a clue	87
5.8	Sequence diagram illustrating the system solving a clue	88
5.9	Package overview of the Cryptic Crossword Solver system	89
5.10	Config package class diagram	90
5.11	Core package class diagram	92
5.12	NLP package class diagram	93
5.13	Resource package class diagram	95
5.14	Servlet package class diagram	96
5.15	Solver package class diagram	97
5.16	Util package class diagram	98
5.17	Testing database entity-relationship diagram	99
5.18	The input form to be completed by the end user	100
5.19	The input form to be completed by the end user	101
5.20	The input form indicating a validation error	102

5.21 Results list displaying the top answer in blue	103
5.22 Results list displaying alternative solutions	104
6.1 Message informing the user that a clue is currently being solved	110
6.2 Live validation feedback ensures users are entering correct data	110

List of Tables

2.1	A comparison of available NLP libraries	31
2.2	A comparison of mobile platforms	35
2.3	A further comparison of mobile platforms	35
2.4	More comparisons of mobile platforms	36
2.5	A comparison of cross platform development tools	36
2.6	HTTP verbs mapped to the associated CRUD operation.	46
4.1	MoSCoW analysis of the project's functional requirements	63
4.2	MoSCoW analysis of the project's usability requirements	66
4.3	Probability-impact table for the project's risks	68
4.4	Feasibility Study for Clue Types	73

Chapter 1

Problem Analysis

An initial problem analysis has been conducted to ensure that the overall project remains focused upon the original problem.

This chapter will discuss key topics and will recommend that these topics are researched further to help support the project. The problem analysis will also define the problem in more detail, to help with understanding the project and its purpose.

1.1 Problem

A typical crossword involves a grid with black squares, not to be filled in by the solver, and white squares used by the solver to input their answers. The input comes from the solver working out the answer to given clues. These clues have an orientation, a length and a number associated with its related position within the grid. The fundamental difference between a typical crossword and a cryptic crossword are the clues themselves.

Cryptic crosswords are a popular type of puzzles found in many parts of the world. Most common-wealth national newspapers will print cryptic crosswords of varying difficulty on a daily basis.

Cryptic crosswords are a unique style of crosswords, in which the answer to each given clue is a word puzzle. An answer can only be obtained if the cryptic clue is read in the correct way. Often when the clue is surface read, the clue makes no sense at all. The challenge is to find a way in which the reading of the clue leads to a solution. To aid with solving cryptic crosswords, the clues are written to be within specific categories, such as reversals and anagrams, which have individual characteristics.

Many users can often become frustrated when a clue appears to be unsolvable. It is the vast range of possible clues that often makes solving not only challenging but interesting as well.

Fundamentally, the overall aim of this project is to develop a piece of software that is able to solve any given type of cryptic crossword clue.

1.2 Product

Within this group project, three components will be delivered. The first deliverable is the final, working piece of software. Whilst the second and third deliverables are written reports. The second deliverable is a group written report comprising of the all research and implementation details of the software product. The final deliverable will be each member's individual analysis and evolution of the project as a whole.

Based upon the given background and problem information it could be possible to develop a product that is able to solve the given problem.

The final product would be a piece of software that is able to understand a given clue and try to deduce what the answer to the clue is. This would require the software to have some form of natural language processing component as well as one or more cryptic crossword algorithms. Once a clue has been correctly "guessed" it can simply be returned to the user. It is the "guessing" of the answer that this project will primarily focus upon.

In order to gain maximum user coverage, the software must have an easy to use interface. The main reason for this is that the computer literacy of the intended users is not known - although basic computer literacy is assumed.

1.3 Client & Stakeholders

Dr Hugh Osborne, a lecturer from the University of Huddersfield will be the client for the group project. Dr Hugh Osborne has a keen interest in cryptic crosswords and the problem in the area which the group intends to help to eliminate. The role of the client for the group project will be to input ideas and potential requirements which Dr Hugh Osborne, as an experienced solver of cryptic crosswords, would consider to be necessary. As a client for the project, Dr Hugh Osborne will also be present for academic demonstrations.

Dr Gary Allen, Sotirios Batsakis and Colin Venters, all lecturers at the University of Huddersfield, will act as stakeholders for the group project. Dr Gary Allen will be the most involved external individual as the project supervisor. The role of project supervisor requires frequent meetings with the team to monitor the development of the project, provide guidance as well as opinion on certain aspects of the life cycle.

Sotirios Batsakis and Colin Venters, will have a less active role within the project during the project life cycle than the role of the client or the project supervisor. These stakeholders will play active roles at particular milestones of the group project such as providing guidance

for the proposal of the project and at the project demonstration approximately half way through the life cycle.

1.4 Users

Kathryn Friedlander and Philip Fine (Friedlander and Fine, 2009) carried out an investigation into whether the amount of cryptic crosswords completed by a solver determined how successful they were at solving them. To complete this study they gathered data from 241 people and have deduced the following facts about the user base (Friedlander and Fine, 2009):

- “209 M, 32 F”
- “mean age=53 years, range=23 – 83”
- “mean time spent=8 hours per week, range=1-30”

To support decisions made within the project life cycle an additional quantitative research method has been utilized to gain a larger understanding of the types of users the deliverable will attract. The research method used by the team is in the form of a survey.

The survey results gathered were seen as additional justifications for the purpose of the project. Moreover, data collected from the survey was expected to indicate the locations in which users complete cryptic crosswords to understand the potential need for the deliverable to be of a transportable nature.

The following questions were asked:

1. Do you play Cryptic Crosswords?
2. How often do you play?
3. Where do you play?
4. Do you often finish them?
5. If no to the previous question, what reason don't you finish them?
6. What is your age group?
7. When do you play Cryptic Crosswords?
8. What gender are you?
9. What is the Highest qualification you have?
10. What platform is your mobile phone on?

The survey was conducted between (date) and (date). It was distributed across the the University of Huddersfield portal message board, Facebook and Twitter. The results of this survey are shown in figure 1.1 and figure 1.2.

1.4.1 Results

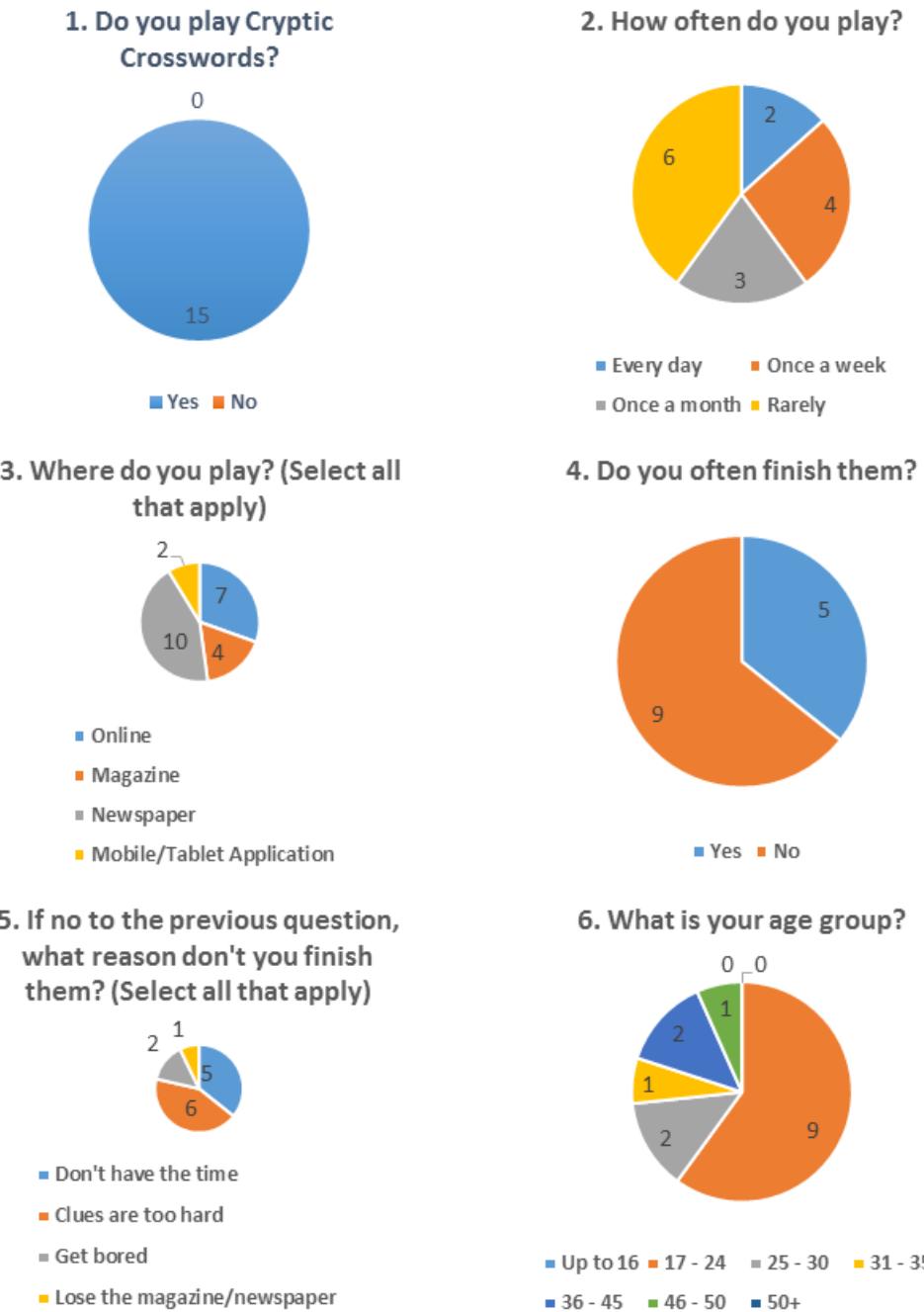


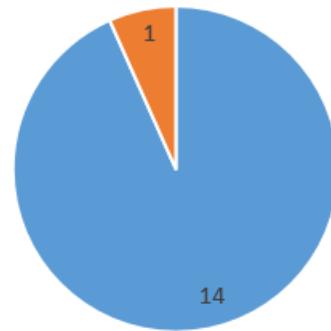
Figure 1.1: Results for survey questions 1 - 6

7. When do you play Cryptic Crosswords? (Select all that apply)



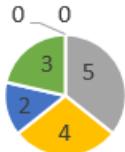
- When travelling
- When I am at work
- During my lunch breaks
- When I am bored
- When I am studying

8. What gender are you?



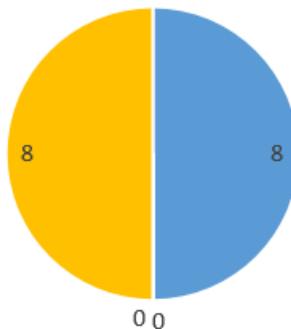
- Male
- Female

9. What is the highest qualification you have?



- None
- GCSE (O Levels) or Equivalent
- A Levels
- Undergraduate Degree
- Postgraduate Degree
- PhD

10. What platform is your mobile phone on? (Select all that apply)



- Android
- Blackberry
- Windows
- Apple

Figure 1.2: Results for survey questions 7 - 10

From these results it can be deduced that

1.5 Research Areas

Before undertaking the project an initial review was conducted. The review's objective was to determine the feasibility of the project as a whole. The review also covered whether or not the project has been completed before.

From the outlined background and problem information it is clear that cryptic crosswords are a popular form of entertainment. It is also clear that some clues are particularly difficult to solve, and users may often ask other people for help in solving a given clue.

1.5.1 Cryptic Crosswords

A review of the national UK newspapers was conducted to determine whether or not there is a pattern in cryptic crosswords. Of all the newspaper's websites that were reviewed (The Guardian, The Times, The Independent and The Mirror) it was clear that all cryptic crosswords are of the same style.

Each clue is categorised as being either 'across' or 'down' with its corresponding grid number, as well as containing the number of letters the answer should be. An example is show below:

12. The seamstress's sensation? (4, 3, 7) =>PINS AND NEEDLES

The Guardian's website utilises web standard technologies such as HTML and CSS, and also provides an option to solve a clue. The Mirror's website follows a similar approach to the Guardian's website; however solutions can only be obtained by dialing a premium telephone number.

The Times and the Independent both utilise a different approach and that is to serve a Java applet. Both Java applets allow the user to solve a clue should they get stuck. The Times provides puzzles as part of their paid subscription service.

All of the above newspapers publish cryptic crosswords upon a daily basis, with the solutions to the crosswords appearing in the next day's newspaper.

Following from the crossword review, a second review into cryptic crossword solvers was undertaken. The objective of this review was to determine whether or not computerised cryptic crossword solvers exist. The three cryptic crossword solvers that were identified were One Across, Crossword Tools and Cryptic Solver.

Each of the solvers manages to solve some clues with the same answers, with other clues providing a range of possible answers.

Crossword Tools (Crossword Tools, 2013) is a paid subscription based service, which allows users to enter a clue and a pattern. A pattern can contain part of the answer or the number

of letters the answer has. If multiple answers are available, they are displayed. An example is shown below:

Kind of dog (10) =>the answer is 10 letters long.

Kind of dog (?????????r) =>the answer is 10 letters long, final letter is ‘r’.

Cryptic Solver (Cryptic Solver, 2013) is a free service that offers the same functionality as Crossword Tools. Although Cryptic Solver does provide the correct answer, it does not necessarily provide the correct answer at the top of the list.

Finally One Across (One Across, 2013) provides all the same functionality as the previous two solvers, along with a score. The score is linked to the number of people who have used the given answer (effectively it’s a ratings system). One Across uniquely highlights how it has managed to deduce the answer, showing the break downs of each sentence. As with Cryptic Solver, One Across is a free service that doesn’t require a subscription.

1.5.2 Natural Language Processing

In order to correctly solve a clue, some form of natural language processing may be required. It is the natural language processing that could try to deduce the meaning of a clue. It is the meaning that can then be aligned with possible answers.

An example of natural language processing can be found within the One Across application. Given a clue (and a pattern) it will try to provide an accurate solution:

Spin broken shingle (7) =>ENGLISH

In order for the answer to be obtained, One Across will follow a natural language processing path and will provide it’s trace path. The trace path shows how the clue has been broken down to get to the answer. The trace path for the above clue can be found below:

‘spin’ is the definition.

‘broken’ means to anagram ‘shingle’ to get ENGLISH.

ENGLISH matches ‘spin’ with confidence score 100%.

1.5.3 Application Platform

The existing products that have been discussed within this problem analysis have all been accessible via a browser. Although this is an acceptable platform, there could be a better platform that allows users to utilise the technology easier.

As previously mentioned, most crosswords are designed for users who have a few minutes to spare on the move. With the recent trends in owning a smartphone or tablet, there may be a gap in the market for a high quality mobile cryptic crossword solver.

An in-depth review will need to be conducted in order to deduce the viability of this proposal.

Chapter 2

Research

In the previous chapter a detailed approach to what the project will cover was discussed. In order to correctly implement a working piece of a software a number of areas will need to be researched thoroughly.

Based upon the problem analysis, the the main topics that have been identified are:

- Cryptic Crosswords
- Natural Language Processing
- Application Platform
- Web Services

2.1 Crosswords

Arthur Wynne produced the first crossword puzzle which was printed on December 21st 1913. A crossword is a puzzle which involves the solver resolving the answer to a clue and placing it in the correct space within the grid.

A grid is made up of black and white squares, the black squares are blanks and the white squares are where the solver must place the answers. A crossword grid comes with a set of clues. The clues are usually arranged based upon their position within the grid. Clues that appear downwards in the grid are kept separate to the clues that appear across the grid.

The white squares which are used for the first letter of an answer to a clue usually have a number in the top left hand corner to indicate the clue which links to this area of the grid.

There are different types of crosswords such as quick, cryptic and double-clue. A quick crossword has clues which simply define the answer. A cryptic crossword is more complex as it has word play as well as simple definitions and many different types of clues. A double-clue crossword combines the two and allows for a simpler option for the solver when the cryptic clues become too difficult.

2.1.1 Cryptic Crosswords

Most cryptic clues consist of two different parts, the word play and the definition itself. The definition is like a clue found within a quick crossword and the word play is an indication to the answer.

Clue types such as double definition and purely cryptic break the usual format of cryptic clues. Double definitions miss out the word play whereas the purely cryptic clues miss out the usual simpler definition and become a fully cryptic definition. Other types of clue add to the usual format when smaller clues are embedded within the larger clue to assist with the word play.

Punctuation within clues should always be disregarded unless it is a question mark or an exclamation mark. Punctuation such as commas and hyphens are used to distract the solver from the answer usually by attempting to dictate how the clue is read.

A question mark tells the solver that the clue requires creative thinking to work out the answer which could be witty in nature. An exclamation mark can mean that the word play and the definition may intersect which is otherwise known as a clue of the type “& lit”. Articles within clues can also be very important and should not be disregarded when reading the clue.

Cryptic clues which have a particular tense will always be for a clue with the same tense. Similarly a plural clue determines that the answer will also be plural.

2.1.2 Crossword Clue Types

An important skill needed to solve a cryptic crossword is to be able to spot the type of clue given. Below is a list of the most common types of cryptic clue and expected rules the clue should follow so they are identifiable.

Purely Cryptic

Although clues within a cryptic crossword usually include both a definition and word play, this type of clue is an exception because the whole clue is a definition written in an unusual way.

Word play within other clue types assist the user in being able to determine their answer is correct as well as solve them, because there is no word play more than one answer could be found which are incorrect. Possible indicators that denote a purely cryptic clue are:

- Question mark
- Exclamation mark

Example: *Frames for summer's activities? (5)* - (Upadhyay, 2008c)

Answer: ABACI

- “Summer” as in a person who does mathematical sums
- An abacus is a frame which holds moving beads
- As the clue is plural so must the answer be, hence abaci

Hidden

The answer for a hidden clue is concealed within the clue itself and can be spread over more than one word as well as possibly being hidden in reverse. The clue will have a definition, an indicator that the answer is hidden within the clue and a word or set of words which have the answer in them. Possible indicators that denote a hidden clue are:

- Word/s e.g. contains, in, within, held by, from
- Large words with a hidden word indicator before it may have the answer inside them
- A clue which seems inelegantly written or a clue which contains proper nouns

Example 1: *Metal concealed by environmentalist* (4) - (Upadhyay, 2008e)

Answer 1: IRON

- “Concealed by” is a phrase indicator for hidden clues
- “Environmentalist” is a large word with an indicator in front
- “Metal” is the definition so the answer is a type of metal which can be found within the word “environmentalist”, hence iron

Example 2: *Mountain range in central Taiwan* (5) - (Upadhyay, 2008e)

Answer 2: ALTAI

- “in” is a word indicator for hidden clues
- “central Taiwan” contains a proper noun which indicates the answer is hidden here
- “Mountain range” is the definition
- Without knowledge of mountain ranges the answer could be narrowed down to the following words (assuming “central” would not be within the clue without a purpose):
 - TRALT
 - RALTA
 - ALTAI
 - LTAIW
- If some of the crossword is completed within the area this clue is placed, the correct answer could be found through trial and error.

Charades

A charade clue forms its answer with the use of smaller answers to smaller clues within the main clue. Abbreviations and first/last letters of words are common within charade clues to make up the complete answer.

Two or three parts are usually within the charade clue to solve the correct answer, they may not be in the right order however, and word indicators will be used to warn the solver.

Other types of clues can also be used within charade clues for the different parts such as reversals and homophones, if this is the case there will be indicators for the specific type. Possible indicators that denote a charade clue are:

- Words e.g. with, follows, behind, after to indicate joining of answers to parts of the clue

Example: *Prior belted one that is ultimately right (7)* - (Upadhyay, 2008b)

Answer: EARLIER

- “Prior” is the definition
- “belted one” gives earl
- “that is” gives the abbreviation for i.e. or ie
- “right” gives the abbreviation for r
- All the segments put together give the word “earlier” which can also mean “prior”

Anagrams

Anagram clue types have a definition, a word or phrase to indicate the clue is of this type and an element called “fodder”.

An anagram is a word whose letters can be rearranged to form another word; within an anagram clue the letters to rearrange are known as “fodder” and are placed next to the indicator. Possible indicators that denote an anagram clue are:

- Words which could mean change or shifting
- A clue which seems inelegantly written or a clue which contains proper nouns

Example: *Toy breeds trained to find out a place for pearls (6,3)* - (Upadhyay, 2008f)

Answer: OYSTER BED

- “trained” indicates the clue is an anagram
- “Toy breeds” is an abnormal phrase and is the “fodder” of the clue
- to find out a place for “pearls” is left to become the definition
- “Toy breeds” is then moved around to give oyster bed

Homophones

A homophone is a word which sounds like another word but has a separate meaning. This type of clue has a definition, a word or phrase which means the same as the homophone to find and an indicator. Possible indicators that denote a homophone clue are:

- Words which indicate hearing or sound e.g. said, heard
- Normally the indicator for a homophone is next to the word or phrase which is to be used to find a homophone

Example: *Refer to a location, reportedly (4)* - (Upadhyay, 2008j)

Answer: CITE

- “reportedly” is the homophone indicator
- “a location” is the phrase which needs to be used to find a homophone
- “Refer to” is the definition
- A location can be otherwise known as a “site”, hence cite

Acrostics

An acrostic clue commonly involves picking the first letter from a group of words and putting them together to form the answer. It is possible that the clue will require the last or middle letters from words to solve them or that the letters should be put together in reverse order.

