

#### 1. INTRODUCTION TO REGULAR EXPRESSIONS

- 01. What are Regular Expressions?
- 02. History of Regular Expressions
- 03. Software that uses RegExes
- 04.RegExes in Python
- 05.RegExes in Java

#### 2. SIMPLE MATCHING

- 06. TERMINOLOGY: Strings and Things
- 07. Using the Period to Match
- 08. More Details on the Period
- 09. Using the Vertical Bar to Match
- 10. Using Brackets to Group

#### 3. REGEX EXAMPLES 1

- 11. Examples with the Period
- 12. Another Example with the Period
- 13. Examples with the Vertical Bar
- 14. Examples with Brackets
- 15. Examples with Letter Case

#### 4. DEBUGGING REGULAR EXPRESSIONS

- 16. Errors in RegExes
- 17. What is Debugging?
- 18. Patient Debugging
- 19. How to Find a Bug
- 20. How to Fix a Bug

#### 5. QUALIFICATIONS WITH REGEXES

- 21. Character Classes: Ranges
- 22. Using the Question Mark
- 23. Using the Wildcard Star
- 24. Using the Plus Sign
- 25. Using Curly Braces

# What are Regular Expressions?

#### Introduction

Regular expressions are a compact way of searching for text in a document. So let's imagine we have a document that has millions of lines of text in it, including a number of different email addresses, and we want to extract those addresses from the document. Unfortunately, the addresses come in different formats, including:

- DamianTGordon@MyMail.com
- Damian.Gordon@MyMail.com
- Damian.T.Gordon@MyMail.co.uk

So they all have the "@" symbol in their text, and they also have zero, one, or more full stops (".") before the "@" symbol. They also have at least one, but maybe more full stops after the "@" symbol.

To add to the complexity, there are also other phrases in the document that look like emails but aren't, that we don't want to extract, including incorrect email addresses:

- DamianTGordon@MyMail (No domain name, e.g. . COM)
- @MyMail.com (No name before the @ symbol)

As well as some typical text that looks like an email address:

- Can we meet@2pm?
- We are meeting@boardroom.

So a regular expression is a special code we can use to describe the rules of what defines a valid email address (and we know there's a few accepted email address formats), and also how to recognise something that isn't a valid address. The good news is ... the regular expression below describes a valid email address according to the rules we stated above, and <u>I know it looks complicated</u>, but don't worry we'll be explaining each part of this code over the following weeks:

#### ^[A-Za-z+.]+@[A-Za-z.]+[A-Za-z]+\$

This code should work in almost any programming language, but it is worth noting that different programming languages implement regular expressions in slightly different ways, so a regular expression that works in one language might not always work exactly the same way in another, but we'll specify them to work in as many languages as possible, and we'll note when different languages behave differently.

## **Some Terminology**

A Regular Expression is often called a *RegEx*, and a collection of Regular Expressions are called *RegExes*. There are two types of characters used in regular expressions:

- **Metacharacters**: These are characters that have a special meaning.
- Literal characters: These are the actual characters we want to match.

# **History of Regular Expressions**

# **History of Regular Expressions**

In 1943, Warren McCulloch (an American cybernetician) and Walter Pitts (an American logician) developed a computer-based model of learning, modelled on how human brains learn. These types of models are generally referred to as artificial neural networks (ANNs), and their specific model is called a McCulloch-Pitts neural network, whose goal was to learn and recognise patterns. In 1951, American mathematician Stephen Cole Kleene created a formal mathematical language to describe these neural networks, and this language eventually evolved into a general pattern matching notation, which we know as Regular Expressions. This notation can be used to explore the structure of any type of text-based patterns.

#### BIOGRAPHY: Stephen Cole Kleene



Stephen Cole Kleene was born on January 5th, 1909 in Hartford, Connecticut, and died on January 25th, 1994 in Madison, Wisconsin. He is a notable mathematician who helped develop some of the foundations of theoretical computer science (often together with his thesis supervisor, Alonzo Church). He is a founder of the branch of mathematical logic known as recursion theory, and as we know, he invented Regular Expressions in 1951 to describe the McCulloch-Pitts neural net.

#### The Church–Turing thesis

In the 1930s, Alonzo Church (Kleene's dissertation supervisor) developed a general system of logic, that he called Lambda calculus ( $\lambda$ -calculus), to explore the limits of what it is possible to calculate. At around the same time, British mathematician, Alan Mathison Turing, was working on the same problem using his computation model, that later became known as the "Turing Machine", and he identified an approach that would also help explore the limits of what can be calculated, or computed. In 1952, Stephen Cole Kleene published "Introduction to Metamathematics" where he showed that Lambda calculus and Turing Machines are strictly equivalent, and that they are also equivalent to a third approach by Kurt Gödel called "Recursive Functions". He called this the "Church–Turing" thesis, and it serves as a foundational principle in computer science and helps to establish the limits of computability.

