

# Solving Cryptic Crosswords through Functional Programming

Michael Skelly  
Department of Computer Science  
Imperial College

April 25<sup>th</sup> 2014

## Part I Literature Review

### 1 Summary of Cryptic Crosswords

#### 1.1 Cryptic Crosswords

To start, let us provide some basic definitions around crosswords and their taxonomies.

A crossword is a puzzle, usually published in newspapers or magazines. They consist of a grid of squares, often 15 x 15. Some of the squares are white (i.e. blank) and some are blacked out. Any contiguous run of more than one white square, either down (vertically) or across (horizontally, left to right) is a space for a word, to be written. These are marked by numbers in the initial square (the top-leftmost one), and referred to by those numbers, and the direction (e.g. '5 down', '8 across'). Horizontal runs can overlap vertical runs, and at the points at which they do, each of the two words, when written in, must have the same letter in that square. Along with the grid are a set of clues, which the solver can use to determine which word to write in each space (the 'answer' or 'solution'). The aim of the puzzle is to find the set of solution words such that each clue's solution is correct for that clue, and fits in the grid correctly, with respect to the overlapping words.

Grids can be very densely white, with few black squares and most squares shared by two words (usually called AMERICAN STYLE) or more sparse, with fewer overlapped clues (called BRITISH STYLE). Clues can also be in two styles. STRAIGHT or QUICK crossword clues usually provide a single straightforward indicator as to what the correct word might be - often a synonym for the clue ('Joyful' = 'Happy') or a missing word ("Stitch in \_\_\_\_ saves nine" = "Time"). CRYPTIC clues are less straight-forward, appearing on the surface to be a valid

syntactic utterance in English, but actually consisting of a definition (as in the Straight clue) and some wordplay which the solver can use to arrive at the same answer as with the definition by apply a series of transformations and operations. The challenge is that the definition and the wordplay are not clearly separated, and that there are multiple ways to apply to the transformations, but with only one yielding the correct answer.

It is the task of determining the correct answer for this type of clue that this paper will address.

## 1.2 Cryptic Crosswords in the Literature

While not a topic well covered in scientific literature in general, what few analytical studies around cryptic crosswords there are tend to be classifiable into three main groups

### 1.2.1 Generation of Cryptic Clues

The largest body of work that exists is centered around the generation of cryptic clues, focused largely around analysis of how string literals from a pre-determined answer can be transformed by set clueing patterns, as well as some work around measures of the quality of generated clues.

### 1.2.2 Interpreting Clues

The next set are the select few who have done prior, similar investigations into interpreting cryptic clues, with some work put into formalizing definitions and notation for the sorts of clue types that appear in the majority of cryptic crosswords, and some attempts at solving based on these interpretations.

### 1.2.3 Other Work

There has also been some work done towards solving non-cryptic crosswords probabilistically, working on whole-grid solutions rather than individual clues. There are also some more left-of-field studies done: statistical studies into errors made during manual solving, and psychological studies into solving.

## 1.3 Complexity

A variety of factors make solving cryptic crosswords a difficult problem:

**Ambiguity** Cryptic crosswords are deliberately ambiguous. Instruction indicators are indistinguishable from string literals, which are identical to words' semantic meanings. Often, the setter will deliberately chose words to give rise to further ambiguities. For example, the *Telegraph* printed

Bug starts to move in dark, glowing endlessly (5)

cluing for 'MIDGE'. Usually "endlessly" and similar mean "remove the last letter", but here it is one of five consecutive words to form an acronym from, with the word "starts" as an indicator.

**State Space** Even with only a few different clue types, the number of different readings of one clue based on those grows exponentially with the length of the clue. This means that unless heuristics are applied, the evaluation time for a whole grid longer clues may be unfeasibly long. Even longer than it takes me to do the Times Crossword.