This clue type has a definition and an indicator as well as “fodder” which in this case means the group of words the necessary letters will come from. Possible indicators that denote an acrostic cryptic clue are:

- Words which could mean start or beginning
- If the clue is unusually long and so is the number indicator to determine how long the answer should be

Example: *Some URLs recommended for beginners to explore online (4)* - (Upadhyay, 2008a)

Answer: SURF

- “beginners” is the indicator
- As the clue states the answer should be of length four it is assumed “to explore online” is the definition as there are only three possible letters for an acrostic in the phrase
- “Some URLs recommended for” is the “fodder” for the clue. The indicator implies that the first letters should be taken from the first letters of the words within the “fodder”, hence surf

Palindromes

A palindrome is a word which reads and looks the same when it is reversed. This type of clue has an indicator and a definition. Possible indicators that denote a palindrome clue are:

- Phrases which may mean either way or going around in circles

Example: *Unacceptable, going up or down (3,2)* - (Connor, 2012a)

Answer: NOT ON

- “going up or down” is an indicator as it could mean “in either direction” like the format of a palindrome
- “Unacceptable” is then left as the definition which gives the answer, not on

Reversals

A reversal requires the solver to reverse a number of letters to give a new word. The clue consists of a definition, an indicator that the clue is a reversal clue and some “fodder” which is a phrase or word which could contain the letters to be reversed or a smaller clue which leads to the letters which need to be reversed. Possible indicators that denote a reversal clue are:

- Words/phrases used for directional purposes e.g. left, up
- The word indicators may also be relative to the direction the clue should be placed within the crossword (down or across)

Example: *Stop the flow in crazy get-up (3)* - (Upadhyay, 2008h)

Answer: DAM

- “Stop the flow” is the definition
- “get-up” is the indicator that the clue is a reversal
- Another word for “crazy” is mad which reversed gives the answer dam

“& lit”

“& lit” clues, which means “and literally so”, is a type of clue where the definition and the word play are the same and are not split out into separate phrases or words as with other clues. The definition is the whole clue and the word play can be one or more of any of the normal clue types such as anagrams and charades. Possible indicators that denote a “& lit” clue are:

- Exclamation mark

Example: *Cop in male form* (9) - (Upadhyay, 2008g)

Answer: POLICEMAN

- The whole clue is the definition
- “form” is an indicator for an anagram clue
- “Cop in male” can be rearranged to policeman which is also the answer to a “Cop in male form”

Double Definition

A double definition clue has no word play and is purely a clue with two (or possibly more) definitions which lead to the same answer. Possible indicators that denote a double definition clue are:

- Possibly shorter than most clues (2 or 3 words)
- Although it is advisable to ignore punctuation when solving cryptic clues, a double definition may have a piece of punctuation separating the definitions.

Example: *Robust author* (5) - (Upadhyay, 2008d)

Answer: HARDY

- Both words are separate definitions which could both mean hardy

Containers

A container clue includes three definitions and an indicator. One of the definitions is for the final solution whereas the other definitions describe two separate answers where one is contained within the other.

They are similar to charade clues in the way that other types of clues can be used within container clues such as anagrams and charades themselves. Possible indicators that denote a container clue are:

- Word indicators which could be used to indicate the inner word (e.g. inside, held) or the outer word (e.g. outside, external)

Example: *Building for the workers in principle* (8) - (Upadhyay, 2009a)

Answer: TENEMENT

- “principle” is the definition for the outer word which gives tenet

- “in” is the indicator that the answer for the inner word will be placed within the outer word (tenet)
- “the workers” is the definition for the inner word which is men
- ‘Building’ is the overall definition which could mean tenement which is also given when men is put within tenet

Deletions

Deletion clues require the solver to retrieve the answer by looking at the definition and the word play and removing the correct letters from the correct word. The clue will not usually have the word which needs letters removing from it directly within the clue. Possible indicators that denote a deletion clue are:

- Words to indicate letters should be removed from a certain place within the word such as its first or last letter
- Words to indicate certain letters should be removed from the word found. These could be words that can be abbreviated .

Example: *Little shark edges away from diver’s equipment (3) - (Upadhyay, 2009b)*

Answer: CUB

- “edges away from” indicates that the first and the last letter should be removed from the answer which comes from “divers equipment”
- “divers equipment” is otherwise known as scuba
- Removing the “s” and the “a” from scuba leaves the word cub which is also the answer given from the definition “Little shark”

Spoonerisms

A spoonerism clue has a definition and word play which is usually a phrase which describes another. This phrase, which is usually two words long, is taken and the first letter of each is swapped around to gain a new phrase which is in turn the answer to the clue. Possible indicators that denote a spoonerism clue are:

- The word “Spooner” or “Spooner’s” is placed within the clue

Example: *Spooner’s cheerful enthusiast? He’ll get you across (8) - (Connor, 2012b)*

Answer: FERRYMAN

- “Spooner” indicates that this clue is a spoonerism
- “Hell get you across” is used as the definition
- “cheerful enthusiast” could also be known as a merry fan
- Swapping the first letters from the words merry and fan give the answer ferryman

Pattern

This type of clue has a definition, an indicator that the clue is of the type pattern and a phrase or word which has the answer within them arranged as a pattern. These patterns can be odd or even letters joined together to make a word or possibly letters picked from regular intervals.

Pattern word play can also be joined together with other types of clues such as charades to form a more complex answer. Possible indicators that denote a pattern clue are:

- Words which could mean even, odd or routine

Example: *Beasts, free, ginned, we hear regular losses there! (8)* - (Upadhyay, 2009d)

Answer: REINDEER

- “regular” indicates it is a pattern clue as well as “losses” to indicate the dropping of letters.
- “Beasts” is the definition.
- “free, ginned, we hear” holds the answer reindeer by picking out the first “r” within free and each letter alternately from then on.

Substitutions

A substitution clue involves removing letters from a word and replacing it with another to retrieve the answer. There are two definitions within the clue, one definition to retrieve the word to substitute letters from and another to define the final answer.

The letter or letters to substitute are usually an abbreviation which can be found within the clue itself. Possible indicators that denote a substitution clue are:

- Words which mean substitution e.g. replace, switch, exchange

Example: *Unexciting story gets mark for length (4)* - (Upadhyay, 2008i)

Answer: TAME

- “mark for length” is an indicator for a substitution clue
- “mark” can be abbreviated to “m” and “length” can be abbreviated to “l”, therefore replace “l” with “m”
- “story” is the definition for the word which needs the substitution and could be defined as a tale
- Replacing the “l” in tale with “m” gives the word tame which can also mean unexciting

Shifting

A shifting clue has an indicator, a definition of the final answer and another definition for the word which needs to be used to shift a letter to a different position within the word to find the final answer.

The shifting of a letter could be moving the first letter to the last position in the word or in a more complex clue letters could be shifted within the middle of the word. Possible indicators that denote a shifting clue are:

- Words e.g. shift, change, move
- Phrases e.g. head to foot

Example: *Character needs help, head to foot (4)* - (Upadhyay, 2009e)

Answer: BETA

- “head to foot” indicates moving a letter from the front of a word to the end
- “help” is the definition for the word which requires letter shifting and can also be defined as abet
- Moving the first letter of “abet” to the end gives the word beta which is a “Character”

Exchange

An exchange clue is similar to a shifting clue however instead of only one letter shifting positions within a word, two letters within a word exchange places to form a new word.

Typically the letters to exchange will be the first and last letters of a word or two letters next to each other, however it is possible more than one letter on each side will need to be swapped.

For example, the word “rage” can be split into two sections “ra” and “ge” which can then be exchanged to make the word “gear”. Possible indicators that denote an exchange clue are:

- Words e.g. swap, exchange, change

Example: *Doomed king switching sides? True (4)* - (Upadhyay, 2009c)

Answer: REAL

- “switching sides” indicates that this clue is an exchange clue
- A “Doomed king” can also be known as a “Lear”
- “True” is the definition which can also be defined as “real” which can be gained by exchanging the first letter of “Lear” (“L”) with the last letter (“r”)

2.2 Natural Language Processing

An inherent disparity between the languages used and understood by human beings and those which are interpreted by machines introduces the requirement for a system which is able to provide a form of compatibility between the two. Natural language can be described as a system of communication that has not consciously been invented (Collins, 2004) and typical examples of this are the languages used by human beings to communicate with each other. English and Korean are two such examples of natural languages and these evolve over time as words are added and removed and the rules that define the language are refined or adapted, meaning there is no immutable grammar for the language.

Conversely, formal (or artificial) languages are those which have been fabricated for a specific purpose, and can be described using precise and unambiguous mathematical rules (Jiang et al., 2010). An example of such a mathematical rule is the BackusNaur Form (BNF) grammar. Programming languages such as Java or Python are examples of artificial languages, as these have been created with a specific purpose in mind. The vocabularies of these languages are well defined and free from ambiguities, making interpretation by the languages corresponding compiler a black or white task where the given language to compile is either valid or invalid.

Natural language processing (NLP) allows for the manipulation of natural languages by a computational device (Bird, 2009) and there are a number of identified steps of the NLP process which have been developed, from scanning and parsing, to meaning extraction and integration. In the context of solving cryptic crosswords, a given clue would need to be broken down into meaningful components which would allow, as an example, the word play and the definition components to be identified and passed forward to the appropriate algorithms to identify the correct answer.

```
Reversible fasteners pinning baker, European
<TOP <NP <NP <NNP Reversible> <NN S fasteners>> <VP <UBG pinning> <NP <NN baker,>
>> <. European>>>
```

Figure 2.1: Creating a parse tree for a given sentence

(Halpern, 2013)

2.2.1 The Steps of Natural Language Processing

Scanning

A number of steps are typically involved in the processing of a natural language. Assuming the input format of the natural language is in a textual form as sentences (rather than as speech or handwriting); these chains of text may be broken down into their separate components - the words which make up the sentences. This process is known as scanning, and may itself comprise of several sub-steps. The first step could likely involve breaking down the given textual input into separate sentences (Apache, 2013).

```
C:\Users\u0955187\Downloads\apache-opennlp-1.5.3-bin\apache-opennlp-1.5.3\bin>opennlp SentenceDetector en-sent.bin < input.txt
Loading Sentence Detector model ... done <0.068s>
Reversible fasteners pinning baker, European
```

Figure 2.2: A single sentence is already in the desired state

Once the separation of any sentences present in the input has occurred, these can be further divided into the separate words to achieve an array-like data structure, such as [*Reversible, fasteners, pinning, baker, European*], a process which is referred to as tokenising. It is important to consider the inclusion and placement of any punctuation provided in the input text, as this may provide some indication as to the separation of the cryptic crossword components and increase the chances of successfully calculating the answer. In the example provided, European corresponds to the definition component of the clue, where Reversible fasteners pinning baker relates to the word play component. Both components have been clearly separated by a comma, though not all cryptic clues are presented in this fashion.

Parsing

Parsing is a process used to ensure that a given sentence conforms to a grammar and thus showing that it is syntactically correct (McCluskey, 1999), and this process can be demonstrated by creating a parse tree to represent the structure of a given sentence. It may become apparent when attempting to solve the clues of a cryptic crossword that certain elements of the provided clue sentences bear some relation to the correct answer. To give an example, it may prove more likely that the adjectives or verbs present in the input are synonyms of the correct answer, where there may be little probability that the answer stems from a synonym of any present determiners (the, this, a, some) or prepositions (on, beneath, over, under). Likewise, adjectives or adverbs that hint at a particular action such as *reversible* or *backwards* could reveal themselves as more likely to be indicators that describe in what way the clue should be manipulated.

Referencing

The ambiguous nature of cryptic clues may not allow for a straightforward parse of an already ambiguous (natural) language, possibly contributing to the complexity of the process of solving. It may be the case that multiple parses of the same clue have to be taken forward to be processed before the solution can be found. Once a parse tree has been selected for the given input, a number of steps can be performed, which attempt to reduce the level of ambiguity. Lee McCluskey outlines four methods which can be applied to achieve this:

- **Contextual Disambiguation:** Such as referring to segments of a previous piece of text. Fred Bloggs may have been declared in a previous sentence, and from that point referred to as he.
- **Physical Constraints:** This refers to the literal meaning of a sentence or phrase, and uses known facts to assist in the reduction of ambiguity. The phrase I answered the door in my pyjamas could be taken to mean that the subject was either wearing their pyjamas when they answered the door, or that the door was wearing the subjects pyjamas when they answered it. As it is unlikely that the door was physically wearing pyjamas, the parse suggesting this could be disregarded.
- **Default Roles in Known Verb Structures:** Some verbs may only be applied to a subset of subjects in order to make reasonable sense. In the example sentence fruit flies like a banana, one parse may take fruit as the noun, and flies as the verb. In this scenario, fruit is physically incapable of flying, and so an alternative parse which takes fruit flies as the noun-phrase is more likely to be the correct parse of the sentence.
- **General Defaults:** To say that one is going to visit an old friend may imply that the friend has been known for a long time, or that the friend is of old age. While the correct interpretation may be context specific, the likelihood is that one particular interpretation will hold precedence over another. In this case, the friend is likely to have been known for a long time.

(McCluskey, 1999)

A parsing algorithm is used to obtain a complete parse tree for a given sentence, providing insight into the complete structure of the sentences. It is possible that this may not be required for the task at hand; a full parse may provide more information than what is required to solve a cryptic clue. An alternative mechanism exists, chunking, which may also be referred to as *light parsing*. This technique divides a sentence into a series of non-overlapping chunks of text, which provide an overview of the input sentence (Litman, 2003). One of the major benefits of this technique is the removal of the need to resolve ambiguity, as chunks of the sentence are non-recursive, unlike parsing which produces components of a sentence which may be nested in each other.

`Reversible_NNP fasteners_NNS pinning_VBG baker,_NN European._.`

Figure 2.3: Assigning *Parts of Speech* tokens to an input sentence

Meaning Extraction

Meaning extraction refers to the process of taking the parse and storing it in a way which accurately represents the information it models (McCluskey, 1999). For example, the following cryptic crossword clues all map to the same answer, and whilst they are written differently, they each imply the same meaning of *DRILL*:

- *Doctor gives patients exercise*
- *Doctor needing a doctor for practice*
- *GP not fit for practice*

(Gordius, 2003)

Effective meaning extraction will ensure that each distinct clue may help to ensure that variations of the same will be mapped to a single, internal representation and this will allow for the answer to be retrieved from a database without having to recalculate what has already been processed once.

2.2.2 Natural Language Processing Libraries

There exists a range of NLP libraries that offer many natural language processing functions for a variety of programming languages. Development on a number of these libraries has become stale, and others are available only for use after purchasing a commercial licence. Below are a selection of available libraries which have recently been updated and are available under some form of free software licence.

Library	Licence	Last Updated	Supported Languages
OpenNLP	Apache 2.0	April 2013	Java
NLTK	Apache 2.0	November 2012	Python
Stanford CoreNLP	GNU General Public Licence	June 2013	Java

Table 2.1: A comparison of available NLP libraries

OpenNLP

OpenNLP provides the ability to carry out a range of natural language processing functions, including parsing, chunking and name finding, where the latter aims to recognise proper nouns in specified text such as person names. The features of the library utilise a maximum entropy model, and this allows the performance of the software to be enhanced as it uses training data to learn (Apache, 2013). This allows for each component of the toolkit to refine as the subset of training data increases over time, with the aim of increasing the accuracy of the corresponding NLP component. The library is available for use in a Java environment or can be accessed directly through a command line interface, but the general usage of the library requires providing an applicable model and language data as input, for which the model will be utilised for the language processing.

Models A number of models are provided for use with the NLP library. These include, but are not limited to, models to assist in the process of parsing, chunking and detecting sentences. A number of libraries also exist which allow for the identification of person names, company names, times, dates or location in input text. These existing models may be further trained in a bid to increase their effectiveness.

DocumentCategoriser Another prominent feature of the library is the ability to classify input text into a range of predetermined categories using the aforementioned maximum entropy model. Once a model has been created, which contains a series of example inputs along with their corresponding categories, further inputs to the system will be paired with a *best outcome* classification.

```
<terminated> NLP [Java Application] C:\Program Files (x86)\Java\jre7\bin\javaw.exe (Nov 1, 2013 3:18:08 PM)
Category for: Posy says no, according to Spooner is: SPOONERISM (0.18832631809541703)
Category for: Preparation for documents once Spooner's transporting luggage is: SPOONERISM (0.005477603261500726)
Category for: Preserve holy one appointed by Spooner is: SPOONERISM (0.003470437596580894)
Category for: Theme's proper order reversed is: REVERSAL (0.016981012268904402)
Category for: This compiler going to the stake? Quite the reverse! is: REVERSAL (3.1981000778711125E-4)
Category for: This vehicle reverses into its club is: REVERSAL (0.01006333130189762)
```

Figure 2.4: A demonstration of the OpenNLP DocumentCategorizer

NLTK

The Natural Language Toolkit is an alternative library which exists with an open-source licence variant. NLTK exists for use with the Python programming language and possesses a similar infrastructure to OpenNLP, where trained models are used to allow each component of the software package to function correctly (NLTK Project, 2012). The package is

supplemented by a comprehensive e-book, detailing and providing examples of the usage of each component available.

2.3 Mobile Platforms and the Market

The demand for applications to be portable on mobile phones and tablets has become a fashion as well as accessibility on the go. In the current market Cryptic crosswords are available in a variety of newspapers, magazines and online websites. These are everyday media in which a person can access their favourite Cryptic Crosswords and participate in. With the expansion of mobile phone and tablet applications research has shown that a unique Cryptic Crossword solver has not been incorporated in any of the mobile operating systems. There are applications, which can be downloaded, that are pre-solved crosswords but there is no real solver, which solves Cryptic Crosswords in real time. In fact research has shown that there is only a very small handful of Cryptic Crossword applications across all mobile operating systems.

2.3.1 Mobile Operating Systems

The Oxford English dictionary states an “Operating System” as ‘the low-level software that supports a computers basic functions, such as scheduling tasks and controlling peripherals’ (Dictionaries, 2011).

A mobile operating system has the same definition but supports the basic functions for handheld devices such as mobile phones and tablets. In the current market the term mobile operating system is associated with smartphones rather than mobile phones due to the powerful processors, which are embedded in the mobile phone devices.

The official first mobile operating system for smartphones was by IBM introduced in 2000 as the ‘Simon’ but it was Ericsson who created the first all in one smartphone with the Symbian OS, which incorporated a keyboard hence, the term smartphone was introduced. This ran on various mobile phones by companies such as Ericsson, Samsung but predominantly on smartphones by Nokia. This created a path for a market, which was to be dominated by others in the coming future (Akman, 2013).

By 2007 from the back of a successful campaign selling music devices known as the iPod, Apple introduced the Apple iPhone which came with their first mobile operating system the IOS which was a full operating system used from the Mac OS X 10 (Apple, 2007).

By July 2008 Apple released IOS 2.0 and with this came the app store. This was a revolution to the mobile market allowing a platform for third party developers to sell and market their own applications for the mobile operating system. On the 7th January 2013 Apple announced that they have had more than 40 Billion downloads of apps through their app store (Apple, 2013).

In September 2008 Google released their own version of a mobile operating system Android which had its similarities of marketing with its very own app store known as the android

market which since has been rebranded to the Google play store.

In April 2009 Blackberry also launched its own application store called the Blackberry World, which works with their mobile operating system the Blackberry 10 (Cha, 2009).

Finally windows released their version of an app store for distribution of applications on October 26th 2012 (Businessline, 2012).

Although there are now several platforms that run on a mobile device, Android and iOS combine for 91.1% of the Worldwide Smartphone OS Market (IDC, 2013).

This shows that although blackberry has been in the market since 2009 there isn't much of a rise to interests in their apps and while Windows is fairly new it has a lot of ground to cover to catch up with their main competitors.

For the purpose of this project it is pretty clear to what is demanded from consumers in the real world with Apple and Google being the two main competitors for mobile applications. In order to decide on what platform the project will be suitable for to design, maintain and deploy a good working product below is the following page contains a table which covers some of the reasons which could be possible to allow the team members to come to a decision to what pathway the project will be going in.

Platform	Programming Language	Open Source	Open API/SDK	License
Android	Java	Yes	Yes	Apache
IOS	Objective C	N/A	Yes	Proprietary
Blackberry	Java, C++	N/A	Yes	Proprietary
Windows	C, C++, C#	N/A	Yes	Microsoft

Table 2.2: A comparison of mobile platforms

Platform	Latest Version	Debugging	Hardware / Software Requirements	Emulator
Android	4.3 Jelly Bean	Yes	Windows / Mac	Yes
IOS	IOS 7	Yes	Intel Based Mac	Yes
Blackberry	Blackberry 10	Yes	Windows / Mac	Yes
Windows	Windows Phone 8	Yes	Windows	Yes

Table 2.3: A further comparison of mobile platforms

Platform	Underlying OS	Development Environment	Submission To App Store	Development Cost
Android	Linux	Eclipse	Unlimited	\$25 One Off cost
IOS	Darwin	XCode	Unlimited	£60 Yearly / \$99
Blackberry	Blackberry OS	Eclipse	100 Per Year	\$100 One off cost
Windows	Windows	Visual Studio	Unlimited	\$19 Yearly / Free for Dreamspark Students

Table 2.4: More comparisons of mobile platforms

Although these are the four main types of platforms available in the market for mobile development there is a growing amount of organisations, which are developing tools to allow developers to create cross platform application with ease. Two of these are Appcelerator and Adobe AIR.

Platform	Programming Language	Open Source	Open API/SDK	Underlying OS and License	Development Environment
Appcelerator	JavaScript	Yes	Yes	Linux Apache 2.0	Eclipse Based IDE / Titanium Studio
Adobe AIR	ActionScript, HTML, CSS, JavaScript	No	Yes	Darwin Proprietary	Adobe AIR

Table 2.5: A comparison of cross platform development tools

2.3.2 Appcelerator

Appcelerator is a platform created by Appcelerator Inc to allow developers to create cross platform native applications for Android and IOS. They later introduced compatibility for Blackberry 10 and are in the process of developing for Windows Phone. The main development environment used to create applications is an Eclipse based IDE known as Titanium Studio. Developers can create great looking native apps using JavaScript. The use and license of Appcelerator falls under Apache 2.0.

2.3.3 Adobe AIR

Adobe Integrated Runtime (Adobe AIR) uses adobe tools such as flash to allow developers to create platform independent web apps. Unlike Appcelerator this means that applications created can only be web based and not native. For this reason a lot of developers avoid using Adobe AIR. It supports all the major vendors for mobile applications but apps created in Adobe AIR are a lot slower.

2.3.4 Native Apps

A native app is a platform independent application designed to work with a particular mobile OS. The app is installed on the device whether its a smartphone or a tablet.

Advantages

- Faster at accessing device features such as camera and accelerometer
- Easy to find in app stores such as Google Play and Apple app store
- Secure as they go through an approval phase from vendors

Disadvantages

- Expensive to develop
- Expensive to maintain
- Approval process can take from days to months
- Support of the app can be hard to maintain due to different people using different versions of operating system installed.

2.3.5 Web Apps

A web application is really a website designed to look and feel like an app. This is wrapped by the web browser, which means an Internet connection is required to use the app. Google Chrome and Safari are examples of Web apps.

Advantages

- Easy to maintain
- Can be compatible with any device
- Approval is not required
- Easy to maintain and update without affecting the user to update the software

Disadvantages

- Limited to what can be accessed on the devices
- Hard to provide support over various Web browsers
- Not easy to find and promote, users will have to browse websites
- Can be insecure due to the application being web based

There are various platforms available for mobile devices. The most popular platforms are the as previously mentioned. Although they are in contest the use and popularity of platforms vary from region to region. An article published by Sabri (2013) states that windows phone is the most popular platform in Latin America. Although the platform was produced later than the other platforms this has become popular due to the lower prices of devices. Devices such as the IPhone and the IPad are a lot dearer and can cost a user a couple of hundred pounds. In September Forbes reported that the most popular platform on mobile devices is Android.

Germany	3 m/e Aug 2012	3 m/e Aug 2013	% pt. Change	USA	3 m/e Aug 2012	3 m/e Aug 2013	% pt. Change
Android	78.7	78.7	0.0	Android	60.7	60.1	-0.6
BlackBerry	0.6	0.5	-0.1	BlackBerry	2.1	1.8	-0.3
iOS	11.1	9.5	-1.6	iOS	33.9	39.3	5.4
Windows	3.8	8.8	5.0	Windows	2.6	3.0	0.4
Other	5.7	2.6	-3.1	Other	0.8	0.8	0.0
GB	3 m/e Aug 2012	3 m/e Aug 2013	% pt. Change	China	3 m/e Aug 2012	3 m/e Aug 2013	% pt. Change
Android	62.7	56.3	-6.4	Android	63.2	72.4	9.2
BlackBerry	10.1	3.7	-6.4	BlackBerry	0.2	0.0	-0.2
iOS	21.4	27.5	6.1	iOS	29.5	20.8	-2.7
Windows	4.5	12.0	7.5	Windows	4.7	2.1	-2.6
Other	1.3	0.5	-0.8	Other	8.5	4.7	-3.8
France	3 m/e Aug 2012	3 m/e Aug 2013	% pt. Change	Japan	3 m/e Aug 2012	3 m/e Aug 2013	% pt. Change
Android	62.7	63.3	0.6	Android		48.6	
BlackBerry	8.2	4.2	-4.0	BlackBerry		0.3	
iOS	13.5	17.5	4.0	iOS		47.4	
Windows	5.6	10.8	5.2	Windows		0.8	
Other	10.2	4.1	-6.1	Other		2.9	
Italy	3 m/e Aug 2012	3 m/e Aug 2013	% pt. Change	Australia	3 m/e Aug 2012	3 m/e Aug 2013	% pt. Change
Android	58.6	71.6	13.0	Android	65.9	62.1	-3.8
BlackBerry	3.9	2.1	-1.8	BlackBerry	1.6	0.5	-1.1
iOS	15.0	14.4	-0.6	iOS	25.8	28.7	2.9
Windows	10.3	9.5	-0.8	Windows	3.7	6.5	2.8
Other	12.2	2.4	-9.8	Other	3.1	2.2	-0.9
Spain	3 m/e Aug 2012	3 m/e Aug 2013	% pt. Change	EUS	3 m/e Aug 2012	3 m/e Aug 2013	% pt. Change
Android	84.9	90.8	5.9	Android	68.8	70.1	1.3
BlackBerry	5.6	0.7	-4.9	BlackBerry	5.8	2.4	-3.4
iOS	2.8	6.3	2.6	iOS	14.1	16.1	2.0
Windows	2.1	2.2	0.1	Windows	5.1	9.2	4.2
Other	4.7	1.0	-3.7	Other	8.1	2.1	-4.0

Figure 2.5: Mobile phone market share September 2013

(Jones, 2013)

2.3.6 Currently available Cryptic Crossword apps

After performing various research on the apple app store, the Google play store and blackberry world, it was discovered that there is not many Cryptic Crosswords available to download. There were two, which could be clearly defined as Cryptic Crossword applications, and these have been analysed in the following sections.

2.3.7 Puzzler Super Cryptic Crosswords

Platform: Apple iOS

Price: £3.99

Compatibility: iOS 4.3 or Later

Website: <https://itunes.apple.com/gb/app/puzzler-super-cryptic-crosswords/id616060420?mt=8>

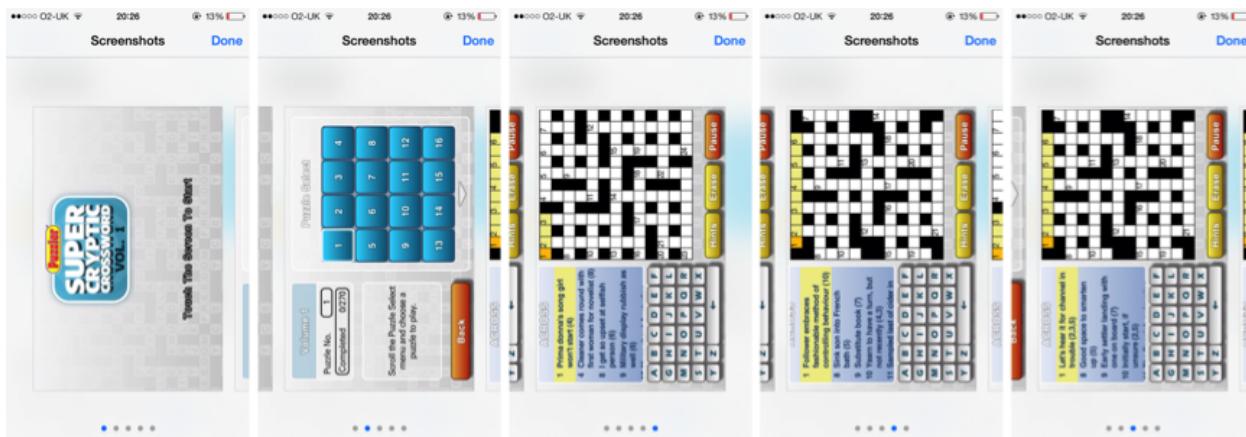


Figure 2.6: Screenshots of Puzzler Super Cryptic Crosswords on iPhone

This application contains 270 pre built crosswords. The interface is really easy to use and it clearly shows what has been completed and what needs to be completed. What's nice about this application is that the crosswords are in various levels and as you complete one crossword you can move on to the next. What was noticed is that the application already stores the answers and also the application allows hints, which makes it a little easier to use. The other noticeable thing about this application is that it didn't fully use the native features of the mobile phone. Like the keyboard is a custom keyboard.

2.3.8 Cryptic Crossword

Platform: Apple iOS, Android

Price: Free with 2 puzzles on Apple devices- In app purchase available - £1.99 on Android Devices

Compatibility: iOS 6.0 or Later

Website: <https://itunes.apple.com/gb/app/cryptic-crossword/id661608021?mt=8>

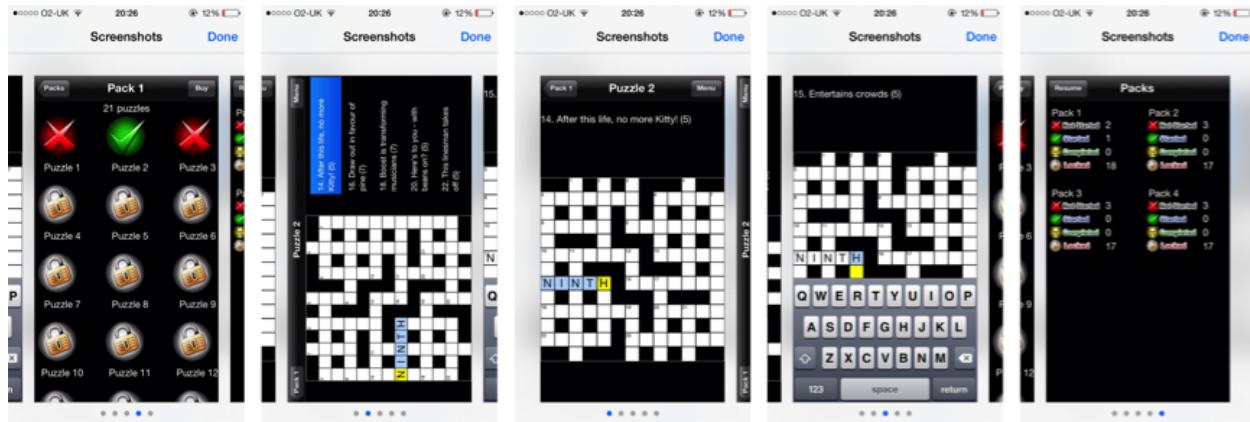


Figure 2.7: Screenshots of Cryptic Crossword on iPhone

The Cryptic Crosswords app comes with 4 different packages and a total of 81 puzzles. Each pack can be bought for £0.69 or all 4 for a discounted price of £1.99. Each package contains various cryptic crossword puzzles. After having a play around with this application at first point of contact it is noticed that look and feel of the application is not that great. The use of colours and the layout styles of the application can be better. The app consists of 81 puzzles, which can be played only after each crossword has been completed. Some of the features of the app are:

- Checking answers
- Revealing a letter, a word or the whole puzzle
- Clearing the puzzle
- Moving the character bar to next box
- Greying out completed clues

These are some of the features the application uses but there are plenty more.

2.4 Web Services

Over the past decade web services have exploded into the computing space. However the concepts that underlie web services are not new. Web services originally evolved from the Remote Procedure Call mechanism that was found in a software development framework used in the 1990s (Kalin, 2013).

During the late 1990s, XML-RPC was developed, which was a stripped down, light weight version of the Remote Procedure Call mechanism. The XML-RPC system only supported a small number of data types along with a number of simple commands. XML-RPC contained two key features, which are use of XML serialise/deserialise for data types and the reliance on HTTP for transport XML-RPC (Kalin, 2013).

XML-RPC is designed to be as lightweight as possible, and thus can be supported on a wide range of devices. XML-RPC was ultimately implemented fully and became known as SOAP. As well as SOAP, another implementation of XML-RPC occurred, which was entitled REST. Both of these technologies fall under the term ‘web services’.

Since 2001 a vast range of companies have adopted the web services movement including (but not limited to) IBM, Oracle, Hewlett-Packard, Amazon, Google, Facebook and Twitter (Sullivan, 2001; Kalin, 2013).

Web services generally tend to reside upon public networks such as the Internet. However it is possible for a web service to run within a private network, such as a company’s internal Intranet.

2.4.1 What are Web Services?

Although there are many companies adopting web services, the term web service has a diverse and loose definition (Kalin, 2013). During the initial explosion, many providers created heavily detailed plans upon the direction of their web service, but failed to exactly define what a web service is.

It was only until the explosion subsided that authors were able to define what a web service is (Kalin, 2013). Kalin (2013) highlighted three common characteristics between web service providers:

1. Can be thought of as a ‘webified application’
2. Typically delivered over Hyper Text Transport Protocol (HTTP)
3. Typically has some form of distributed nature allowing for components to be deployed and executed across multiple devices.

For the purposes of this project a web service will be defined as:

A service that contains one or more software components that are designed to allow machine-to-machine interaction over a network using standard protocols.

Web services follow the client-server model, which is the standard architecture for accessing a website. However unlike the traditional approaches to client/server models (such as a web server/web page setup), web services do not provide the end user with a Graphical User Interface (GUI).

The web service will provide the end user with machine readable data — i.e. the data must be put into a pre-defined GUI. This architectural design concept is not new and has been around for a number of years. Web services often can be thought of as imitating mainframes — i.e. a ‘dumb terminal’ sends a request to a service hosted upon a central computer system.

Web services can be broadly categorised into the distributed software systems category (Kalin, 2013). Broadly speaking a distributed software system is a system that is often split up into various components. Each component can run upon a separate physical machine, and is able to communicate with other parts of the system by passing ‘messages’ around. Although a web service does fit into that broad definition, there are several features that are unique to a web service.

Firstly web services heavily depend upon open, industry-standard, vendor-independent protocols such as HTTP, JSON and XML. By adding networking, data formatting and security features, web services can effectively lower start-up costs and promote interoperability between new and existing services (Kalin, 2013).

It is the interoperability that allows web services to promote language transparency. This means that web services and client programs do not need to be developed in the same

language. Many popular languages (e.g. C/C++, Java, and Python) provide inbuilt libraries or frameworks in support of web services (Kalin, 2013).

Finally web services are designed to be modular in design. This allows new services to be brought online in staggered stages, as well as allowing for laying of existing services. Again as previously mentioned each new service, can be written in the same language as the previous service, or use a completely new language (Kalin, 2013).

2.4.2 Web Service Categories

Web services can be divided into two distinct groups — SOAP based and REST-style (Kalin, 2013). Interestingly the distinction could be described as being little at most, but they are not necessarily directly compatible with each other.

SOAP

SOAP originally stood for Simple Object Access Protocol, but is often referred to as Service Oriented Architecture (SOA) Protocol (Kalin, 2013). At a glance the name change doesn't appear to be too trivial, but it is acutely an example of the technology becoming better defined (Kalin, 2013).

SOAP utilises concepts that can be seen throughout the industry, but none more so than the use of XML. One of the major advantages of XML, is that it is able to provide flexible, self-describing data structures that can easily be produced and read.

SOAP tries to imitate the postal system — i.e. allowing two machines to send and receive letters. In this analogy, the letter is the raw XML data, and the envelope is an additional data layer that wraps around the letter. The envelope adds additional information to the request, such as which operation is being requested, and may also include authentication and session information in envelope headers (Gershon, 2004).

In order to ensure one client or service can ‘talk’ to another service, SOAP responses must use a Web Services Definition Language (WSDL). The WSDL defines the inputs (e.g. parameters), the outputs, the operations, the protocols and the network addresses that are required and used by the service (Gershon, 2004).

The underlying implementation is loosely coupled with WSDL, which means the provider is able to change the implementation, without negatively impacting the end service users. It is the configurable services aspect that is the central concept behind all service oriented architectures (Gershon, 2004).

REST

REST stands for Representational State Transfer, and is a relatively new architecture for creating web services. Despite its relatively new architecture it is actively used by some of the larger vendors such as Google and Amazon (DOSPINESCUM and PERCA, 2013).

REST relies upon the emerging architecture known as resource-oriented architecture. Essentially, these resources are a number software components that can be combined together to create reusable functionality.

As well as using a resource-oriented architecture, REST makes clever and effective use of open standard web technologies, such as the Hypertext Transfer Protocol (HTTP), the Uniform Resource Identifier (URL) and the Extensible Mark-up Language (XML) (DOSPINESCUM and PERCA, 2013).

Although not all of the features have been implemented, (mostly because they are layout properties rather than data properties) the major concepts found in web technologies have been implemented and the most notable of these features are:

1. Data from the client is transmitted to the server via the URI
2. The server will perform the operation described by the HTTP method (such as GET, DELETE)
3. The URI for each resource will contain the server name and address

As previously mentioned, HTTP methods are widely used within REST. A HTTP method will describe the necessary action (Create, Read, Update and Delete — CRUD) that is required to be performed by the server (DOSPINESCUM and PERCA, 2013).

The HTTP methods follow another standard in terms of the basic functions of a database management system. It must be said that REST and the HTTP protocol are mutually exclusive — REST doesn't require HTTP (DOSPINESCUM and PERCA, 2013).

Table 2.6 describes common HTTP verbs and the associated CRUD operation.

HTTP verb	CRUD operation
POST	Create a new resource
GET	Read a resource
PUT	Update an existing resource
DELETE	Delete a given resource

Table 2.6: HTTP verbs mapped to the associated CRUD operation.

There are additional optional verbs, such as HEAD, TRACE, CONNECT, OPTIONS and INFO, but these may not be implemented by the server and/or service for security reasons. Every HTTP request will include a verb to indicate which CRUD operation should be performed upon the resource (DOSPINESCUM and PERCA, 2013).

SOAP vs REST

Both REST and SOAP utilise standard protocols when communicating, and also originate from a similar specification. The real difference between the two technologies is that SOAP utilises its own application protocol by extending current protocols.

This causes a number of issues, such as protocol standardisation. Although SOAP is based upon the HTTP protocol, each client will have to correctly understand the new extended protocol — via an additional layer of software or libraries. This adds weight to the overall technology.

SOAP describes functions, and the types of data, which requires large amount of documentation in order to use the service. As well as this there are several protocols and technologies that directly relate to it, such as Web Services Description Language, Web Servicing Addressing, XML Schema Definitions.

All binary data that is to be transmitted must be first encoded in a supported format (e.g. base64), which increases processing power at both the client and server ends. All requests are transmitted via XML, which is much slower to parse and interpret than other text-based human readable data, such as JavaScript Object Notation (JSON).

REST on the other hand is based upon uniform interfaces. This means the various clients will have a small understanding of the web service, but not necessarily how it operates or what it will return.

REST doesn't need to operate over HTTP, and doesn't contain the complexity that SOAP provides. Rather than utilising XML, REST uses the standard HTTP methods to describe what a service should do. For example obtaining a resource would use GET, and for creating a resource POST would be used.

Clients do not require additional REST supporting libraries. As long as the language supports HTTP, the client will be able to consume a REST HTTP service easily.

Unlike SOAP, REST can deliver binary data without having to encode, and responses can be formatted to either XML or the more popular JSON (due to speed increases).

2.4.3 Clients

As previously mentioned, the broader web server architecture follows the client-server application model. When designing a client-server application, a decision has to be made as to which operations (or parts of operations) should be performed upon the client and the server.

This decision is vitally important as it can affect the speed to which a system can be brought to market. It might also affect any additional extensions or updates that the system might receive in the future, as well as affecting the design flexibility.

In order to simplify the design the client will need to fall into one of the two categories — ‘thin’ client or a ‘thick’ client.

Thin

A thin client is a computer system that depends largely upon a main server, or a number of servers in order to complete any computation tasks. The client has no knowledge of how to process data, it simply knows how to pass data to another entity, and receive data from another entity.

A recent example of a thin client is Google’s Chromebook. Unlike typical computers where by the applications are installed locally upon the computer, the Chromebook allows for applications to be installed within the cloud — upon an external server.

The thin client design presents a number of advantages and disadvantages. Firstly an application that is hosted upon a central server can be easily updated — as there is only one code base. Once the application has been updated, this will be pushed immediately to all thin clients.

This obviously provides an advantage in some use cases such as trying to sell goods over the Internet. For example if a product’s price changes, the update will only need to be applied once to the central server, rather than having to update all clients wishing to purchase the product.

Thin clients will utilise powerful servers to do the majority of the processing. This allows for the thin clients to be less powerful, and hence the overall costing to reduce.

However this will mean that thin clients will have poorer response times. The main reason for this is the fact that the majority of the operations are being complete upon another machine (potentially many miles away). Simple operations such as populating a menu, might require a request to the main server, thus increasing the overall time to achieve something.

Resources within a thin client network will need to be managed more effectively. Thin clients will use more bandwidth upon the network, and will make more connections to the server. This would require the server to be able to handle lots of potentially fast and slow connections, with each connection using a wide range of internal server resources (CPU, Memory etc).

Thick

A thick client is a computer system that has little dependency upon a main server, or a number of servers in order to complete computational tasks. The client will still require a limited connection to a server, but will not use the connection as often in comparison to a

thin client. A thick click will often be able to perform many operations without a connection to a network.

An example of a thick client would be a standard desktop installation. The desktop installation might provide various pieces of software that are installed locally upon the computer. For example the computer would be able to produce various documents regardless of the state of the network connection.

The thick client design presents a number of advantages and disadvantages. Firstly due to the fact that clients are able to do more of the computational work, server specifications do not need to be as high. This allows for cheaper servers to be purchased, and few overheads in terms of running and maintenance costs.

This will also lead to an increase in server capacity, again due to the fact that the client is carrying out more work. This ultimately means that the server is required to do less work, and can hence support a larger number of users.

Thick clients have an increased advantage over thin clients in terms of network connectivity. Thick clients do not require a constant connection to a server. This in turn frees up bandwidth that is being used upon the network, as well as reducing server loads.

Finally the end user is able to store files and applications locally upon the machine. This in turn allows for a faster application start up time, and a reduced file access time. Hence increasing the speed of operations, as well as reducing bandwidth upon the network.

However thick clients are more expensive to purchase, deploy and maintain. The reason being is that there will be more computers with higher specifications. This can lead to more expensive repair bills, should systems fail.

Fixing and troubleshooting become more difficult, simply because there are more machines to troubleshoot and fix should problems occur. This is obviously not a problem if there was a central server, such as found within the thin client model.

Chapter 3

Development Methodologies

In order to ensure all objectives and goals that have been set within this project are completed to the highest quality and upon time, a software development methodology will need to be chosen. Dividing a larger project into a set number of defined processes may seem like additional unnecessary work, but the advantages of this process far outweigh the disadvantages (Knott and Dawson, 1999).

The defined processes combine together to form part of a process model. The process model will allow for the following achievements (Knott and Dawson, 1999):

- Adding an element of control and planning
- Allowing for progress to be mapped visually
- Providing a structured approach to development
- Allowing for a higher quality of code and documentation to be produced

The Systems Development Life Cycle (SDLC) was one of the first formalised methodologies for building software. The SDLC utilises a methodical and structured approach to analysing, designing, building and testing software, to which many methodologies follow this rigid structure (Elliott, 2004).

3.1 Waterfall

One of the main aspects to the waterfall model is the fact that the project is expected to progress down the primary path (Cadle et al., 2010).

The waterfall model takes the major components of any project (requirements, design, implementation, testing and maintenance) and assigns each component a stage of its own. Each component is delivered as the flow down the primary path is completed (Cadle et al., 2010).

The waterfall model also supports backtracking (i.e. reverting back to a previous deliverable). This allows for project managers to check that the project has not expanded its defined scope, and to also ensure that each deliverable flows into the next correctly. It also allows for slight modifications to be made, however making many large changes might affect the project in the long run (Cadle et al., 2010).

3.1.1 Advantages

Cadle et al. (2010) states that the waterfall model houses a number of advantages, including:

- Provides a rigid project structure, that is easy to follow and review
- Deliverables are delivered in project order, one at a time
- Can work well for smaller projects, or for projects where the requirements will not change.

3.1.2 Disadvantages

However Cadle et al. (2010) also goes on to state that the waterfall model houses a number of potential problems, including:

- Changes are difficult to implement the further a project is down its primary path
- Large projects may not benefit from the rigid structure
- A working piece of software is not delivered until late into the project

3.2 Spiral

A common feature found in the waterfall model is that all requirements are stated at the start of the project. It is these requirements that will form the basis of all work, along with any project planning (Cadle et al., 2010).

The spiral model forms its basis around iteration and prototyping to try to explore the requirements and develop the solution. During each turn around the spiral, a set of requirements are analysed and developed using prototyping (Cadle et al., 2010).