**Lack of Standardization** Although all cryptic crossword share some common conventions, there are no fixed rules shared between publications for what can and can't be a clue, indicator etc.. Although most publications have internal guidelines or style-guides, these are not accessible to the solver, and some publications (such as the *Guardian*) have named setters whose styles and self-imposed rulesets differ, even between one publication. Alistair Ferguson Ritchie, who set for *Listener* for many years, referenced the concept of fairness in his book *Armchair Crosswords* in 1946. He defers the judgement of fairness to a notional rulebook:

We must expect the composer to play tricks, but we shall insist that he play fair. *The Book of the Crossword* lays this injunction upon him: "You need not mean what you say, but you must say what you mean." This is a superior way of saying that he can't have it both ways. He may attempt to mislead by employing a form of words which can be taken in more than one way, and it is your fault if you take it the wrong way, but it is his fault if you can't logically take it the right way.

Although *The Book of the Crossword* there have been many books written on the subject of what should and should not constitute a valid cryptic crossword clue. One of the most notable and influential was written by *Observer* setter Derrick Somerset Macnutt, both cluing and writing under the name Ximenes, in his book *Ximenes on the Art of the Crossword Puzzle*. The book contains many in-depth guidelines about what a fair clue entails, summed up by his successor Azed (Jonathan Crowther, born 1942):

A good cryptic clue contains three elements:

1. a precise definition
2. a fair subsidiary indication
3. nothing else

A crossword setter following these rules is said to adhere to 'Ximenean principles' and their produced work to be Ximenean. Most mainstream crosswords exist on a continuum between being more closely Ximenean (examples include *The Times*, the *Independent*) to being very libertarian (e.g. *Guardian*). No crossword in a major UK newspaper is 'strictly Ximenean'.

**Knowledge Base** As well as being made up of encrypted and hidden meanings, cryptic crosswords also draws on a diverse knowledge base of synonyms, abbreviations, facts etc. These can include information as diverse as names of capital cities, common sayings, and the fact that one may carry a wallet in ones pocket.

In order to run a fully working cryptic crossword solver against any arbitrary clue, all of these pieces of information must be encoded, stored and accessible to the solver in a machine readable form. Understandably, this is subject to an entire field of study itself.

## 1.4 Programming Language Analogues

Much of the current work on interpreting crosswords draws on work by Backus, Naur and Chomsky in creating a specification for the grammar of crosswords. While these frameworks are useful for describing many different languages, interpretations of the grammar of cryptic crosswords seem to be perversely somewhat closer to mathematical and programming languages than to natural language. In some ways, the cryptic clue as a whole can be thought of as a program that generates the output string as its answer. The wordplay section is analogous to a program, and the definition section of a clue could be thought of as a checksum to verify the final answer.

### 1.4.1 Lexing, Parsing, Evaluating

The steps for compiling and running a computer program apply also to solving (or 'running') a crossword as a program. Each word in the input string needs to be tokenized, parsed into a relevant structure and then that structure evaluated to produce the final answer. Unusually for a programming language, however, the grammar of a cryptic crossword is highly ambiguous, and requires complex parsing. Firstly, programming language are only usually required to output the one valid abstract syntax tree, however here we may need to output many thousands in order to evaluate them to see which yields the correct answer. Furthermore, the grammar cannot be expressed without using complex context-sensitive features such as lookbacks, lookaheads and backtracking. Most major programming languages are parsed without these features, allowing information to flow in one direction from the lexer to the parser. To parse a cryptic crossword, lexing and parsing need to take place simultaneously in a process referred to as called "Scannerless Parsing".

## 1.5 There's No Accounting for Wit

Along with clearly defined and program-like cryptic crossword clues, there exist other clues that rely on humour, imagery and wit, rather than following the regimented classical structure, as set out by Ximenes. Some examples include:

Flower of London? (6)

```
(= THAMES, flower = that which flows)
In which you can get three couples together and have sex (5)
(= LATIN, 'sex' is 6 in Latin)
```

Clues such as these, and the question of computerised wit and humour, unfortunately exist out of the scope of this project.

## 2 Parsing Frameworks and Notation

Some different notations for denoting parsing of cryptic clues have come out of previous work – in order to properly provide a rigorous analysis of the structures and conventions of cryptic crosswords, it is necessary to analyze and choose a framework in which to do it.