3.2.1 Advantages

Cadle et al. (2010) states that the spiral model houses a number of advantages, including:

- A high amount of risk analysis is conducted, and thus risk is more likely to be avoided
- The model allows for approval from clients, and large amounts of documentation to be produced
- Software can start to be produced earlier, in comparison to the waterfall methodology
- Additional functionality can be added on at any time during or after the project

3.2.2 Disadvantages

The spiral model allows for a high level of control, without too much restriction. However Cadle et al. (2010) states that this can cause difficulties such as:

- A thorough investigation into all of the requirements cannot be achieved early, therefore some requirements (and their priorities) may get completely missed
- The spiral model is based upon the clients knowing exactly what they want, which is unlikely
- A risk analysis must be conducted, and requires highly specific expertise to complete. If a risk analysis is not completed, then the project may completely fail

3.3 Agile

The agile software development methodology is designed to “reduce risk by delivering software systems in short bursts or releases” (Dawson, 2009).

Each release (sometimes referred to as iterations) will involve minimal planning and will cover all the major SDLC components: analysis, design, implementation and testing. The agile development model also heavily promotes collaboration/development between team members (Dawson, 2009).

3.3.1 Advantages

One of the main advantages of using the agile development model is that software is developed in rapid cycles, which ultimately results in smaller constant incremental releases of software. As well as this major advantage, Dawson (2009) states the following advantages of using the agile development model:

- The methodology surrounds the concept of regular face-to-face meetings as opposed to in-depth documentation
- Utilises a close working relationship between the client and the developers, thus providing continuous delivery of useful software
- Uses shorter, iterative time scales (usually weeks rather than months or years), which results in working software being delivered frequently
- Easily able to change the requirements at any stage (however late the changes are)

3.3.2 Disadvantages

However many of the disadvantages of agile development model are surrounded by the lack of a rigid documentation, as Dawson (2009) also suggests the following disadvantages:

- There is often a lack of emphasis on necessary documentation (user documentation, design documentation etc.), which is normally skipped to save time
- The uncertainty of a specification may lead to poor code and/or structure
- The project can become confused if the original specification is not clear from the start
- Some software deliverables can be difficult to allocate the correct amount of resources (time, effort etc.) at the start of the project

3.4 Rapid Application Development

The Rapid Application Development (RAD) model is an extension to the incremental development methodology. The RAD model states that all requirements should be treated as mini projects, and that they should be completed in parallel. Each of the mini projects are ran like a normal project, and hence time scales need to be adhered to (ISTQB Exam Certification, 2013).

Upon completion of the mini project, the customer is able to review the output, and provide value feedback regarding to the delivery and the requirements. RAD will follow a somewhat simpler primary path, allowing for business modelling, data modelling, process modelling, application generation, testing and turnover (ISTQB Exam Certification, 2013).

3.4.1 Advantages

ISTQB Exam Certification (2013) states that there are many advantages of adopting the RAD model within a team:

- A reduced development time, due to the fact that the business modelling and data modelling processes should cover all aspects
- The combination of Data modelling and Process modelling should allow for the increased ability to reuse components
- Reviews of delivered outputs are constantly reviewed by the customer, allows for early feedback to be gained
- Parts of the system are integrated at an earlier stage, which allows for fewer integration issues towards the end of the project

3.4.2 Disadvantages

However, ISTQB Exam Certification (2013) also states that there can be disadvantages of adopting the RAD model within a team:

- There is a high dependency upon an overall strong team and strong individual performances for identifying business requirements
- The model will only work for systems that can be modularised
- The model assumes that the team members are highly skills designers and developers, with an even higher dependency upon modelling skills

3.5 Summary

In order to achieve the best possible product, it is clearly evident that the project should be developed utilising a feature driven approach. This will allow for any revisions, modifications, and changes to be considered and implemented with as little delay as possible, as well as little impact upon the rest of the project.

The projects requirements are not set directly by an external client, and hence it is possible for the requirements to be changed. It is because of these uncertainties that an agile development methodology would be best adopted by this project.

This methodology will not only allow for the requirements to change, but can allow for substantial research to be able to take place upon new topic areas if needed. Agile development methodologies allow for multiple releases of software, which fundamentally means that the team is able to use prototyping techniques to find the best outcome to a given problem.

Chapter 4

Software Specification

Robertson and Robertson (2013) state that ‘a requirement is something the product must do to support its owner’s business, or a quality it must have to make it acceptable and attractive to the owner.’ The aim of gathering the requirements is ensure that all ambiguities are removed before a product is developed.

Within this chapter the project’s aims and objectives will be outlined, as well as project constraints and assumptions. Software requirements will be categorised as functional and non-functional requirements, to which a full MoSCoW analysis will be conducted.

Finally a risk assessment into the various project related risks will be conducted to help identify any potential risks. The risk assessment will also identify how the risks can be reduced.

4.1 Aims & Objectives

4.1.1 Aims

- Produce a software application that accepts the input of cryptic crossword clues from a user
- Produce a software application that will present a set of results to a user for the cryptic crossword clue which has been input
- Effectively adhere to a project plan as a team
- Effectively abide by the guidelines of a chosen software methodology
- Construct systematic and comprehensive documentation for the academic report

4.1.2 Objectives

- To research cryptic crosswords themselves and the various types of cryptic crossword clue
- To study relevant software architectures associated with the group project
- To investigate the field of artificial intelligence, in particular natural language processing and justify the need for it within the project to solve cryptic crossword clues
- To research the various software methodologies which exist and select the most appropriate for the project
- To understand and effectively follow a collection of requirements
- Evaluate the progress of the project at numerous instances throughout the life cycle of the project

4.2 Constraints

The following section describes the constraints that effect the design of the product. The product that is to be developed cannot be successful unless these constraints have been accomplished.

- The project will have four developers operational throughout the life cycle
- The user interface to be developed for the software application must be developed with considerations regarding its simplicity for the usage on mobile phones
- Resources to test the software application on mobile phones will cover the following platforms:
 - Android
 - BlackBerry
 - iOS
- The product shall require an Internet connection
- The product shall make use of the Apache OpenNLP Library
- Possible external resources to be used for the project, such as a dictionary and a thesaurus, may have limits on the usage
- The product shall be completed on or before the 9th May 2014
- The product shall be developed using all freely available tools and the only constraint of budget shall be the time the project resources put in to the product

4.3 Facts and Assumptions

- Answers to Cryptic Crosswords are usually published the following day
- Two clues that are identical, may not lead to the same solution
- Equipment needed to test the application on different platforms will be available from the University
- A server to host the web service will be available from the University
- An external resource (The Guardian) has given permission to use their cryptic crossword data to use for our test data
- Four team members working on the project throughout the year at approximately 25% of their academic study time
- Project scope will remain the same throughout the project

4.4 Scope

4.4.1 Organisation Structure

The project team is based largely upon democratic discussions and decisions, however to ensure that team deadlocks do not occur a project leader has been chosen. Figure 4.1 reflects the hierarchy of the team.

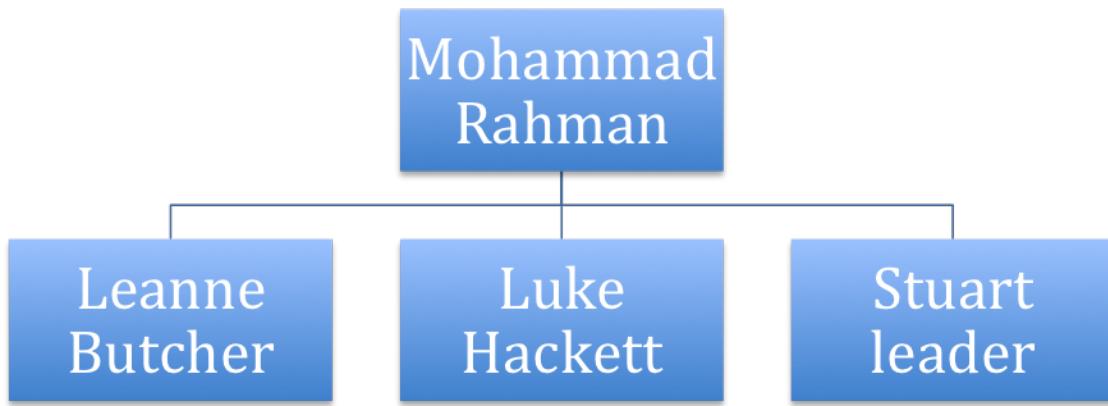


Figure 4.1: Hierarchical Structure of the team

4.4.2 Methodology

It has been decided that the team will follow an Agile software development model. This will allow the team to split the larger task down into smaller, more manageable ‘chunks’, that allows for a good quality analysis, evaluation, development and planning (on to the next ‘chunk’). An iterative approach is best suited for this project due to the nature of changes and updates the product will require in future builds.

This method of development also allows for a more feature-driven approach to the project. Ultimately this allows for more important features and aspects of the project to be completed first.

4.4.3 Meetings

A weekly meeting will take place between all project members to discuss all aspects of the project. This includes (but is not limited to) project issues, software development issues, research findings, possible improvements and code reviews.

4.4.4 Product

Within the cryptic crossword area, there is a need for a cryptic crossword clue solver in the form of a software application. This has been justified previously within the document through research and investigating existing applications within the field. The main purposes of the software deliverable are to allow the input of a cryptic crossword clue by a user and output a potential results or number of results.

Once the project has been completed by the team, the following deliverables will have been accomplished:

- A software application which accepts input from the user and outputs appropriate solutions.
- Two written reports that:
 - document the entire software development process from a group perspective
 - analyse and evaluate the project as a whole

Furthermore the subsequent internal components will be implemented to aid with the goals of the project:

- A storage area to collect and store the cryptic crossword clues and their details used for test data as well as, clues successfully solved by the application which a user has input.
- A service which will allow connections from web browsers and potentially mobile phones to ensure the running of the software application.

There are specific criteria which have been identified as important features which must be completed to deem the project as adequately finished:

- A document outlining the full software development life cycle of the project
- A software application which can be accessed via either:
 - A web browser **and/or**
 - A portable device (e.g. smartphone or tablet)
- A service which can be connected to by either a mobile phone or a web browser which will sufficiently solve a cryptic crossword clue and output a result or number of results

- A storage area accessible by the service which stores the cryptic crossword clues and their associated details

Furthermore, there are specific criteria which have been deemed as out of the scope of the project:

- The user will not be able to access the storage area to browse through data
- The software application will not be implemented to have the ability to generate cryptic crosswords from the data stored

The restrictions listed below are the justifications for the project scope:

- The time scale of the project
- External priorities from the academic year such as other module assignments and examinations
- Limited resources for testing the software application on restrict the platforms which the software application will be fully compatible with

4.5 Functional Requirements

4.5.1 MoSCoW Analysis

Functionality	Must	Should	Could	Won't
Retrieve and process cryptic clues and output possible solutions	Yes	-	-	-
Obtain characteristics of the clue's solution (e.g. number of words, word lengths and known letters) and match these against possible solutions	Yes	-	-	-
Provide a web interface for users to interact with the system	Yes	-	-	-
Provide a mobile-friendly interface for users to interact with the system	-	Yes	-	-
Store clues and their corresponding solutions for future retrieval, including the clue type and solution length	-	Yes	-	-
Take feedback from a user on a solution's accuracy and use this to rank solutions for a given clue	-	Yes	-	-
Relay the process (solution trace) followed to arrive at a proposed solution to the user	-	Yes	-	-
Holistically determine a confidence score for each proposed solution, and relay this to the user	-	Yes	-	-
Document and store the process (solution trace) followed to solve a given clue	-	-	Yes	-
Solve clues which take cues from the clue's orientation or clue number in the crossword	-	-	Yes	-
Provide real-time feedback of the processes being followed by the application in an attempt to calculate the correct solution	-	-	Yes	-
Take annotations provided by the user of the known aspects of a clue, such as the definition word(s), fodder or indicator word(s)	-	-	Yes	-
Store a clue's orientation and clue number in the database	-	-	Yes	-
Use stored data to generate complete cryptic crosswords	-	-	-	Yes
Provide a personalised service, including a history of a user's interactions with the system	-	-	-	Yes

Table 4.1: MoSCoW analysis of the project's functional requirements

4.6 Non-Functional Requirements

4.6.1 Performance

The performance of the system should be good, however it must be noted that the system is intended to run as a web service upon a server. The quality of the connection between the end user and the web service will be out of the scope of this project.

4.6.2 Legal and Access

Within the project a number of data sources will be used. Firstly an English dictionary will need to be accessed in order to allow the system to deduce if a given arrangement of characters formulates an English word.

In order to achieve this, an offline dictionary and an online dictionary will need to be used. The offline dictionary is a short dictionary and is provided as part of the GNU/Linux operating system under the GNU Licence Agreement.

A number of cryptic crossword clues and solutions have been obtained from the Guardian's website. The data will be used as a training set of data, when testing the final software. The Guardian has given us permission to use their data as long as the project remains within an academic environment, and that no profit is made from the project.

4.6.3 Backup and Recovery

The project will make use of the git revision control system over a secure shell connection. The service provider is github, who offer a distributed setup to ensure that data is always available. Using a secure revision control system will ensure that:

1. A copy of the project is stored upon a remote, secure server;
2. Changes are able to be tracked;
3. Issues and comments are able to be raised in a secure environment.

The above measures will in effect be able to reduce against the loss, damage or theft of all project information.

4.6.4 Archiving and Retention

The University will retain a copy of the project (software product and the final report) for assessment and evaluation purposes.

4.6.5 Maintainability

The project – including the software product – will be maintained up until the project hand in deadline: 9th May 2014.

4.6.6 Availability

The software will not contain any request or time limitations for users, and thus should be available at all times. However it must be stated that the server availability that will be hosting the web service is out of the scope of the project.

4.6.7 Capacity

The end product should be able to handle a number of requests from different users at the same time. As the project is only intended as an academic project, large scale user support will be out of the scope of this project.

4.6.8 Usability Requirements

Functionality	Must	Should	Could	Won't
Provide a text field for the user to input the cryptic clue to be solved	Yes	-	-	-
Provide a drop-down box allowing the user to input the number of words in the solution, with an initial upper-limit of 10 words	Yes	-	-	-
Dynamically provide individual text boxes for each word of the solution, which allow the length of the words to be specified. Also provide check-boxes between these text-boxes to define whether they are separated by a space or a hyphen	Yes	-	-	-
Dynamically provide text boxes to represent each character of each word of the solution, allowing the user to input any known characters	Yes	-	-	-
Provide a table of results which allow the user to view the possible solutions	Yes	-	-	-
Alert the user to required fields with red asterisks	Yes	-	-	-
Have a consistent layout avoid unnecessary scrolling	-	Yes	-	-
Alert the user to invalid input through validation checks with error messages	-	Yes	-	-
Provide a confidence rating in the table of possible solutions	-	Yes	-	-
Allow the user to select a proposed solution in the corresponding table and mark this as correct	-	Yes	-	-
Display help buttons to indicate to the user the purpose and use of each control	-	-	Yes	-
Provide a button to submit the user input to the application for processing	-	-	Yes	-
Provide a group of radio buttons for the user to select the clue's orientation in its containing crossword	-	-	Yes	-
Provide a text box to input the clue's number within its containing crossword	-	-	Yes	-
Provide mechanisms to accommodate for users with difficulties, such as colour-blindness or poor eyesight	-	-	Yes	-

Table 4.2: MoSCoW analysis of the project's usability requirements

4.7 Risk Assessment

An investigation into the potential risks that may present themselves during the course of this project will allow the group to effectively minimise their impact should they become a reality.

Many facets of risk management contribute to the potential for a project that will be completed on time and with little interruption from unexpected occurrences. According to McManus (2003), not only will risk management improve the flow through the project due to a reduced overall project risk, but fewer surprises will allow for a more accurate schedule, and hence a greater likelihood of finishing on time.

4.7.1 Risk Identification

Risk (Threat)	Description	Prob (%)	Impact (/10)	P-I Value
Workload constraints	Fluctuations in the workload from university modules mean that certain times of the year are typically busier than others, as time and effort have to be delegated to the tasks at hand. This is especially applicable when assignment deadlines are imminent.	80	7	5.6
Underestimated system complexity	Without a thorough investigation into the project's requirements and comprehensive designs, it may become evident that the team has not accurately ascertained the complexity of the task at hand.	60	6	3.6
Project scope incorrect	A misunderstanding of the scale of the project may present significant problems and have adverse effects on the team's ability to produce a quality product on time.	40	6	2.4
Lack of understanding	The team may simply not understand a given task and be able to continue with the project in the expected fashion.	40	4	1.6
Inadequate facilities/resources	Potential resources that may be required include web servers, database management systems, phone app development kits and dictionaries / thesauruses. Seamless access to each and every resource may prove difficult.	70	2	1.4
Ineffective team structure	Explained elsewhere in this document, the team have adopted a largely democratic model. As pressures and workloads increase, conflicts may arise and the team structure may hamper the group's ability to resolve any issues effectively.	40	3	1.2
Skills mismatch	The team may find that the skills required to complete the project are beyond those which are currently held. This has been identified as a risk as the group's collective experience is relatively low.	15	4	0.6
Workforce reduction	The most likely cause for this risk to become a reality is if a team member leaves the university course. As there are only 4 members working on the project, this would equate to a 25% reduction in the workforce.	5	8	0.4

Table 4.3: Probability-impact table for the project's risks

4.7.2 Risk Classification

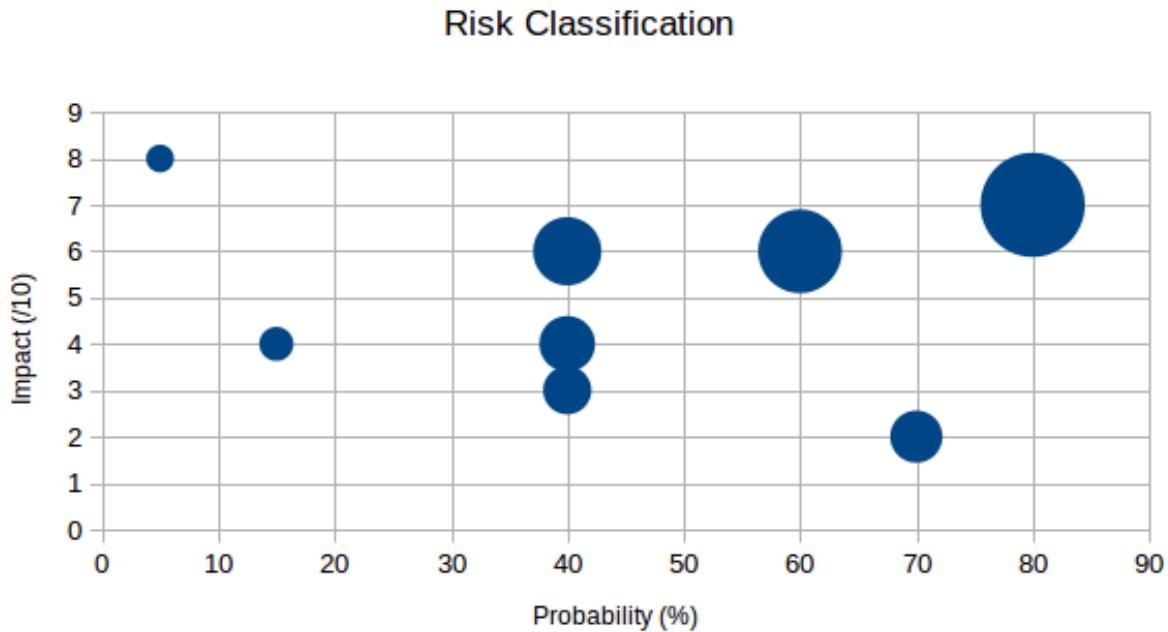


Figure 4.2: A visualisation of the identified risks of the project

Figure 4.2 demonstrates the impact and likelihood of the identified risks occurring. The risk with the greatest PI (*Probability / Impact*) score is the risk of constraints from other university commitments interfering with the expected progression of the project. This is further compounded by the fact that each member of the team is enrolled on the same course of study, and so the impact from such commitments will be magnified across each and every member.

In contrast, the risk with the lowest identified PI score is that of a group member leaving the course, and hence the project. While this has the largest impact score identified, the risk of such an occurrence is minimal. This is categorised by Stoeller (2003) as an *alligator*. This is defined as a risk with a high impact, but a low probability of becoming a reality.

4.7.3 Risk Prevention

Workload constraints The appearance of this risk balances on each team member's ability to effectively manage their time. If a consistent effort is maintained throughout the course

of the academic year, there will likely be a greater chance of this risk being prevented. Otherwise an sudden increase in workload will detract from the team's collective ability to work on the project. Project planning should include contingencies for unexpected occurrences, especially around the time of pressures or deadlines from other course modules.

Underestimated system complexity A thorough set of design documents based on the gathered requirements will minimise the potential for this risk to show itself. A range of diagrams from use-case diagrams to sequence diagrams will allow for a greater understanding of the identified tasks, thus reducing the risk of encountering unexpected complexity.

Project scope incorrect & Lack of understanding A recurring theme of this project is that a lack of planning may turn into magnified issues further down the line. The team members should be able to gather a reasonable understanding to the project's scope during the earlier phases. Conclusions may be drawn from the research and requirements analysis phases which hold some indication to the potential size of the project.

Inadequate facilities/resources The team will need to draw on a host of resources in order to efficiently work on the project. Repository and versioning systems will be used to maintain documentation and a code base. Development environments and development kits will be used to aid the software development, and servers and frameworks may be required to host any applicable software. To prevent the risk of issues arising from any of these software packages, it may be best to forecast which software will be required along with an estimation of when this is likely to be. Consequently provisions can be made in advance to ensure the software is set up.

Ineffective team structure Despite a unanimous vote for a democratic structure within the group, a team leader has been selected (*Mohammad*) who will hold the capacity to make an overriding decision on the project if the team cannot come to an agreement unaided. It is also the responsibility of Mohammad to delegate work where he feels it necessary. Such actions may collectively contribute to minimising the impact of this risk, but they may not allow for it to be entirely prevented.

Skills mismatch To ensure the team possess all the skills required to create the proposed solution, the gathered requirements, and in particular the technical requirements will outline the technologies earmarked to be used. As this phase is positioned relatively early on in the project, there should be relatively little chance of realising we require skills we do not possess further on in the project.

Workforce reduction A decision to leave the course before completion will likely be one that will not be taken lightly by the respective student. Consequently this is a threat for which there is little that could possibly be done to prevent it, other than allowing for an amiable working environment amongst each team member.

4.7.4 Risk Mitigation

Workload constraints Effective project planning should prevent this risk from manifesting itself in the project, but should it occur, it may be possible to balance the project workload more effectively amongst the team members. It may be possible that the anticipated workload of each team member can be mapped, and planned work on the project can be temporarily postponed until these constraints no longer present an issue.

Underestimated system complexity An assessment of the implications on the project should be conducted, and the effects investigated. If it would not be possible to incorporate the additional complexity using the time and resources available, the project scope may be reduced to allow for a smaller, but fully functional product. Research into development methodologies suggest that an agile methodology will be adopted. This will allow for a fast reaction to any change seen in the scope or structure of the project.

Project scope incorrect The group should draw on experiences from encountering this risk, and reflect on why it became a reality. A revised project scope should then be created, using these deductions to ensure the risk doesn't manifest itself twice.

Lack of understanding & Skills mismatch The obvious solution to mitigate this risk is to acquire the understanding necessary to carry forward with the work. Contingencies in the project plan may allow for such tasks to be carried out with little impact on the overall schedule. A variety of resources should be utilised to ensure the most effective acquisition of the knowledge necessary to proceed.

Inadequate facilities/resources In the case that a particular resource is unavailable, the team should conduct an assessment of why this is the case and what can be done to rectify the issue. If the issue is a delay in the time it takes to acquire a particular resource, the most viable option could be to search for alternatives or rearrange the schedule to complete other tasks which are non-dependent on this particular task.

Ineffective team structure If the chosen team structure is proving ineffective, it may be necessary to adopt a different set-up for the benefit of the group. With the team members

working with one another over a span of months, each members' working style in a group environment should become apparent, and this may help to shape any new structure.

Workforce reduction The time frame for the project's completion is fixed, but a reduction in workforce may necessitate a respective reduction in project scope or complexity. For the remaining three team members to attempt to complete the work set for four would possibly result in a product which is either unfinished or below standard. While it may prove frustrating to reduce the scope, this should allow for a subset of the complete solution, but which is just as refined.

4.8 Feasibility Study

A feasibility study will be carried out to identify the difficulty of the implementation needed for each type of cryptic crossword clue with the aid of the research carried out previously. The study will focus on the resources needed for each implementation which will contribute to the level of difficulty for each clue type. Furthermore, using the previous research resources within the cryptic crossword field, a measurement of how common each clue type is within a cryptic crossword will be presented to assist in the understanding of the necessity of the clue type within the project.

Aspects such as how common clue types are featured within research resources and how regular indicators are found in the project test data will contribute to the judgement made on how regular a type of clue is. The justification of the use of specific resources for a clue type will come from research resources and inspecting project test data.

Clue Type	Possible Resources	Difficulty	Regularity
Hidden	Dictionary	Low	Common
Anagrams	Dictionary	Low	Common
Acrostics	Dictionary	Low	Intermediate
Pattern	Dictionary	Low	Intermediate
Homophones	Homonym Dictionary	Medium	Common
Charades	Abbreviations, Dictionary, Thesaurus	Medium	Common
Deletions	Abbreviations, Dictionary, Thesaurus	Medium	Common
Reversals	Abbreviations, Thesaurus	Medium	Common
Palindromes	Dictionary, Thesaurus	Medium	Intermediate
Double Definition	Dictionary, Thesaurus	Medium	Intermediate
Substitutions	Abbreviations, Dictionary	Medium	Rare
Shifting	Dictionary, Thesaurus	Medium	Rare
Exchange	Dictionary, Thesaurus	Medium	Rare
Spoonerisms	Dictionary, Thesaurus	Medium	Rare
Containers	Abbreviations, Dictionary, Thesaurus	High	Common
Purely Cryptic	-	High	Common
& lit	Abbreviations, Dictionary, Thesaurus	High	Intermediate

Table 4.4: Feasibility Study for Clue Types

Chapter 5

Design

This chapter will describe the various design decisions that were taken in order to allow for the system to be developed. The chapter will present a number of artifacts that are expressed using the Unified Modeling Language (UML).

Additional non-UML diagrams will also be presented to describe the database and user requirement design decisions.

5.1 Primary Path

Before being able to correctly design a large, complex system a good understanding of the main scenario through through the system will be required. The main scenario is “a sequence of event or actions” (Lunn, 2003) in which many smaller scenarios may occur.

The main scenario of the proposed cryptic crossword solver system is outlined below:

1. Input the cryptic clue to be solved
2. Input the pattern of the solution
3. Process the data based upon a number of algorithms
4. Output the most confident results

Based upon the previously given main scenario, the primary path of the entire system can be proposed. The primary path is defined as the most commonly used ‘route’ though a system with as few variances as possible (Lunn, 2003).

The proposed primary path for the cryptic crossword solver system is outlined below:

1. Input the cryptic clue to be solved
2. Input the solution length
3. Input the pattern of the solution
4. Determine the type of clue that has been given (e.g. anagram)
5. Run the given clue type algorithm
6. Rank each of the results based upon a pre-defined criteria
7. Output the top ranking results

The primary path outlined above takes a relatively generic approach with regards to how a clue should be solved. Ideally the clue would be categorised into a certain type of clue, thus requiring that particular solver to be run.

However the primary path does not illustrate the steps that are required to be taken if the clue can not be categorised — does the system simply stop or would it try an alternative approach?

These additional steps may often be ‘optional’, and hence are described as alternative paths. The alternative path will contain small and subtle changes to the primary path (Lunn, 2003).

If the system is unable to determine the type of clue that has been given then an alternative approach is to be used. In this instance, the primary path would house the following additions and changes:

3 Determine the type of clue that has been given (e.g. anagram)

3.1 Start a new instance of each solver

4 Run the each of the solvers separately

The above alternative path shows that if the clue can not be categorised, then the system will attempt to solve the clue utilising a brute force approach. The above alternative path, could also be regarded as a ‘worst case scenario’.

The outlined primary and alternative paths are a simple way of illustrating the the general ideas behind the system. These paths will be actively referred to throughout the design process.

5.2 Activity Diagrams

Within the Primary Path section a primary path analysis was conducted, ensuring that common paths throughout the proposed system are identified. Within this section a number of activity diagrams will be presented that will break the primary path into a number of simple processes.

An activity is a task in a process that allows for small meaningful work to happen. A process may have a number of activities linked together to form a work flow (Lunn, 2003).

Figure 5.1 on page 78, illustrates the main activity diagram for the complete system. The activity diagram follows what was discussed within the Primary Path section, and expands upon some of the generalised points. Over the following subsections key parts of the main activity diagram will be discussed in more detail.

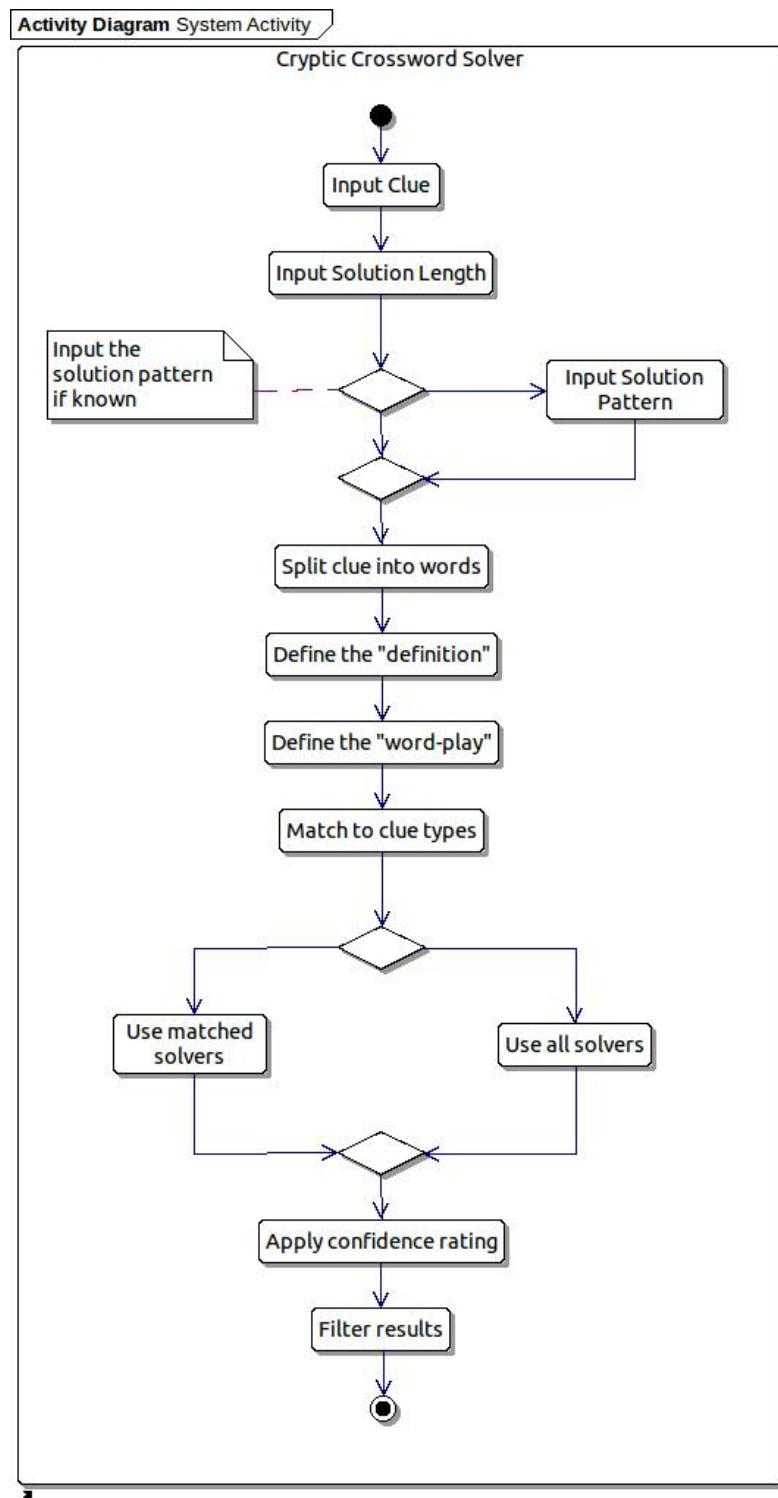


Figure 5.1: The complete Activity diagram

5.2.1 Clue Input

Figure 5.2 illustrates how the data should be input into the system. In order for the system to solve a given clue it requires the following three inputs:

- the clue
- the solution length
- the solution pattern

At first it may seem that the solution length is unnecessary, as it could be computed from the solution pattern. The reason for including the solution length is that it provides an additional validation element — much like having to type in a password twice when creating a new online account.

The solution pattern allows a user to predefine the required answer, based upon a number of characters. The solution pattern will accept all letters as well as three predefined ‘special characters’, as outlined below:

- ? ⇒ wild card character
- ⇒ hyphenated word
- , ⇒ word separator

For example, the pattern “????e?-???” could match to “mother-in-law” or “father-in-law”, just as “?????? d????” could match to “sitting ducks”. As previously mentioned, a user may simply present the correct number of wild card characters to denote an unknown pattern.

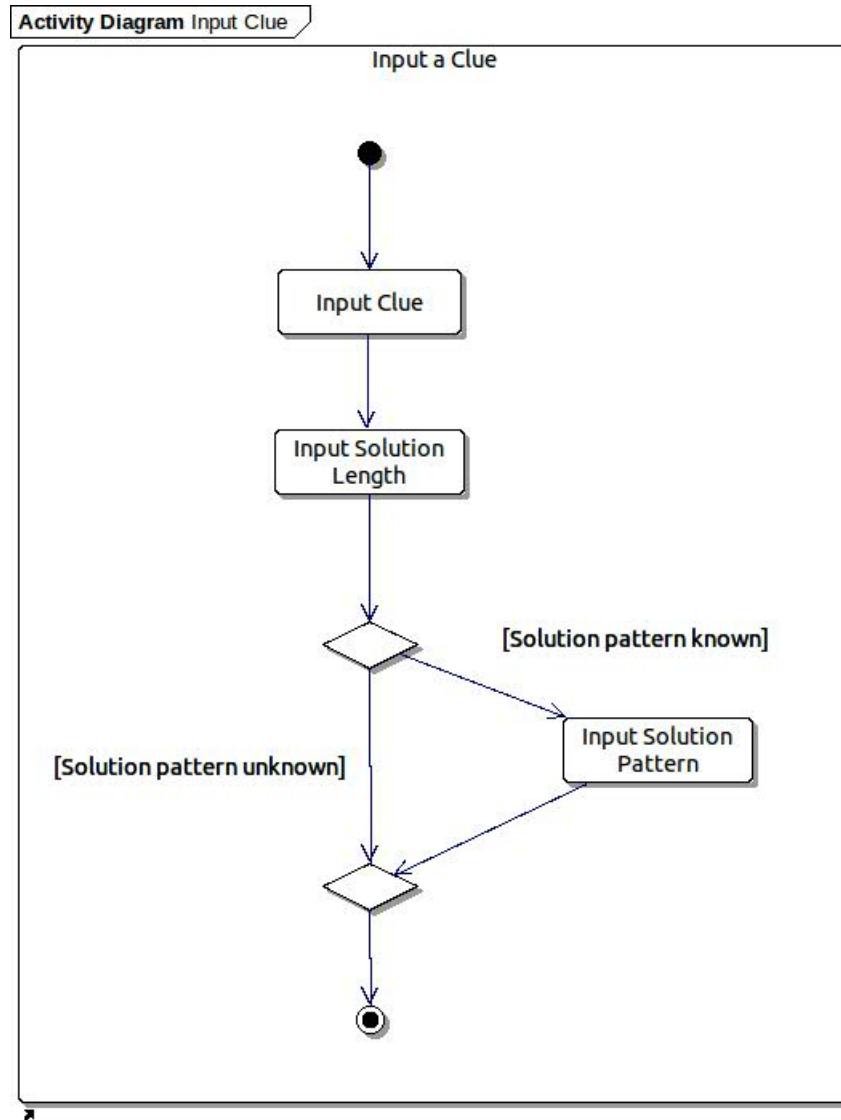


Figure 5.2: Activity diagram showing how a user would input a clue

5.2.2 Choosing a Solver

Once the user has successfully input the three variables, the system will then try to determine the type of clue. There are a number of actions that must be performed in order to deduce the type of clue, as shown in figure 5.3.

The system will try to match the given clue to a list of properties that each clue has. For example spoonerism clue types, will have the word “spooner” within the clue.

If the match is successful then the system will bump the ranking of that particular clue type. Once the system has finished applying the categorisation ‘tests’ to the clues, it will select

the top solver to try to solve the clue.

If the chosen solver is unable to solve the clue, then the next ranking solver will try to solve the clue. This process will continue until a set of results have been returned, or if all solvers have been used.

If the categorisation fails for some reason, then the system will try to brute force it's way to answer, by using all solvers.

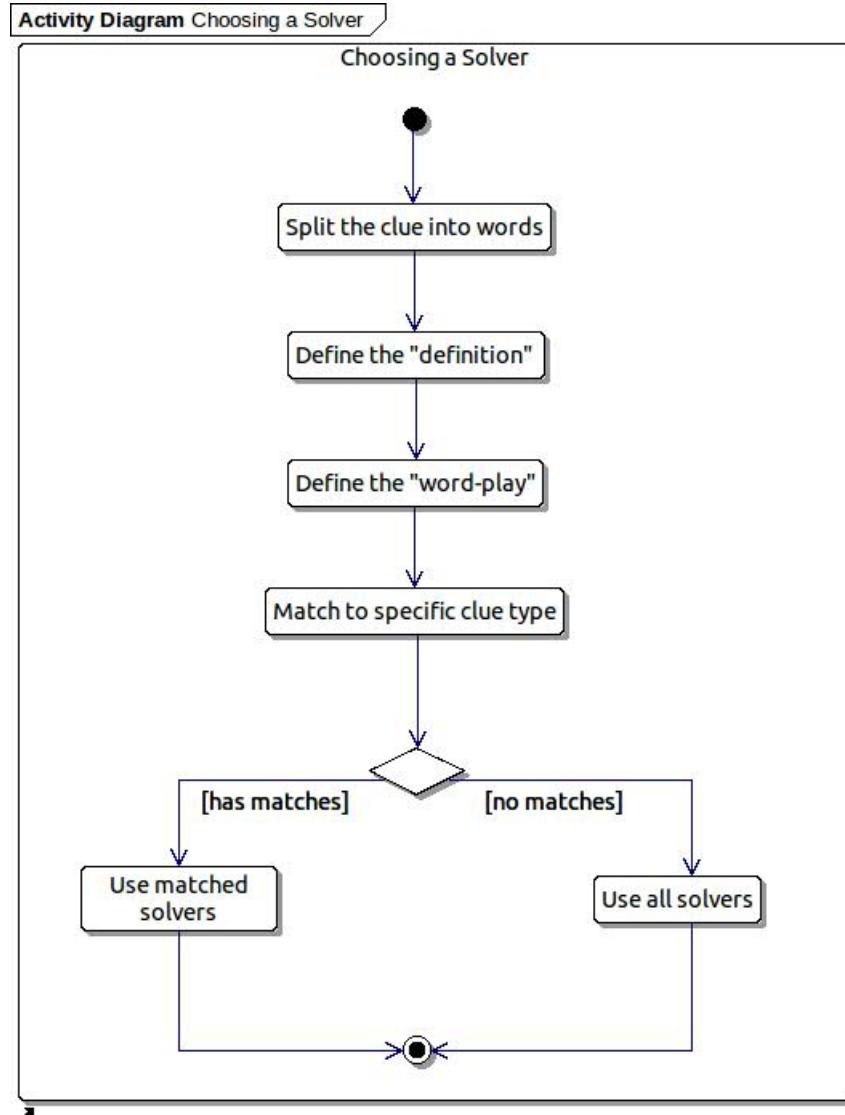


Figure 5.3: Activity diagram showing how a solver is to be chosen

5.2.3 Handling Results

Once the solving process has finished, the system will then attempt to apply a ‘confidence rating’. The confidence rating of a solution is how close the system ‘thinks’ the solution is to the correct answer.

For example, if a solution had a confidence rating of 100% it would be deemed as the only answer, whilst 0% would indicate there is no chance of the solution being correct.

It is intended that if there are a large number of results, then the system may choose to filter the results down. For example if thousands of results are generated, then the system may choose to only return results that have a confidence rating higher than 70%.

Figure 5.4 shows the handling of results as an activity diagram.

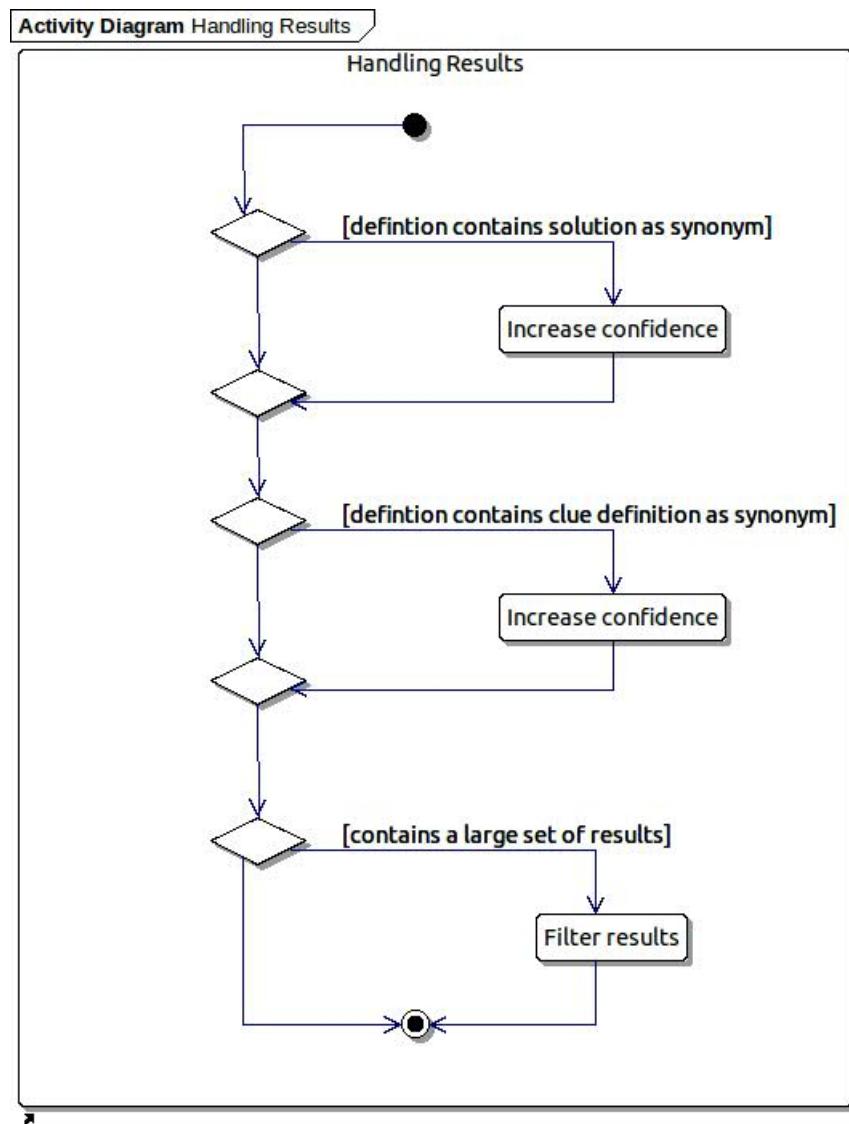


Figure 5.4: Activity diagram showing how results should be handled

5.3 Use Case

A Use Case diagram illustrates “a meaningful interaction with a computer system”, and provides various approaches to analysing system design requirements (Lunn, 2003). Within this section the use case of the entire proposed system will be presented and discussed. The use case has been split into two diagrams to increase the readability of the diagram.

The primary intended use of the system is to be able to submit a clue to be solved, which requires the clue, solution length and solution pattern. Figure 5.5 illustrates the three parameter requirements in order to submit a clue to be solved.