### 2.1 LACROSS

William and Woodhead produced language called LACROSS, which forms a sort of calculus for describing crossword clues. They also provide a BNF definition of this grammar. Their clues are of the general form

```
Clue := Δ = G | G = Δ
```

the orientation of which corresponds to the order in which we find the definition ( $\Delta$ ) and the wordplay ( $G$ ) in the clue. The wordplay may be further expanded out – the wordplay section of the clue is expressed as a sequential annotation for the constituent parts, either as ‘text’ ( $t$ ), ‘shortening’ ( $S$ ) (etc.) or as placeholders for the operators ( $*$ ), which are detailed afterwards, including a reference to the substituted indicator. So for instance:

```
Get in odd bit of colour (5) [= tinge]
t* = Δ, a (odd, a)
```

There are several issues with this grammar. Firstly, all unitary operators are treated the same, as are all binary operators, and there is some issue with binding and precedence which they address with an underlining notation, in addition to brackets. Secondly, the grammar attempts to include both the structure of the parsed cluing and how that structure relates to the original sequence of words at the same time. As a result, we end up with complex grammar that does not aid human parsing of the solution well, nor does it lend itself easily to computer or mathematical manipulation

Still, they have provided the basis for future work, and begun a basic enumeration of clue types.

### 2.2 Simple Clue Markup Language

Proposed by Hall and Rapanotti, Simple Clue Markup Language (SCML) attempts to notate the structure of the solution directly onto the clue.

Double underlining is used to denote the definition, underlining denotes an operator, with its class as an optional subscript, with scope provided by brackets and concatenation (and definition/wordplay separation) given by a semi-colon. Thus in their given example:

**Note** the shuddering appliance Bill regularly installed, noisy thing (6,7)  
Note;(the)shuddering<sub>a</sub>;(appliance,(Bill)regularly<sub>t</sub>)installed<sub>e</sub>;noisy thing

Note' often indicates a musical note, resolving to one of 'a' to 'g',  
 'do', 're', 'mi', etc;  
 'the shuddering' may be an anagram indicator applied to 'the';  
 the 'regularly' of 'Bill regularly' may indicate alternate letters ('t');  
 i.e., 'bl' or 'il'; and  
 'installed' suggests the embedding ('e') of those letters within something meaning 'appliance'.

In this, we have no markup differentiation for literal strings ('Bill') against words with their semantic context ('appliance'), and we also take certain words that reduce to abbreviations ('Note') to be non-deterministic nullary operators. With some changes and additions (tagging of string vs. semantic word, for example), this markup serves as a good way to represent a parsing of a clue in a human readable way. It even has the advantage that a printed clue could be annotated (carefully) by hand, as a teaching aid, for example. Unfortunately, the language as it stands is not expressible as a BNF grammar, nor is it a particularly good format for representing the clue and its parsings internally in a program (as it would need to be re-parsed to use!)

## 2.3 Clue-answer notation

There are several emergent solutions within online cryptic crossword communities for notation to explain solutions derived from clues. From <http://cryptics.wikia.com>:

Consider the down clue A message from the setter, hauled up with broken arm after heroin withdrawal (8) yielding the answer TELEGRAM. The corresponding wordplay, having the prolix and possibly ambiguous explanation THE next to LEG reversed next to an anagram of ARM, all with H (heroin) removed could be concisely represented in clue-answer notation simply as T[h]E,GEL<=,(ARM)\*.

These meanings are not fixed, but some definitions are given here:

**ABC<= or ABC (rev.) ABC reversed.** The (rev.) notation is most commonly used when the wordplay consists of a single reversal.

**[abc] or -abc or (abc)** Letters abc removed, as in[c]OUNT to represent 'count' with c removed; the convention is to use lower case for the removed letters.

**(ABC)** Letters placed inside others, as inC(AND)ID to mean 'and' inside 'cid'.

**"ABC"** Homophone of ABC.

**(ABC)\*** Anagram of ABC.

**A+B or A,B** A concatenated with B. Sometimes both notations are used together where ambiguities may arise.

**aBcDeF** Alternate letters of ABCDEF (shorthand for [ a]B[c]D[e]F).

## 2.4 PICCUP

Hart and Davies define what is currently the most satisfying proposal for a formal syntactical definition of cryptic crossword syntax, in a loosely BNF grammar. There is the only current definition that closely resembles a usable formally defined language.