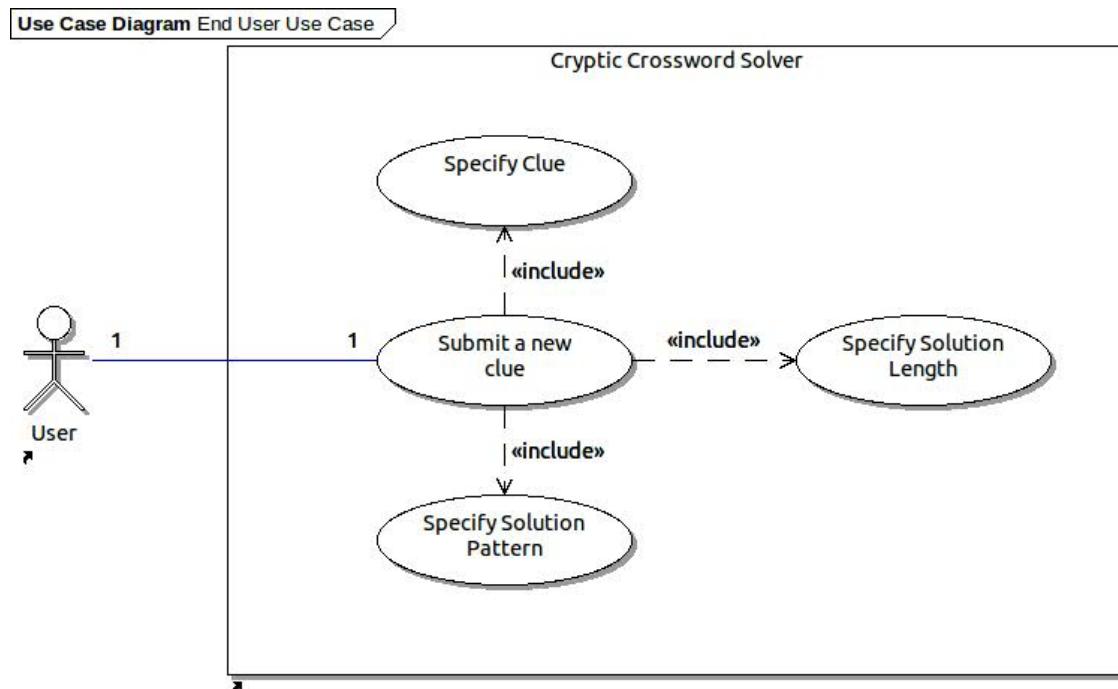


Figure 5.5: Use case illustrating the actor's range of possibilities

It must be stated that the solution pattern could be unknown. In this instance a set wild card characters must be given to the system. This was previously shown in figure 5.2 on page 80.

The remainder of the Use case is shown in figure 5.6. Once the user has input the various data elements the system will then attempt to solve the clue.

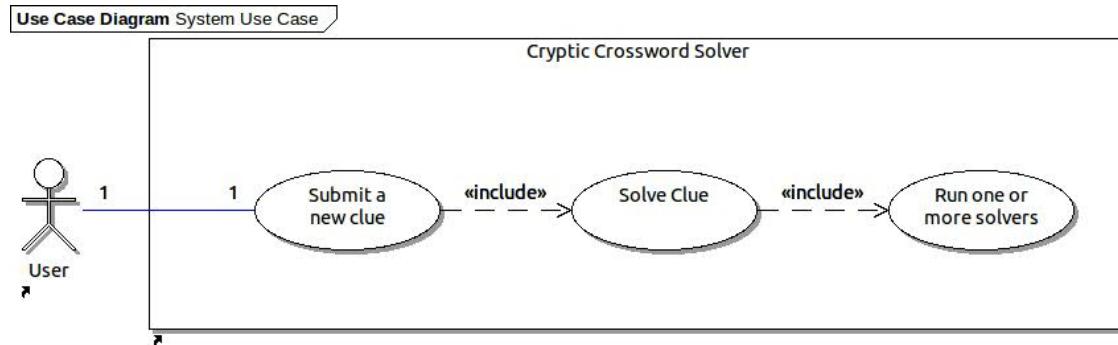


Figure 5.6: System Use Case

The use case demonstrates that there is very little user input required, and once the data parameters have been entered there will be no additional user interaction required.

A design decision has been taken so that there is minimal user interaction with the service. This will be particularly useful for those who are using a mobile device, as inputting data can be laborious in comparison to a desktop machine.

5.4 Sequence Diagrams

Within the previous Use Case section a number of key, generic functions were identified. In order for a sequence diagram to be developed, the generic functions will be converted into a set of objects and interactions. Each of the objects and interactions will eventually form the classes that will be used to build the system. Therefore a sequence diagram can be described as a “way of describing a journey through a system” (Lunn, 2003).

The two sequence diagrams that will be described within this section focus upon submitting a clue and solving a clue.

5.4.1 Submitting a Clue

Figure 5.7 illustrates a user sending data to the web service (servlet). Within this diagram a number of key elements such as how a clue will be solved have been removed. This is to ensure that the diagram is as simple as possible, as this diagram will only be focusing upon how a user sends and receives data.

The Solving a Clue subsection on page 87 describes in detail how a solver will solve the clue.

To simplify the diagram the user is submitting a ‘block’ of data, which in fact would be the three main parameters — the clue, the solution length and the solution pattern. This has been shown upon the diagram as a ‘block’ of data, because all three parameters will be sent to the system at the same time.

Once the data has arrived at the servlet, the data will be validated. If the data is deemed to be invalid, the servlet will send an error response back, finishing the current instance.

If the data is deemed to be valid, then the system will try to categorise the clue, and solve the clue using the required solver(s). This is will be discussed in more detail in the Solving a Clue subsection on page 87.

Once a set of solutions have been generated, the the in-memory set will be converted to either JSON or XML depending upon the client’s request. Once converted the data is returned as part of the response, satisfying the client’s request.

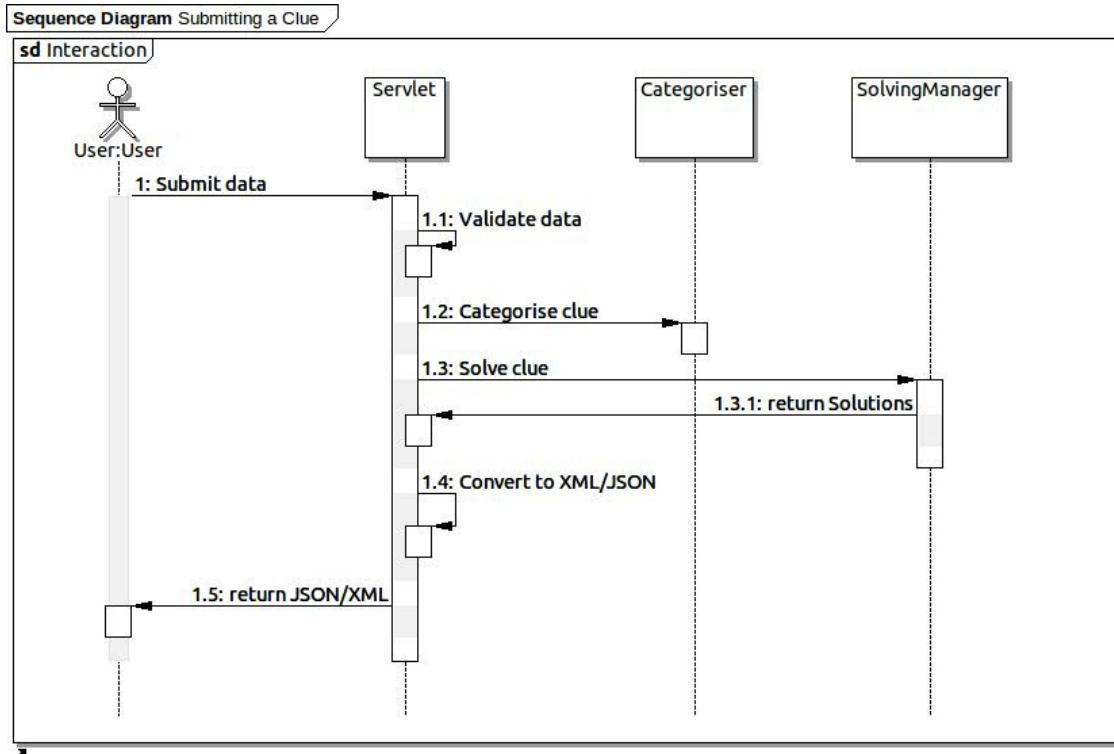


Figure 5.7: Sequence diagram illustrating a user submitting a clue

5.4.2 Solving a Clue

Figure 5.8 illustrates how a given clue will be solved. This sequence diagram assumes that valid data has been passed to the ‘SolvingManager’. This process was described in the previous subsection entitled ‘Submitting a Clue’ on page 86.

The ‘SolvingManager’ object will manage the process of solving a given clue. In order to do this it will have to house several processes, including distributing the clue out to one or more solvers, and merging all results.

Figure 5.8 illustrates the distributing and solving of the clue as a synchronous process, however this was used for illustrative purposes only. The system will in fact distribute and solve the clue upon an asynchronous basis. This will mean that each solver will be running at the same time, and thus reducing the total amount of time needed to solve the clue.

The system will wait however, for all solvers to finish, ensuring that all solutions can be returned to the user. The merging of the results will simply remove any duplicate solutions, upon a first come basis. So if a solution is already in the list, it will not be re-added because a different solver has managed to find the same solution.

Finally the confidence ratings will be adjusted. The adjustments are made upon a number

of factors, and were described in figure 5.4 on page 83.

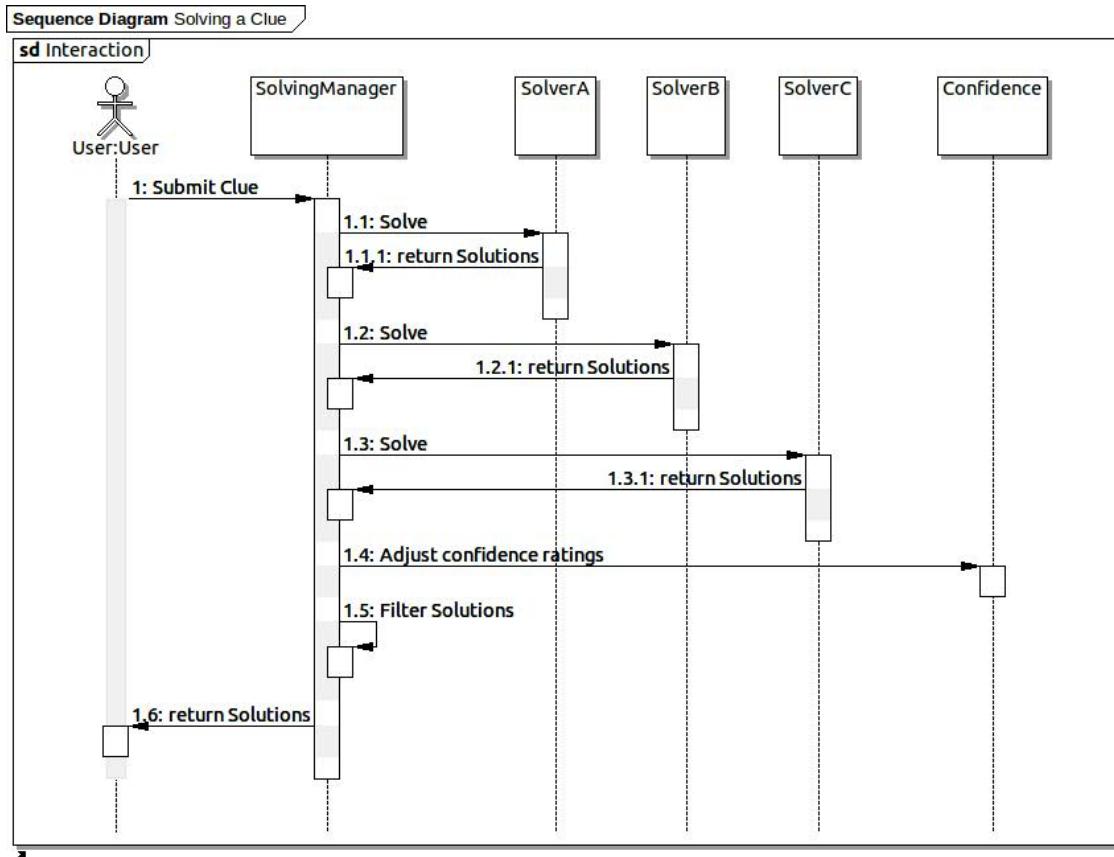


Figure 5.8: Sequence diagram illustrating the system solving a clue

5.5 Class Diagrams

Within this section a range of class diagrams will be presented and discussed. “A class diagram describes the structure of a system by highlighting the system’s classes, attributes, methods, and relationships with other classes” (Lunn, 2003).

Figure 5.9 illustrates the seven packages that make up the cryptic crossword solver. Over the following subsections, each of the packages will be described in more detail.

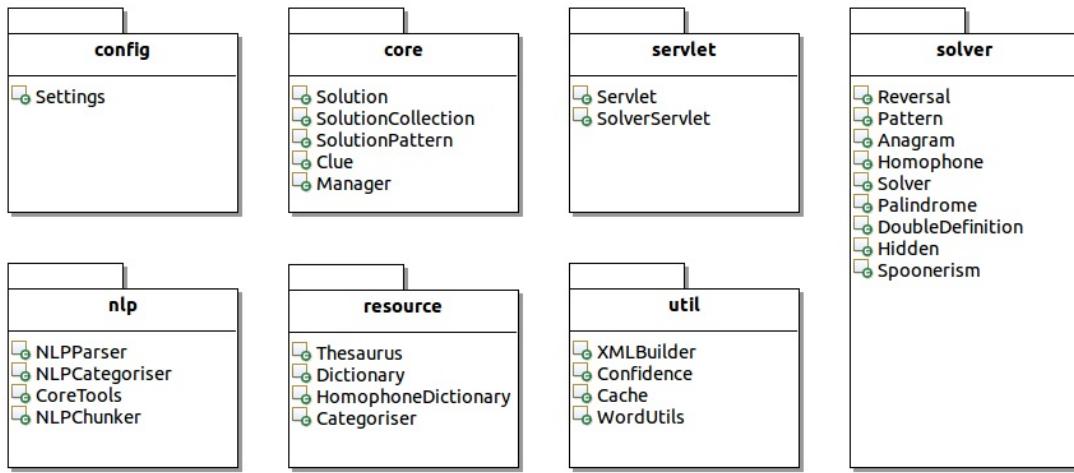


Figure 5.9: Package overview of the Cryptic Crossword Solver system

The architecture of the system will utilise a ‘plug and play’ approach. This approach, will allow for various numbers of solvers to be added at run time, rather than a compile time.

From a development point this will allow for the underlying system to be developed and tested utilising a limited number of solvers. Once the underlying system has been completed, solvers can then be added at any time, without having to alter any previously written code.

5.5.1 Config

The configuration class only contains one class, Settings. The Settings class is designed to hold all application specific settings, and each setting can be altered based upon the mode it is being run under — for example development, testing or production.

The Settings class follows the Singleton design pattern which ensures that the class will have only one instance (Gamma, 1995). This prevents various numbers of the same Settings objects being used in memory, and thus wasting valuable space.

Figure 5.10 illustrates the config package, containing the Settings class.

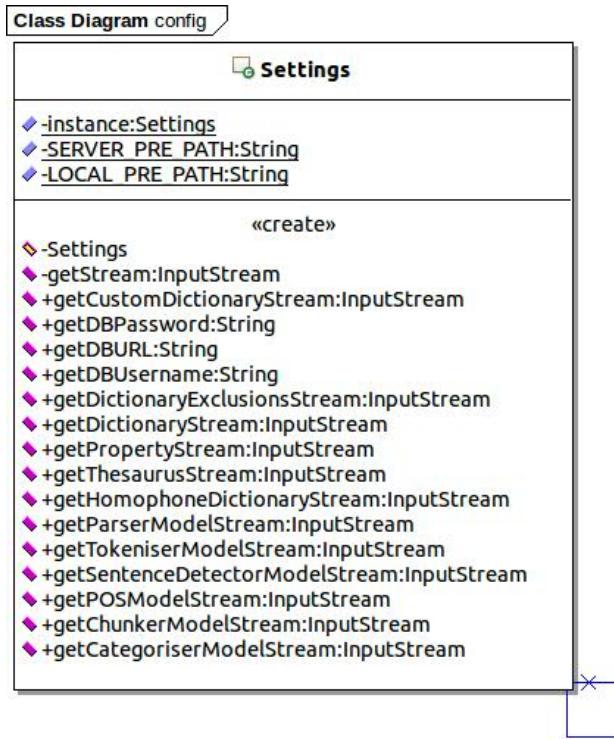


Figure 5.10: Config package class diagram

5.5.2 Core

The core package contains the various classes that are deemed to be at the heart of the system. A class has been devoted to the clue, a solution and a solution pattern providing various specific functionalities.

For example the solution pattern provides methods that allow for words to be computed, without having to work out the length of each word every time. The clue class is able to deduce the best solutions, based upon a collection of solutions.

The SolutionsCollection provides a simple interface to handle any number of solutions. The SolutionsCollection uses the HashSet as it's base class, which prevents duplicate solutions to be stored within the same set. It also provides a faster lookup, in comparison to List based structures.

Finally the manager class is the heart of the system and was originally illustrated in the subsection Solving a Clue on page 87. It will be able distribute a clue and it's solution pattern to the various solvers asynchronously. It also handles the merging of results within the `distributeAndSolve` method.

Figure 5.11 illustrates the core package, containing the various core classes.

Class Diagram core

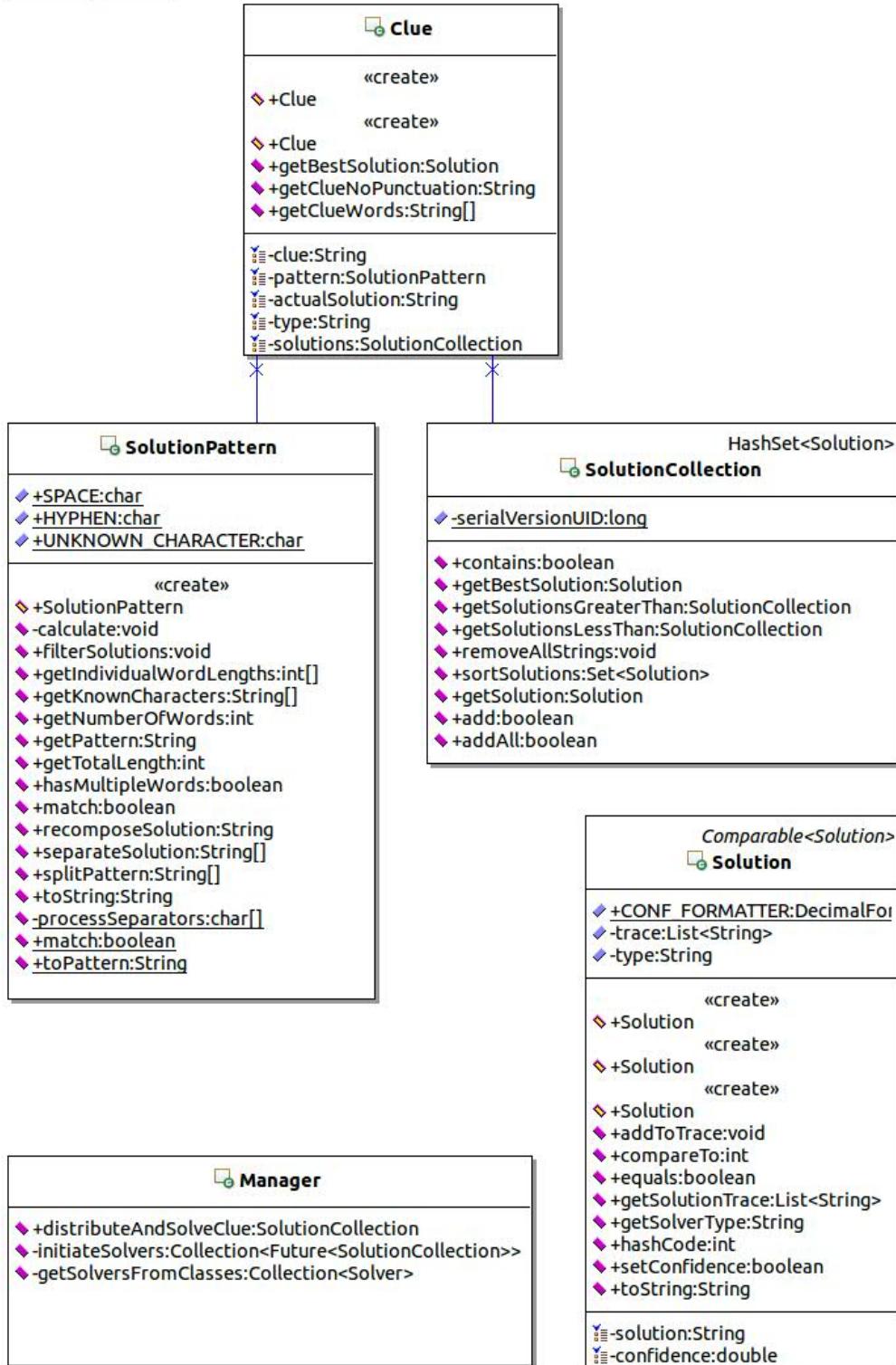


Figure 5.11: Core package class diagram

5.5.3 NLP

The NLP package contains the various classes that focus on providing the system with an interface to the Apache OpenNLP library. This will allow the system to be able to use natural language processing techniques from within the application.

Figure 5.12 illustrates the NLP package, containing the various NLP related classes.

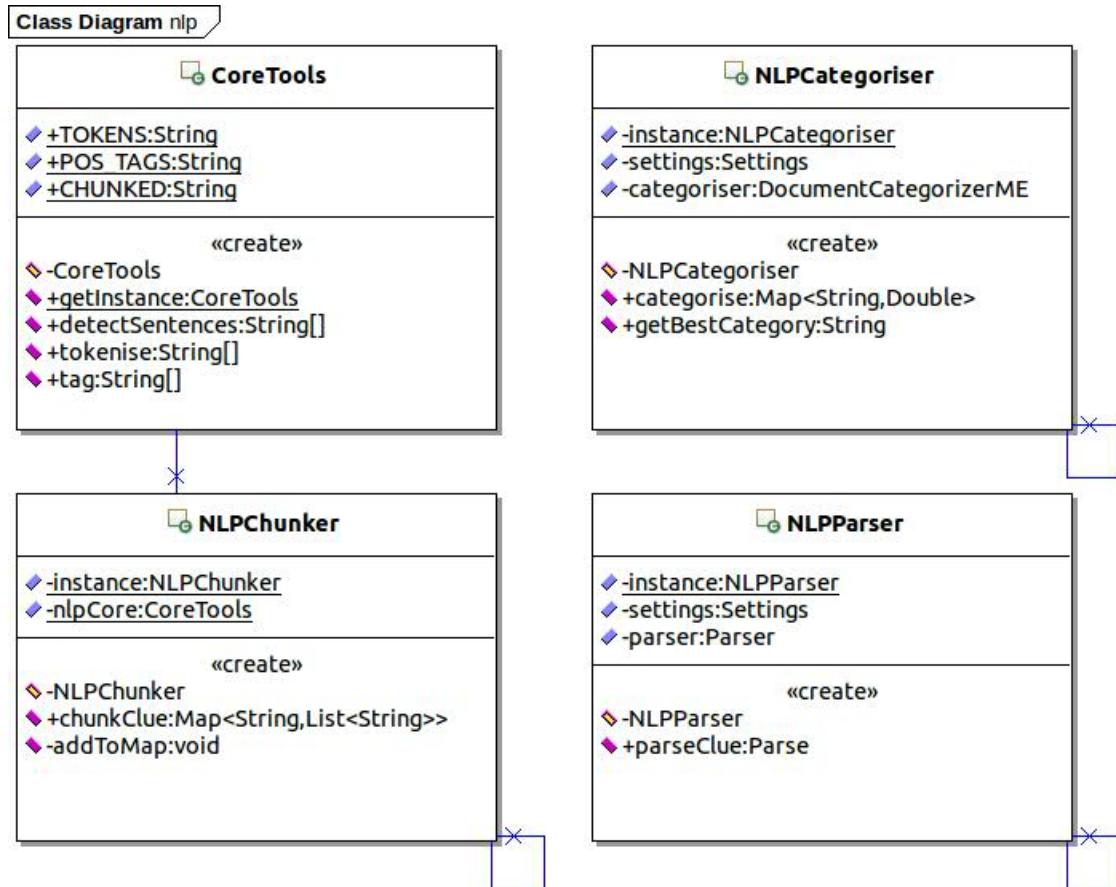


Figure 5.12: NLP package class diagram

5.5.4 Resource

The resource package contains various classes that provide interfaces to the array of resources that the system will use. As with the config package, all classes utilise the Singleton design pattern.

The reason for this is that it is to be expected that these classes will be used multiple times throughout the life of the system. To prevent multiple objects being created and being left

in memory, the singleton design pattern was selected. The pattern will “ensure that a class only has one instance, and provide a global point of access to it” (Gamma, 1995).

The Dictionary, Thesaurus and HomophoneDictionary classes provide the functionality that would be expected. For example the Dictionary has the ability to check to see if a given word is a ‘real word’. Whereas the Thesaurus class has the ability to find synonyms of a given word.

All of the classes are able to utilise the SolutionPattern class (as found in the core package), which allows for words to be matched upon various known and unknown character combinations.

The Categoriser class will map a given clue and solution combination to various lists of clue type indicators. If an indicator is found, then the confidence rating of the solution will be increased. This will form part of the overall ranking of a solution.

Figure 5.13 illustrates the resource package, containing the various resource classes.

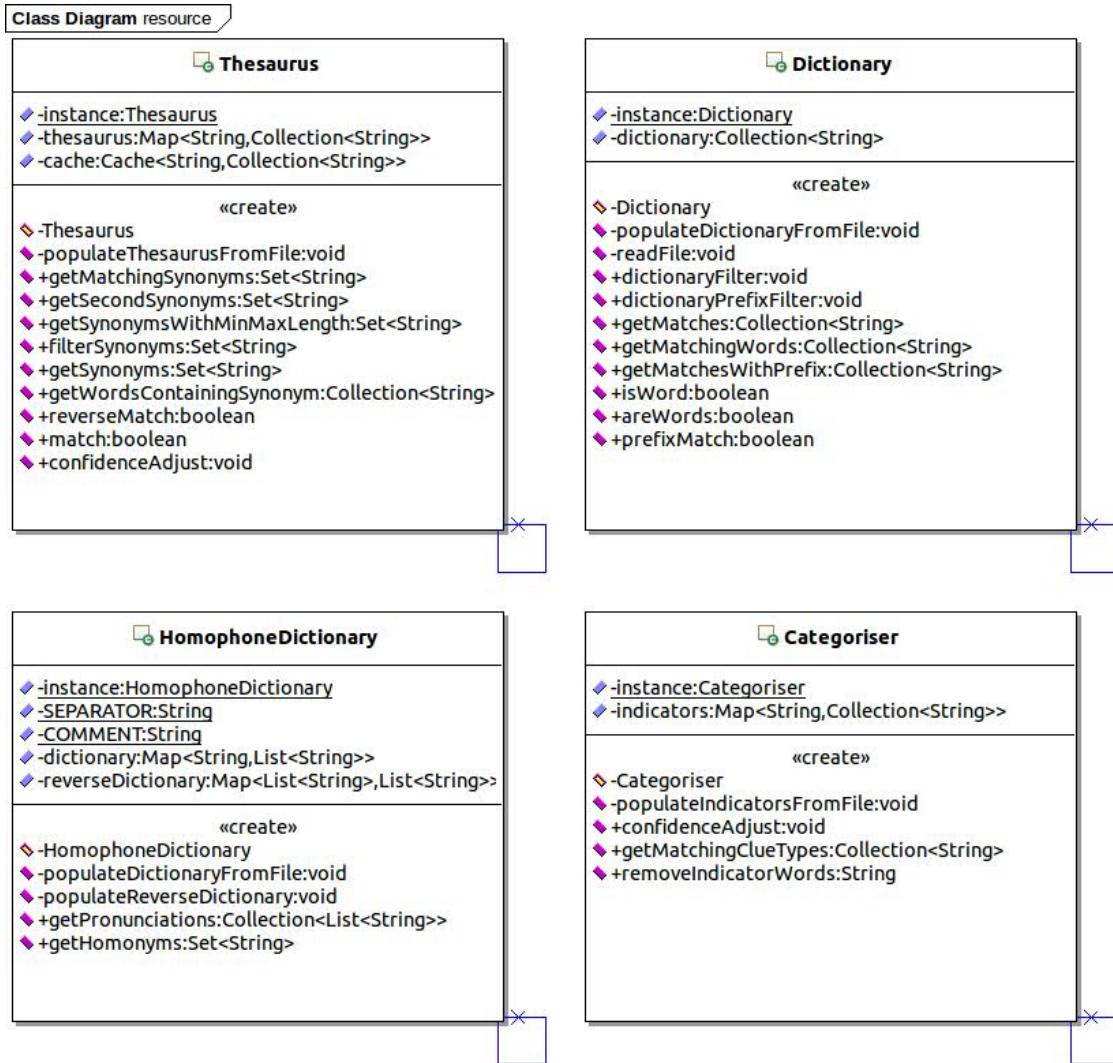


Figure 5.13: Resource package class diagram

5.5.5 Servlet

The servlet package contains two classes that provide the web service interface. Both classes tie into the Tomcat system, which provides a container for the web service to sit in.

The Servlet class is a customised base class that provides various levels of common functionality, that other specific classes will utilise. This functionality includes converting XML into JSON, and determining if a reject was made via AJAX.

The Solver class extends the Servlet class, and provides the functionality required for handling solver related queries. All GET and POST requests to the Solver class made via AJAX will be treated as a request over the web.

However all GET and POST requests that do not pass the AJAX flag, will be treated as if the user intended on viewing a website. All non-AJAX posts will be redirected to the Solver's default HTML page, which will allow users to enter the various parameters via a web form.

It is intended that the solver class will listen upon all requests upon the /solver URL path — for example www.example.com/solver.

Figure 5.14 illustrates the servlet package, containing the public web service related classes.

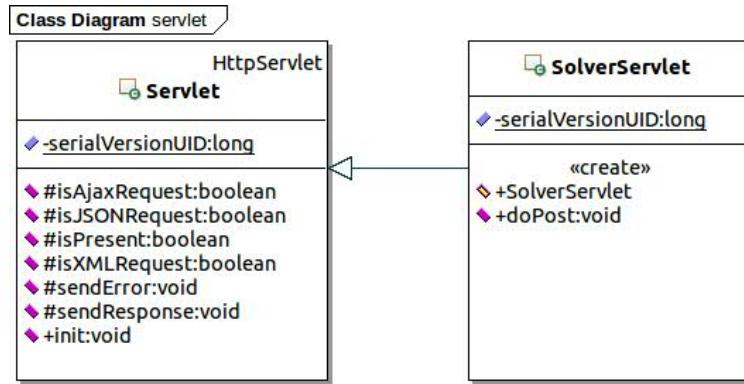


Figure 5.14: Servlet package class diagram

5.5.6 Solver

The solver package contains the all of the solver algorithm classes that are used as part of the over all solving process. Each solver extends a main Solver class which provides a number of base methods, as well as some additional methods that require each solver to override.

The Solver class implements the Callable interface, which provides functionality to run a class in it's own thread. A design decision was taken that the Callable interface was to be used, as once it has finished executing it is able to return an object unlike similar interfaces such as Runnable.

Each of solver classes implement the Abstract Factory design pattern via the base class Solver. This pattern “provides an interface for creating families of related objects without specifying their concrete class” (Gamma, 1995).

Each of the solvers is able to return a concrete product, which in this instance is a a SolutionCollection object. The SolutionCollection will contain all the solutions that the given solver has managed to compute.

Figure 5.15 illustrates the solver package, containing the various solving algorithms used by the system.

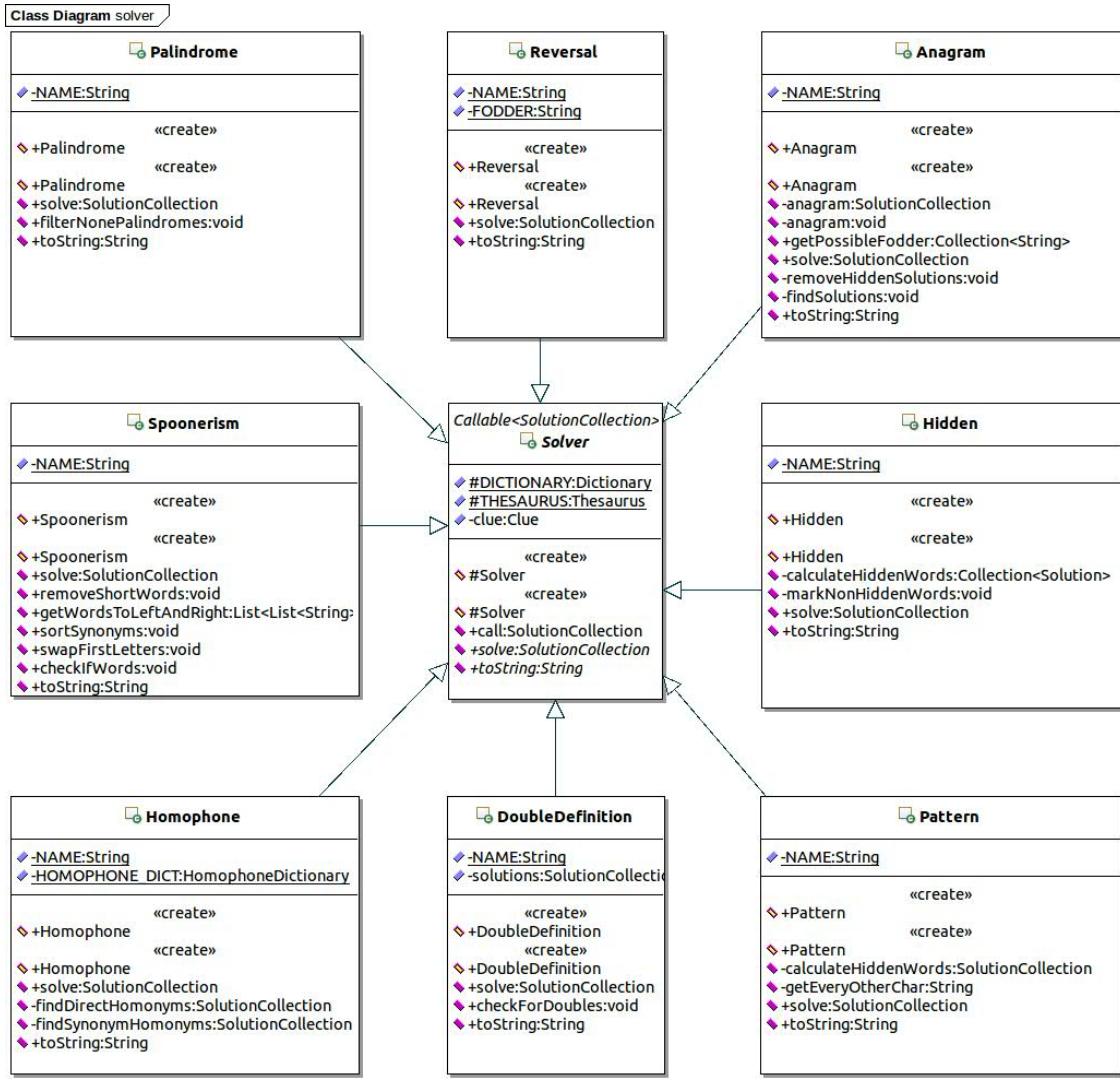


Figure 5.15: Solver package class diagram

5.5.7 Util

The util package contains a number of classes that are unrelated with each other, but provide functionality that is required throughout the system.

The XMLBuilder class provides a data encapsulation layer around an XML library. This class is able to build well-formed XML documents that are suitable for the system's output.

The WordUtils class provides a number of generic methods that are used to manipulate strings as words. This functionality is particularly useful, as the system will be dealing with large amounts of natural language (English) sentences.

The confidence class provides a number of rules, that when applied to a solution is able to deduce it's overall 'correctness'. This value is directly reported back to the user as the confidence rating.

Figure 5.16 illustrates the util package, containing the a number of additional utility classes.

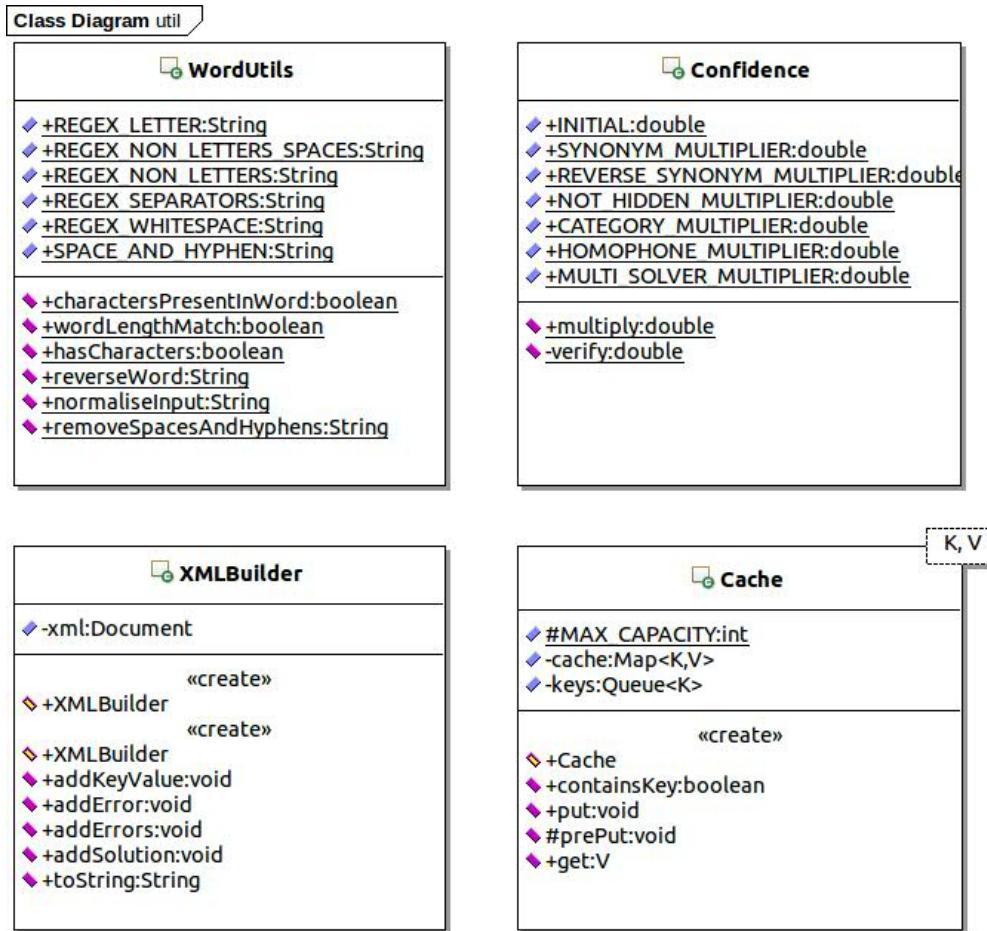


Figure 5.16: Util package class diagram

5.6 Database Design

As part of the testing strategy a list of clues and their solutions are required to be stored. In order to accurately store a large amount of training data, a database will be used.

The database design is of a minimalistic approach as it is not intended to hold a large amount of data, nor will it be used within a production environment. Figure 5.17 shows the proposed database design.

Each record will contain a clue, a solution, and the solutions length. The solutions length has been assigned a VARCHAR data type, because solutions can span across more than one word. For example 3,4 would relate to two words, of three and four characters respectfully. Each of these attributes are required — i.e. not null.

The orientation attribute is an ENUM that contains either ‘down’ or ‘up’, indicating the direction that the solution is arranged. The orientation may be useful for certain clue types, such as reversal, where by the clue references the direction — e.g. north-to-south, or east-to-west.

The clue number attribute stores the number of clue, as some clues use the clue number in obtaining the solution answer. Although this is uncommon, it will be taken into consideration.

The clue type attribute is an ENUM that holds the various types of clues, such as ‘purely cryptic’ or ‘anagram’. It is intended that these values will be ‘mapped’ to the system code, so that clues can be used automatically when training.

cryptic_clues	
• id	INT(10)
•clue	VARCHAR(255)
•solution	VARCHAR(32)
•length	VARCHAR(32)
◦orientation	ENUM
◦clue_number	TINYINT(3)
◦type	ENUM

Figure 5.17: Testing database entity-relationship diagram

5.7 User Interface

The previous sections have presented a number of system designs from a programmatic point of view. Within this section, there will be a focus upon user interface designs, and thus how the end user will interact with the system.

Fundamentally there are two aspects to the user interface design of the system — inputting the clue and retrieving the results. Each of these aspects will be discussed in more detail in the following subsections.

5.7.1 Platform Support

One of the main objectives of the project is to develop a system that can be used upon a number of mobile platforms, as well as trying not to neglect ‘traditional’ desktop users.

In order to complete both of these objectives, a responsive design is required. A responsive design is one that is able to dynamically change based upon the screen size.

This means that there will only be one code based for multiple screen sizes, and thus allows for easier maintainability. All of the designs within this section feature a responsive design. Figure 5.18 illustrates the power of responsive designs.

The left-hand side form shows a ‘traditional’ desktop experience, whilst the form upon the right-hand side, shows a mobile experience.

The figure displays two versions of an input form for a user interface, designed to demonstrate responsive design across different platforms.

Left Form (Desktop Experience):

- Clue:** Punctuation can also be included.
- Solution Length:** Any combination of single words (e.g. 3), multiple words (e.g. 3,5) or hyphenated words (e.g. 3-5) can be entered.
- Solution Pattern:** Provide any known characters, unknown characters (?), word separators (comma) and hyphens (-).

Right Form (Mobile Experience):

- Clue:** Punctuation can also be included.
- Solution Length:** Any combination of single words (e.g. 3), multiple words (e.g. 3,5) or hyphenated words (e.g. 3-5) can be entered.
- Solution Pattern:** Provide any known characters, unknown characters (?), word separators (comma) and hyphens (-).

Both forms include a **Reset** button and a **Solve** button at the bottom.

Figure 5.18: The input form to be completed by the end user

5.7.2 User Input

In order for the system to solve the clue, it must first be given the clue, along with additional supporting information. The additional information such as the solution length and solution pattern is vital to the system in order for it to compute the correct answer.

However the intended users of the system are likely to be operating upon some form of mobile device, meaning that a simple and power interface is required. It also means that space will be at a premium, and thus it can not afford to be wasted.

Figure 5.19 illustrates the design of the user input form. One of the most noticeable elements to the form design is how much space has been allocated to the input boxes.

The design of the inputs utilises a bi-column setup using the ratio of 1:3. This means that for every 1 pixel of space allocated to the left-hand labels, there will be 3 pixels of space allocated to the right-hand input boxes. This allows users of mobile devices to quickly select the input box, rather than having to tap a small area several times.

In order to assist the end user additional help blocks will be included, and give useful information to the user, such as describing the solution pattern format required.

Clue:
Punctuation can also be included.

Solution Length:
Any combination of single words (e.g. 3),
multiple words (e.g. 3,5) or hyphenated
words (e.g. 3-5) can be entered.

Solution Pattern:
Provide any known characters, unknown
characters (?), word separators (comma) and
hyphens (-).

Reset **Solve**

Figure 5.19: The input form to be completed by the end user

As the system is dealing with input from users, it is inevitable that a user will attempt to

input incorrect data. This will require validation to be performed, and if validation has failed then the user should be notified.

As previously mentioned space is a premium upon a mobile device. Therefore displaying a list of validation error messages at the top of the screen would not be utilising limited amount of space wisely.

In order to combat this issue, validation will be displayed ‘inline’ with the form input elements, as seen within figure 5.20. By utilising an ‘inline’ approach it recycles the screen space that has been made available to the form.

For inputs that have passed the validations the outline will change to green, and for validations that have failed the input group will change to red. This will immediately alert the user to the various issues, and thus allows the user to fix the errors in a quicker and more efficient manner.

The image shows a mobile application interface with a light gray background. At the top, there is a text input field containing the text "Found ermine, deer hides damaged". Below this, a green box contains the text "Clue: Found ermine, deer hides damaged". A note below the box says "Punctuation can also be included." To the right of the input field, the text "Solution Length:" is followed by a green box containing the number "10". A note below the box says "Any combination of single words (e.g. 3), multiple words (e.g. 3,5) or hyphenated words (e.g. 3-5) can be entered." Further down, the text "Solution Pattern:" is followed by a red box. A note below the box says "Provide any known characters, unknown characters (?), word separators (comma) and hyphens (-)." At the bottom of the screen are two buttons: "Reset" on the left and "Solve" on the right.