Their interpretation only specifies the grammar in terms of building an abstract syntax tree, rather than attempting to include a notation for clue or answer.

```
Anagram → Synonym(.Equ Indicator).AnagramSentence
/AnagramSentence(.Equ Indicator).Synonym
AnagramSentence → AnagramPointer.AnagramMaterial
/ Anagram Material.Anagram Pointer AnagramPointer~ Word(.Word)*
AnagramMaterial → Word(.Word)*
Synonym → Word(.Word)*
Equ Indicator → Word (.Word)*
```

## 2.5 Syntactic and Metasyntactic Conventions

Here we apply a similar convention to Hart, in using a modified Backus Naur Form (BNF). We will later see that a context-free grammar may not be sufficient to model a cryptic crossword, and may have further deficiencies as a basis for finding a solution. Nevertheless, we will adopt a similar notation:

```
→ = is composed of
, = followed by
| = or
(x) = x is optional
x* = 1 or more occurrences of x
```

$(x)^* = 0$  or more occurrences of  $x$

We also take the BNF conventions

Word = non-terminal symbol  
 “word” = string literal  
 $[x, y, z]$  = list containing  $x$   $y$  and  $z$   
 $(x, y)$  = pair  $x$  and  $y$

For clarity, we additionally define:

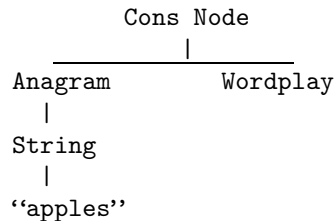
String = [any string literal]

## 2.6 Context Free?

The grammar described in this paper not a regular grammar (for example: any of the binary operators generate two non-terminals), but it can be formulated as a context-free grammar.

We define a CONTEXT-FREE GRAMMAR (CFG) as one in which the expansion of a non-terminal is not affected by the symbols before and after it.

While we can certainly define a working grammar for cryptic crossword clues in terms of a CFG, it may be useful to consider other options as a means of reducing the number of trees generated during the parsing phase to speed up the evaluation phase. We could take, for example, the clue length as a contextual variable: in that case, a 6 letter clue whose parse tree contains an anagram of a 4-letter string cannot yield another anagram of 5 letters.



In this example, the wordplay on the right should not be expanded out to an anagram node featuring a string of 5 characters (to consume, say, the string “mixed pears”).

## 2.7 Syntax vs Semantics

Due to the ambiguous and duplicitous nature of the structure of cryptic crosswords, especially the deliberate challenges in the lexing phase, unclarity between the boundaries between parsing the syntax and evaluating the semantics emerge.

Strings consisting of one or more words can be at once tokens representing different operators, they can be strings, and can be split in multiple ways into combinations. This is especially true when we have token that, in the original



text, represent their semantic meaning in English, and evaluate out to a finite number of equivalent words (roughly, synonyms: see later for discussion about this equivalence relation).

Hall and Rapanotti treated these roughly as their own operators: so the string “rough” would parse to the token `Rough`, which later evaluates to a finite number of definitions. This may, indeed, be tempting if we had a limited number of candidate words. And, indeed, we do need to differentiate these from raw string literals that are subject to Hidden Word or Initial Letter operators, or occasionally concatenated in their raw form (for example the string “it” is sometimes taken as given where necessary) .

I think a more manageable way and satisfying way to consider these options is to consider them subject to an invisible ‘word’ operator. This keeps the semantics and syntax more separate, but certainly poses some challenges for a parser / lexer.

## 3 The Cryptic Crossword Clue

### 3.1 Structure of a cryptic clue

A cryptic crossword differs from a normal crossword in that the clue for each answer consists of two parts. The first is the definition, which performs the same function as a clue in a ‘regular’ crossword. The answer to the clue is usually a synonym for the definition (‘circular’ and ‘round’) or may be an example of the definition (‘farm animal’ and ‘pig’). Other forms that the definition may take will be discussed later on. The second part of the clue is the wordplay. This is an encoded and often ambiguous second method of deriving the answer, using techniques such as anagram, substitution and concatenation. The clue as a whole is presented as a concatenation of the two parts, sometimes with a subsidiary word indicating that one can be derived from the other (for example, ‘from’ or ‘is’). We can present this breakdown as:

```
Clue → Definition, (Indicator), Wordplay
      | Wordplay, (Indicator), Definition
```

The final clue will often resemble a valid English utterance, although this ‘surface reading’ (i.e. {clue} ) very rarely has any relation to the answer. Later on we will consider other information and context within the definition of a clue.

### 3.2 Definition

The definition of the clue consists of one or more English words. The answer to the clue will be a word or phrase that fits an appropriate equivalence function (that we will define later).

The definition carries a variety of linguistic features with it that the overall answer, and so the answer as derived by the wordplay, must match. These include aspect (noun, verb, adjective), plurality (tree, trees), tense (go, going, gone). These features may also be considered as ‘context’ to the clue itself.

### 3.2.1 Formally

We can define the definition as

Definition  $\rightarrow$  Words

## 3.3 Wordplay

The wordplay section of a clue is a set of deliberately ambiguous instructions that allows the solver to arrive at the eventual answer. As the instructions are ambiguous, multiple possible parsings of the instructions are possible. Some of these parsing will not lead to a valid English word:

```
Imbecile, bonkers, in a cult (7)
==> Wordplay 'Imbecile, bonkers = definition 'in a cult'
==> Anagram 'imbecile' [indicator = bonkers] = definition 'in a cult'
==> ??? (no anagrams of imbecile in english language)
(correct reading was anagram of in a cult = lunatic)
```

Others will lead to a valid English word, but one that is not equivalent to the definition:

```
Minder shredded corset (6)
==> Wordplay 'minder shredded' = definition 'corset'
==> Anagram 'minder' [indicator = shredded] = definition 'corset'
==> 'remind' = definition 'corset?' X
(correct reading was anagram 'corset' = escort = minder)
```

The solver must find the correct parsing of the wordplay that yields the correct definition: even though they may not know which part is wordplay and which is definition.

## 3.4 Special Operators

I include these two operators first, as they really form the backbone or basis of other clues. They are also unique in being implicitly clued, rather than requiring an indicator word to signify their presence.

**Word Equivalence** In the most simple of clues, we have the definition, along with a word or phrase that is somehow semantically equivalent to that definition.

<sup>1</sup> A clue that contains just this structure is said to be 'double definition'

Metal guide (4) [= LEAD]

---

<sup>1</sup>In this case, it becomes a difficult task to be precise about exactly which of these is the definition and which is the wordplay! Sometimes there is a defined answer: From 'Oinking tendency? (8)' we get both 'pen chant' and 'penchant', and we can see from the letters required (no space) that the second half is the solution. In other cases, this may not be defined at all!

However, even in this simple example we see that this equivalence relationship is not at all straightforward. While 'guide' and 'lead' are synonyms (as verbs in the present tense), it's not true that 'lead' is a synonym for 'metal'. We must also include 'for example' in this relationship too, which causes us to have to discard reflexivity. Although 'metal' can be a clue for 'lead', it's not the case that 'lead' can be a clue for 'metal' (in that case, we signify 'an example of' by writing 'lead, say' or 'bronze, for instance').

We also include abbreviations, which are perhaps more closely related to synonyms, although not usually found in thesauruses, along with some useful 'setters favourites', where an abbreviation of a synonym or of an example is particularly useful for cluing a difficult letter combination used in a wordplay ('Books' becomes 'NT', for 'New Testament').

```
Words → Synonym | Abbreviation | Example
Synonym → String
Abbreviation → String
Example → String
```

The semantic task of evaluating this will be discussed later.

**Concatenation** While not strictly necessary for this grammar (as we have included a concatenation in our metaseantics, we could define multiple definitions of each operator in the form `Operator → Indicator, Wordplay, (Wordplay)*`), it makes sense to add this explicitly as it mirrors the structure of an explanation of a computer solution (i.e. the parse tree).

```
Concatenation → Wordplay (ConcatIndicator) Wordplay
```

This represents a key tool for cluers to create more complex wordplay clues in the form of a charade, where two or more parts can be split out (sometimes syllabically as in 'bath', 'tub', or sometimes otherwise 'bat','htub') and clued separately, and then later joined to form the overall solution.