Figure 5.20: The input form indicating a validation error

5.7.3 Results

The displaying of the results follows on from some of the previous design decisions. Each potential solution is rendered into it’s own panel, and will contain the confidence rating, and the solver type that managed to deduce the solution (as shown in Figure 5.21).

Within an open panel, additional information is displayed — known as the solution trace. A solution trace provides a step by step account of how the solution was able to be computed. The solution trace may help to teach users how a particular type of clue is solved.

The results are ordered by confidence rating in ascending order, and by default the first panel will be ‘open’, showing the solution trace. The reason for this is that the top answer is likely to be the answer the end user is looking for

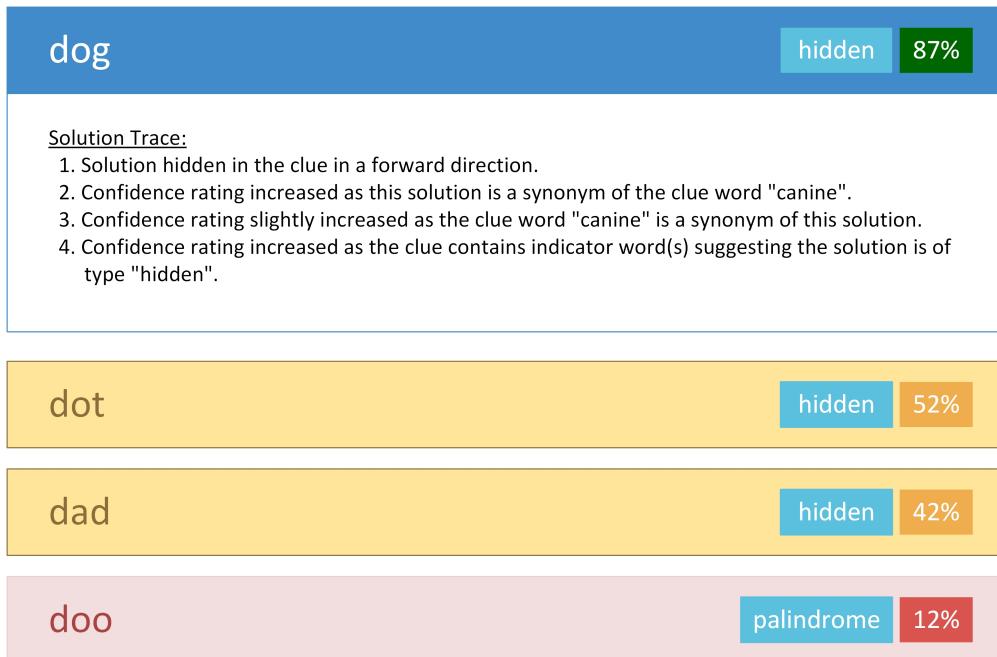


Figure 5.21: Results list displaying the top answer in blue

Each of the panels are “clickable”, meaning that if any given panel is selected then the solution trace will be displayed, whilst previously selected panels will be ‘closed’ (as shown in Figure 5.22). The main reason behind this is to reduce the amount of scrolling a user has to do, especially for those without a mouse.

Each solution is awarded a colour based upon the confidence rating. Blue refers to a top answer, whilst the remaining colours (green, yellow and red) indicate the likelihood of the solution being correct.

This will allow the end user to immediately deduce that ‘green’ results are more likely to be correct in comparison to ‘red results’, and that a ‘blue’ result is likely to be correct.

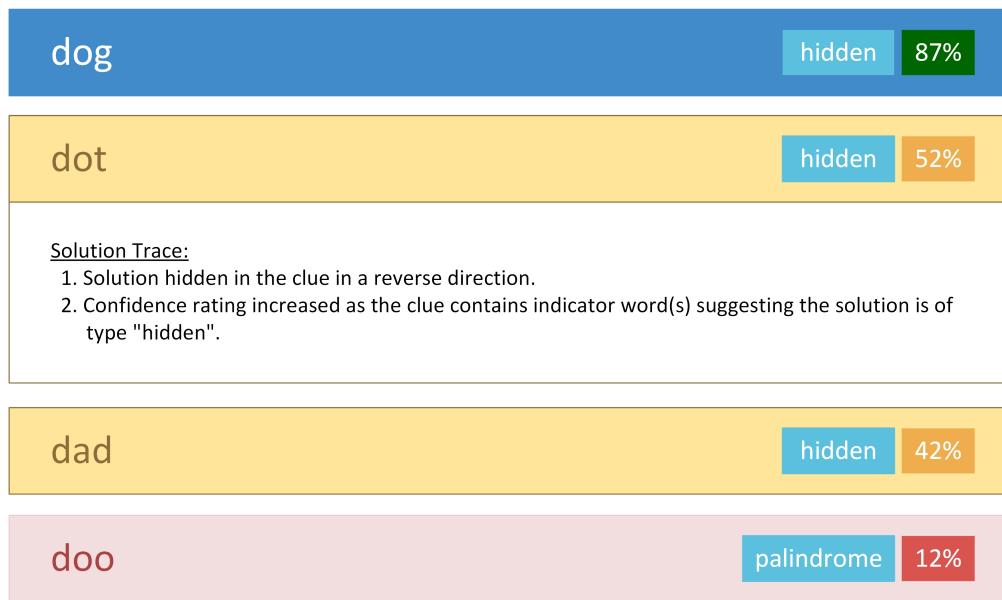


Figure 5.22: Results list displaying alternative solutions

Chapter 6

Development

6.1 Web Service & Servlets

The core system is wrapped around a RESTful web service, that allows users from various devices to submit clues to be solved. Within this section both the web service and the servlet implementations will be discussed and presented.

6.1.1 Web Service

During the project analysis phase the decision was taken that the system's functionality will be delivered via a web service. The web service was developed using Java Enterprise Edition (Java EE) and Tomcat 7.

The main reason for using Java is that the web service could easily make use of the various packages that are provided by the chosen natural language processing library — Apache OpenNLP written in Java.

The web service has solely been produced using the Java EE platform and does not use any additional frameworks or libraries such as Apache Axis. The reason being is that the Java EE platform will run on any machine that is capable of running the Java virtual machine without any additional configuration.

Although Apache Axis (for example) provides additional functionality in configuring the web service, it was decided that the project was to focus upon the solving of a clue. Therefore a 'standard web service' setup would easily meet the requirements of the project.

Another design decision was taken to ensure that the web service followed a RESTful style of communication. The mains reason for this is that some of the target devices (i.e. mobile and tablet platforms) do not support SOAP based communication without additional plug-ins. RESTful web services also provide a number of advantages over their SOAP-based counter parts, as was highlighted within the research section.

6.1.2 Servlets

The servlet design has been split across two classes — `Servlet` and `Solver`.

The `Servlet` class extends the standard `HttpServlet` class and provides common functionality. The `Servlet` class provides a base for all system Servlets to use the common functionality.

The `Servlet` class is able to if a given request is from a JSON, XML or Ajax background. For example if a client was to make a request to the web service through a web browser utilising an Ajax request, then the `isAjaxRequest` method would return `true`.

For illustrative proposes the `isAjaxRequest` method is shown below.

```

protected boolean isAjaxRequest(HttpServletRequest request) {
    String ajax = request.getHeader("x-requested-with");
    return ajax != null && ajax.toLowerCase().contains("xmlhttprequest");
}

```

The servlet also contains two customised methods – one to handle errors, and the other to handle a good response – that are able to send a response back to the requesting client based upon a number of factors.

For example the methods are able to convert the return data into either XML or JSON depending upon what the client has asked for. The `sendError` method will also set the HTTP status code correctly, allowing the client to correctly authenticate the response.

Finally the `Servlet` class overrides the `init` method, which is automatically called as part of the object construction. This method will initialise resources at servlet creation (i.e. first run-time within tomcat) rather than during the first call to the servlet.

In doing this, Tomcat will take more time initially starting up, however the user will notice that their queries are dealt with much quicker. In this project the `init` method has been used to initialise the various in-memory dictionaries and thesauri.

The second class is the `Solver` servlet and handles all request that are specifically for solving a given clue. The `Solver` servlet accepts both GET and POST requests, with each requiring the clue, the length of the solution and the solution pattern.

The `Solver` servlet upon receiving a request will validate the input parameters, based upon a number of criteria including presence checks and regular expressions. The code snippet below is an example validation rule, that will validate the solution pattern against a regular expression.

```

private boolean isPatternValid(String pattern) {
    // Pattern string regular expression
    final String regex = "[0-9A-Za-z?]+((,|-)[0-9A-Za-z?]+)*";
    boolean match = Pattern.matches(regex, pattern);

    // Pattern String must be present and of a valid format
    return isPresent(pattern) && match;
}

```

In order for validation to pass, the solution pattern must not be empty and must match to the regular expression stated in the method.

The `Solver` servlet class will initialise the solving of a clue if the three inputs are deemed to be valid. The `solveClue` method will utilise the Clue manager class, that will handle the

distributing of the clue to the various solvers.

This has been designed so that the servlet and the solving processes are upon separate threads. This prevents tomcat from freezing, and allows it to handle requests from users.

Once all the solvers have finished executing, the **Solver** servlet will produce an XML document based upon the various elements. Once the XML document has been created, will be sent back to the client as either XML or JSON.

6.2 User Interface

In order for users to use the system a simple and powerful user interface was required. The design reasoning's behind the user interface were described in the User Interface design section, which can be found on page 100.

The user interface has been designed utilising a fall-back system, which is a standard web development approach. The majority of the user interface is delivered by the server utilising JavaServer Pages (JSP) language.

The JSP language essentially extends from XML, and allows for html-like code to be written so that when complied with Java, a full server-side page is rendered. Within this project an additional JavaServer Pages Standard Tag Library was used. The library – xml – provided additional functionality so that JSP was able to directly utilise XML within it's rendering technique.

An example code snippet has been presented below from the `solver.jsp` file.

```
<x:choose>
    <x:when select="$solution/trace">
        <p>Solution Trace:</p>
        <ol>
            <x:forEach select="$solution/trace" var="trace">
                <li><x:out select="$trace"/></li>
            </x:forEach>
        </ol>
    </x:when>
    <x:otherwise>
        <p>Solution Trace Unavailable.</p>
    </x:otherwise>
</x:choose>
```

The code snippet above shows use of the XML library, as denoted by the ‘x’ name space to certain elements. The code is utilising XPATH to find all possible traces that are listed within the trace element (the trace element is an array).

For each of the traces found they will be printed out into the standard HTML output. However if a solution has not been found a simple predefined message will be presented.

The server side rendering that has been described above is known as the fall-back option. This option will work on all browsers and operating systems. The reason for this is that the rendering is controlled by the server, and hence can be fully tested.

Many modern browsers will support JavaScript in some form of fashion. This allows for various additional functionality to be provided to enhance the user's web browsing experience.

The cryptic crossword website features a JavaScript override, that will override the default server-side rendering and provide it's owner rendering. For example when a user clicks upon the 'solve' button, the web site will display a message informing the user that the clue is currently being solved as shown in figure 6.1. Under a non-JavaScript supporting browser this would simply look like a normal page request that is taking it's time.

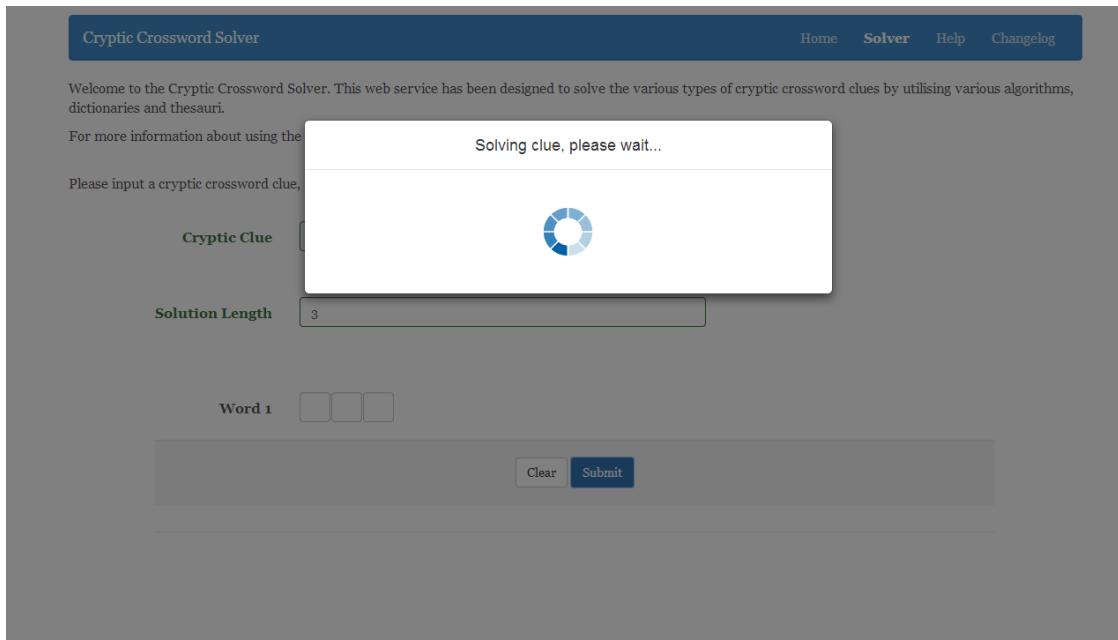


Figure 6.1: Message informing the user that a clue is currently being solved

The JavaSCript engine will also provide some basic forms of validation, to prevent the user from waiting a long time simply to find out that there was a validation error. The validation occurs upon a key press, and hence the feedback is instant as shown in figure 6.2.

 A screenshot of the same website as Figure 6.1, but with validation feedback. The 'Cryptic Clue' input field now has a red border and the placeholder text 'Introduction to do-gooder canine' is in red. The 'Solution Length' dropdown menu has a red border and the value '3' is in red. Below the dropdown, a red note says 'Any combination of single words (e.g. 3), multiple words (e.g. 3,5) or hyphenated words (e.g. 3-5) can be entered.' The rest of the page is identical to Figure 6.1.

Figure 6.2: Live validation feedback ensures users are entering correct data

The JavaScript engine will also make POST requests to the server. This reduces the total amount of work the server will be required to do, as it only has to return the requests (rather than render HTML).

The fall-back system that was previously described is that if for some reason the JavaScript engine fails to load, or is incompatible with the device, then the JavaScript will not load. However, the functionality of the site will still continue to work, as the browser will ‘fall-back’ to the standard server-side validation and rendering.

Chapter 7

Testing

Glossary of Terms

The following section contains a glossary with the meanings of all names, acronyms, and abbreviations used by the stakeholders.

Term/Acronym	Definition
The Guardian	A newspaper with a website featuring cryptic crosswords
Blackberry	A mobile phone platform by Blackberry
iOS	A mobile phone platform by Apple
Android	A mobile phone platform by Google
NLP	Natural Language Processing
SRS	Software Requirements Specification
App	Short for application

Bibliography

- Akman, O. (2013). The smartphone revolution. Retrieved from <http://mail2web.com/blog/2011/05/smartphone-revolution-growth-smartphones-exchange-activesync/>.
- Apache (2013). Apache opennlp developer documentation. Retrieved from <http://opennlp.apache.org/documentation/1.5.3/manual/opennlp.html#tools.sentdetect>.
- Apple (2007). Macworld san francisco 2007 keynote address. Retrieved from <https://itunes.apple.com/gb/podcast/apple-keynotes/id275834665>.
- Apple (2013). App store tops 40 billion downloads with almost half in 2012. Retrieved from <http://www.apple.com/pr/library/2013/01/07App-Store-Tops-40-Billion-Downloads-with-Almost-Half-in-2012.html>.
- Bird, S. (2009). *Natural Language Processing with Python*. O'Reilley Media, Inc.
- Businessline (2012). Microsoft to release windows 8 on oct 26. Retrieved from <http://search.proquest.com/docview/1321630361?accountid=11526>.
- Cadle, J., Eva, M., Hindle, K., Paul, D., Rollaston, C., Tudor, D., and Yeates, D. (2010). *Business Analysis*. British Computer Society, second edition.
- Cha, B. (2009). Rim store crowned blackberry app world. Retrieved from http://news.cnet.com/8301-17938_105-10188400-1.html?tag=mncol;txt.
- Collins, J. (2004). Languages, natural and formal. Retrieved from <http://www.uea.ac.uk/~j108/languages.htm>.
- Connor, A. (2012a). Cryptic crosswords for beginners: palindromes. Retrieved from <http://www.theguardian.com/crosswords/crossword-blog/2012/nov/01/cryptic-crosswords-beginners-palindromes>.
- Connor, A. (2012b). Cryptic crosswords for beginners: spoonerisms. Retrieved from <http://www.theguardian.com/crosswords/crossword-blog/2012/mar/01/cryptic-crosswords-beginners-spoonerisms>.

- Crossword Tools (2013). Clue solver. Retrieved from <http://www.crosswordtools.com/cm/>.
- Cryptic Solver (2013). Solver. Retrieved from <http://cryptic-solver.appspot.com/>.
- Dawson, C. W. (2009). *Projects in Computing and Information Systems: A Student's Guide*. Addison Wesley, second edition.
- Dictionaries, O. (2011). *Concise Oxford English Dictionary*. Oxford University Press, twelfth edition.
- DOSPINESCU, O. and PERCA, M. (2013). Web services in mobile applications. *Informatica Economica*, 17(2):17 – 26.
- Elliott, P. G. (2004). *Global Business Information Technology: An Integrated Systems Approach*. Addison Wesley, first edition.
- Friedlander, K. and Fine, P. (2009). Expertise in cryptic crossword performance. In *An exploratory survey*, Auckland, New Zealand. International Symposium on Performance Science.
- Gamma, E. (1995). *Design patterns: elements of reusable object-oriented software*. Addison-Wesley.
- Gershon, G. (2004). Web services. *e-Doc Magazine*.
- Gordius (2003). Cryptic crossword no 22,717. Retrieved from <http://www.theguardian.com/crosswords/cryptic/22717>.
- Halpern, J. (2013). Cryptic crossword no 26,091. Retrieved from <http://www.theguardian.com/crosswords/cryptic/26091>.
- IDC (2013). Android and ios combine for 91.1% of the worldwide smartphone os market in 4q12 and 87.6% for the year, according to idc. Retrieved from <http://www.idc.com/getdoc.jsp?containerId=prUS23946013>.
- ISTQB Exam Certification (2013). What is rad model- advantages, disadvantages and when to use it? Retrieved from <http://istqbexamcertification.com/what-is-rad-model-advantages-disadvantages-and-when-to-use-it/>.
- Jiang, T., Li, M., Ravikumar, B., and Regan, K. W. (2010). *Algorithms and theory of computation handbook*. Chapman & Hall.
- Jones, C. (2013). Idc: Windows phone is the second most popular mobile platform in latin america. Retrieved from <http://www.forbes.com/sites/chuckjones/2013/09/30/iphones-market-share-down-prior-to-5s-5c-launch-with-windows-almost-double-digits-in-europe/>.

- Kalin, M. (2013). *Java Web Services: Up and Running*. O'Reilly Media, second edition.
- Knott, R. and Dawson, R. (1999). *Software Project Management*. Group D Publications Ltd.
- Litman, D. (2003). Cs 1573: Artificial intelligence application development (spring 2003). Retrieved from <http://people.cs.pitt.edu/~litman/courses/cs1573s03/lec/chunking.ppt>.
- Lunn, K. (2003). *Software development with UML*. Palgrave, first edition.
- McCluskey, T. L. (1999). Cha 2555 - artificial intelligence. Retrieved from http://helios.hud.ac.uk/scomtlm/cha2555/lecture13_18.ppt.
- McManus, J. (2003). *Risk Management in Software Development Projects (Computer Weekly Professional)*. Routledge.
- NLTK Project (2012). Nltk 2.0 documentation. Retrieved from <http://nltk.org/data.html>.
- One Across (2013). Crossword puzzle help. Retrieved from <http://www.oneacross.com/>.
- Robertson, J. and Robertson, S. (2013). *Mastering the requirements process*. Addison-Wesley, third edition.
- Sabri, S. (2013). Idc: Windows phone is the second most popular mobile platform in latin america. Retrieved from <http://www.wpcentral.com/windows-phone-second-most-popular-mobile-platform-latin-america>.
- Stoeller, W. (2003). What is risk management? *Globalization Insider*, XII(3.7).
- Sullivan, T. (2001). Web services. *InfoWorld*, 23(11):38.
- Upadhyay, S. (2008a). Acrostics. Retrieved from <http://www.crosswordunclued.com/2008/09/first-and-last-letters.html>.
- Upadhyay, S. (2008b). Charades. Retrieved from <http://www.crosswordunclued.com/2008/11/charades.html>.
- Upadhyay, S. (2008c). Cryptic definitions. Retrieved from <http://www.crosswordunclued.com/2008/12/cryptic-definitions.html>.
- Upadhyay, S. (2008d). Decoding double definitions. Retrieved from <http://www.crosswordunclued.com/2008/10/decoding-double-definitions.html>.
- Upadhyay, S. (2008e). Digging out hidden words. Retrieved from <http://www.crosswordunclued.com/2008/08/digging-out-hidden-words.html>.

- Upadhyay, S. (2008f). How to spot anagrams. Retrieved from <http://www.crosswordunclued.com/2008/08/how-to-spot-anagram.html>.
- Upadhyay, S. (2008g). &lit, literally so. Retrieved from <http://www.crosswordunclued.com/2008/08/and-literally-so.html>.
- Upadhyay, S. (2008h). Reversals. Retrieved from <http://www.crosswordunclued.com/2008/11/reversals.html>.
- Upadhyay, S. (2008i). Substitutions simplified. Retrieved from <http://www.crosswordunclued.com/2008/12/substitutions.html>.
- Upadhyay, S. (2008j). Tune in to homophones. Retrieved from <http://www.crosswordunclued.com/2008/10/homophones.html>.
- Upadhyay, S. (2009a). Containers. Retrieved from <http://www.crosswordunclued.com/2009/02/containers.html>.
- Upadhyay, S. (2009b). Deletions. Retrieved from <http://www.crosswordunclued.com/2009/03/deletions.html>.
- Upadhyay, S. (2009c). Letter exchange. Retrieved from <http://www.crosswordunclued.com/2009/12/letter-exchange.html>.
- Upadhyay, S. (2009d). Letter picking. Retrieved from <http://www.crosswordunclued.com/2009/04/letter-sequences.html>.
- Upadhyay, S. (2009e). Letter shifting. Retrieved from <http://www.crosswordunclued.com/2009/12/letter-shifting.html>.