### 3.5 Other Wordplay Operators

For the other wordplay operators, we define them in terms of our grammar, as well as discussing their semantic meaning.

```
Wordplay → Words | Concatenation | Anagram | Reversion | Contraction
          | Selection | Hidden Word | Containment | Subtraction
          | Homophone
```

These operators all include an indicator word to show they are being applied (as is far more common with operators in programming language parsing!) Each operator will usually have many different indicators (lists of anagram indicators on the web span multiple hundreds). Only select ones are included in the specification here.

### 3.5.1 Unitary Operators

**Anagram** A very commonly used operator in crossword clues is an anagram. These take the form of an indicator word that denotes that the anagram function is being used (called an ‘anagrind’ within cruciverbalist circles), along with the candidate letters to be anagrammed. The simplest form of this gets the candidate letters verbatim from the clue:

`Anagram → Anagrind, String | String, Anagrind`

Sometimes, however sometimes there is some sort of operation applied to the letters before the anagram is applied. For example:

```
Comic bare for short comedy play (7,5)
==> Wordplay 'Comic bare for short comedy' = Definition 'play'
==> Anagram 'bare for short comedy' [anagrind = 'comic']
==> Anagram ("bare for" + Shorten 'comedy')
==> Anagram ("bare fore" + "comed")
==> Anagram ("bare fore" + "comed")
==> Anagram ("bareforecomed")
==> "Bedroom Farce"
```

In which case we find the more general case one proposed structure:

`Anagram → Anagrind, Wordplay | Wordplay, Anagrind`

Wherein we know that the repeated evaluation of the Wordplay will eventually result in a string literal that can be anagrammed. In *Art of the Crossword Puzzle*, Ximenes argued against this form of indirect anagram:

Secondly – and here, for once, I differ from Afrit – I hate what I call an indirect anagram. By that I mean "Tough form of monster" for HARDY (anagram of HYDRA). There may not be many monsters in five letters; but all the same I think the clue-writer is being mean and withholding information which the solver can reasonably demand. Why should he have to solve something before he can begin to use part of a clue? He has first to find "hydra" – and why shouldn't it be "giant"? – and then use the anagrammatic information to help him think of "hardy". ... My real point is that the secondary part of the clue – other than the definition – is meant to help the solver. The indirect anagram, unless there are virtually no alternatives, hardly ever does. He only sees it after he has got his answer by other means.

Even so, most setters that claim to be Ximenean will allow small abbreviations and contractions (to be defined later) to be included in their clues. We therefore must define a new class which includes String Literals as well as the abbreviation where appropriate.

```

Anagram → Anagrind, StringWordplay* | StringWordplay*, Anagrind
Anagrind → “free” | “novel” | “comic” [...]
StringWordPlay → String | Abbreviation | Contraction

```

**Reversion** Clues can also be reversed. While this is functionally a subset of anagrams, there are some crucial differences. Firstly the ‘directionality’ of the clue (i.e. whether it is a ‘down’ or an ‘across’) comes into effect, in determining the sorts of indicators that can form it: “turned back” may only apply to ‘across’ clues, where “taken up” may only apply to ‘down’ clues. Further, these clues are usually taken to be ‘fairer’ game for subsequent operations to be applied to the target of the reversion. Therefore, a clue with nested wordplay such as ("Stressed, made upside-down pudding (7)" = DESSERT) would be acceptable, where an equivalent clue as an anagram ("Stressed, cooked up pudding") would often not be seen as Ximenean.

```

Anagram → ReversionIndictator, Wordplay | Wordplay, ReversionIndictator
ReversionIndicator → “around” | “turned back” | “taken up” [...]

```

**Contraction** Clues of this form range from specific, such of first/last letters ('first in line' = 'l', 'last of the Mohicans' = 'm') to more general operators ('mostly harmless' can yield 'armless', 'harmles', 'harmle'...) whose definitions are more flexible.

```

Contraction → FirstLetterContraction | LastLetterContraction | GeneralContraction
FirstLetterContraction → PreFLCIndicator, Wordplay | Wordplay, PostFLCIndicator
LastLetterContraction → PreLLCIndicator, Wordplay | Wordplay, PostLLCIndicator
GeneralContraction → PreGCIndicator, Wordplay | Wordplay, PostGCIndicator

```

**Selection** There are three similar operators here: A pair which select even or odd letters respectively, and one which takes initial letters across multiple words. These are rarely, if ever, applied to anything other than pure strings. The initials indicator needs to be applied to an argument consisting of multiple words.

```

Selection → Evens | Odds | Initials
Evens → EvensIndicator, String | String, EvensIndicator
Odds → OddsIndicator, String | String, OddsIndicator
Initials → InitialsIndicator, String, “ “, String* |
          String, “ “, String*, InitialsIndicator

```

**Hidden word** The hidden word clue finds a word which appears as a substring (ignoring spaces) inside its operand. These typically only occur once per puzzle, and are always accompanied by a clear indicator. In this example clue:

`'Smack which appears in East Anglian ports.(4)'`

the solution to this example is 'TANG', (meaning 'smack' in the sense of 'taste'), and which is concealed (indicated by 'which appears') in 'east Anglian ports'.

`HiddenWord → HWIndicator String | String HWIndicator`

### Homophone

Also called 'sounds like', this operator produces homophones of a given word. They may be spelled differently ('right' and 'rite') or the same but said differently ('Polish' and 'polish'). This operator is not applied to words that are both spelled and said the same, but with different meanings ('must' as an imperative and 'must' as a noun).

Often, if clues are straightforward, placement of this operator can determine the spelling of the answer.

`We hear twins shave (4)`

yields 'pare' whereas

`Twins shave, we hear (4)`

yields 'pair'. A formulation with the indicator in the middle, in this case, would result in a strong ambiguity. The homophone indicator is only applied to equivalence words, not to clued wordplay.

`Homophone → HomophoneIndicator Words | Words HomophoneIndicator`

### 3.5.2 Binary Operators

As with the unitary operator, each of the arguments of binary operators can be one or more words.

**Containment** Here are two styles of wordplay which are clued very differently, but are actually the same operator, which places one set of letter inside another. This is either presented as a insertion ('end inside ls') or as a containment ('ls around end'). This operation always preserves letter order, unless some nested indicator allows otherwise.

`Containment → Wordplay ContainmentIndicator Wordplay`  
`ContainmentIndicator → "inside" | "around" [...]`

**Subtraction** In a subtraction clue, a number of letters are removed from the target. Usually, the target is some wordplay itself, although sometimes just a string literal. The letters to be subtracted are also often the product of some sort of cluing, although this is usually fairly limited in scope (abbreviations, contractions, first letters of string literals). There are two constraints on this: all the letters from the subtraction set must be in the target, and the length of the subtraction set must be less than the length of the target.

```

Subtraction → SubPreIndicator1 Wordplay (SubPreIndicator2) Wordplay
              | Wordplay SubMidIndicator Wordplay
              | Wordplay (SubPostIndicator1) Wordplay SubPostIndicator2
SubPreIndicator1 → “took”, “without” [...]
SubPreIndicator2 → “from” [...]
SubMidIndicator  → “without” [...]
SubPostIndicator1 → “with” [...]
SubPostIndicator2 → “removed”, “deleted” [...]

```

Semantically here, we have the difference in pre- and post- as the difference between “wanted ant removed” and “removing ant wanted”

The letters in the set are thought to be removed in the order in which they’re found in order to be a properly clued wordplay. Thus “standing” with “tan” removed, gives “sanding”, whereas “ant” cannot be appropriately removed. Note though that the order in which nested clues are applied can change what the set is applied too. If we also had an anagram indicator, as in “Boy muddled standing missing trap” we can apply the muddled to standing to get “dansting” before removing “sting” to get the answer “Dan”.

### 3.6 Meta-references

Sometimes, clues contain references to that cannot be parsed in isolation, or contain a cluing structure that is incompatible with the main model of cluing.

#### 3.6.1 Self reference

A type of clue called an ‘&lit’ clue allows the setter to not include a definition part if the text that makes up the wordplay also can also be read as the definition. Thus in

Spoil vote! (4)

we have the wordplay Anagram (=spoil) “vote” to give ‘VETO’, as well as the clue as a whole ‘spoil vote’ meaning ‘veto’.

#### 3.6.2 Reference to other clues

Some publications will have clues that reference the answer to other clues (‘8 across. Cake made badly by 7 down.’). Sometimes these may also be cyclical (in this example, 7 down would reference 8 across too).

### 3.6.3 Contextual References

Sometimes references will refer outside of the crossword itself. For example, The *Sunday Telegraph* on Easter Sunday 2014 had an anagram clue whose answer was EASTER SUNDAY, and its definition part was "today". In a crossword by setter *Araucaria*, "Araucaria is" coded for IAM (= "I am") as part of an answer.

## 4 References

**Cryptic crossword clues: generating text with a hidden meaning** David Hardcastle - 2007

**The Generation of Cryptic Crossword Clues** G. W. Smith, and J. B. H. du Boulay - 1986

**Crossword Compiler-Compilation** H. Berghel and C. Yi. - 1989

**PROVERB: The Probabilistic Cruciverbalist** Greg A. Keim, Noam M. Shazeer, Michael L. Littman - 1999

**Computer Assisted Analysis of Cryptic Crosswords** P.W.Williams and D. Woodhead - 1977

**LACROSS language, formal definitions - good building material** Cryptic crossword clue interpreter M Hart, RH Davis - 1992

**Microcomputer compilation and solution of crosswords** RH Davis and E J Juvshol - 1985

**Give Us A Clue** Jon G. Hall and Lucia Rapanotti - 2010

**A Statistical Study of Failures In Solving Crossword Puzzles** Naranana, 2010

**Expertise in cryptic crossword performance** Kathryn Friedlander, Philip Fine, 2009

## Part II

# Progress Report

Throughout this initial phase of the project, I have concentrated primarily in getting someway towards creating a system in Haskell working on a limited



subset of clue types, doing end-to-end lexing parsing and evaluating. The reason for this was mainly to spot systematic errors in the way I was approaching the problem, and to identify potential issues with a large-scale implementation.

Indeed, this seems to be the correct method - my early attempts suffered in two ways.

1. I failed to properly utilise Haskell's excellent type system to correctly describe the clues I was parsing, conflating `CONSTRUCTORS` with `DATA TYPES`, among other issues. This meant I was unable to write subsequent expressions that could evaluate partial clues without having to re-parse.
2. I didn't properly separate out the parsing from the evaluation, trying instead to apply insertions and perform anagrams before having completed the full parse tree. This not only creates messy code, but makes an already complicated and expanding state-space problem (one clue can yield multiple parse trees, each of which can have multiple valid evaluations) even worse.

I am currently able to parse a clue made up of a subset of the possible operation types (cons, anagram, abbreviation, synonym, hidden word, reversal) into an abstract syntax tree, and then correctly evaluate that tree into all possible values, and then compare those against the definition using a simplified equality measure to prune for valid trees.

One caveat that this system runs under is that you need to specify the thesaurus definitions manually in Haskell. I tried to generate Haskell clauses for the synonym definitions, but the 200+mb file was bigger than Haskell could process. I am confident that some Haskell magic can remedy this - if not, then I don't believe that this challenges the integrity of the project: a working system could first create a custom thesaurus file in a Haskell readable format before passing the string to Haskell for parsing and evaluation. I have also not done any work on collecting the sorts of troublesome data I reference in my report. That is a very large task indeed, and way beyond the scope of what I am doing. In place of this, I propose to act as an oracle, ruling on the truth value of facts where appropriate.

The biggest challenges I see before me right now are dealing with the state space as my cluing options expand, and possibly translating this into a logic language like QuLog. I have made some small steps towards an initial implementation of some of the parsing in Prolog, and need to solidify this to work out how viable and useful it will be to the project as a whole.

I also plan to do further reading into context-sensitive grammars (for example via DFGs) to see if they could help in what I see as the real next phase of this project, which is optimizing the process to be less sensitive to growing clue length and complexity.

— Michael Skelly