The Church–Turing thesis is part of the broader theory of Computation, which looks at the general question of whether a problem can be solved by a computer; and if it can be solved, is it an approximate solution or a very precise one? We can express these questions in a mathematical language, and when looking at problems that examine pattern matching, we can use Regular Expressions to represent them.

# **Software that uses Regular Expressions**

## **RegEx in Software Tools**

One of the first occurrences of the incorporation of regular expressions in software tools was in 1968, when Ken Thompson built Kleene's regex notation into the text editor *QED* that ran on an operating system called *CTSS* (*Compatible Time-Sharing System*), which allowed users to search for, and replace, specific formats of text.

Thompson went on to co-develop the *Unix* operating system in 1969, which proved to be a very popular operating system for many years. He and others incorporated regexes into many software tools, including text editors such as *ed*, *sed*, *vi* and *Emacs*. He also developed other tools that uses regexes such as *grep*, which can be used to search files for text that matches a particular pattern, and this led to the development of *AWK*, which is a scripting language for extracting data from files.

#### BIOGRAPHY: Ken Thompson



Ken Thompson was born on February 4th, 1943, in New Orleans, Louisiana. He is a notable Computer Scientist who helped design and implement the original Unix operating system. He also invented a number of programming languages (including B, Bon, and Go), operating systems (Unix, Plan 9, and Inferno), and utilities such as QED, ed, and grep.

By the 1980s regular expressions were incorporated into a series of formal standards including work by the International Organization for Standardization to create the *ISO SGML* standard. And in 1992 they were incorporated into a standard for operating systems by the IEEE Computer Society in their *POSIX* family of standards.

In 1996, the **PostgreSQL** RDBMS (relational database management system) was developed as a free and open-source RDBMS emphasizing extensibility and its compliance with regular SQL. It incorporates regexes as part of its key features.

In the late 2010s, computer companies started to offer hardware implementations of regexes, including ones that are incorporated into *FPGAs* (*Field-Programmable Gate Arrays*) and *GPUs* (*Graphics Processing Units*).

#### **RegEx in Computer Programming Languages**

We can use regexes in many programming languages, and their functions are either built directly into the programming languages (like in *Perl* and *ECMAScript*) or they can be included into the language by using a software library (which is a collection of additional features that can be optionally included into a programming language to extend the functionality of the language), in languages like *C*, *C++*, *Java*, and *Python*.

# **Regular Expressions in Python**

#### Introduction

The Python programming language was developed by Guido van Rossum starting in 1989, and the first version was released in 1991. It is one of the most widely used and popular programming languages, and is considered one of the easiest programming languages to read, because it uses indentation to show blocks of code.

## RegExes in Python

Python doesn't have regular expressions built into the programming language, but it does have a library (or package) called re that can be imported into Python programs, and that gives those programs a range of RegEx functions.

The program below begins with stating the program name <code>RegEx\_email</code>, and next it imports the regular expression package (import re). Following that the specific regular expression is declared, as well as the text we are going to test it against (these are <code>RegEx\_Pattern</code> and <code>Test\_Message</code> respectively). The pattern is then compared with the specific regular expression using <code>re.match</code> function, and since they will match each other, this program will print out the following message: The <code>pattern</code> matches.

```
# PROGRAM RegEx_email:
import re

RegEx_Pattern = "^[A-Za-z+.]+@[A-Za-z.]+[A-Za-z]+$"
Test_Message = "Damian.T.Gordon@mymail.com"

if (re.match(RegEx_Pattern, Test_Message)):
#then
    print("The pattern matches")
else:
    print("The pattern does not match")
# ENDIF;

# END.
```

#### Other Functions in the re Package

Importing the re package means a range of functions related to regular expressions are available to the program. We've seen one function already, re.match(), which compares a regular expression with some text. Two other important functions are:

- re.search(): This will search some text for a pattern, and it will return the first occurrence of that pattern within the text.
- re.findall(): This will find all occurrences of a pattern in some text.

# **Regular Expressions in Java**

#### Introduction

The Java programming language was developed by James Gosling, Mike Sheridan, and Patrick Naughton starting in 1991, and the first version was released in 1996. It is one of the most widely used and popular programming languages, and was designed with a C/C++-style syntax that programmers would be familiar with.

#### RegExes in Java

Java doesn't have regular expressions built into the programming language, but it does have a library (or package) called <code>java.util.regex</code> that can be imported into Java programs, and that gives those programs a range of RegEx functions.

The program below is similar to the Python program we have seen already, with a few differences; it begins by importing two classes from the <code>java.util.regex</code> package (<code>java.util.regex.Pattern</code> and <code>java.util.regex.Matcher</code>), those are <code>Pattern</code> and <code>Matcher</code>. Next the Main class is declared, as is the Main method. The regular expression is specified using the <code>Pattern</code> class, and the test text is specified using the <code>Matcher</code> class. The test text is compared with the specific regular expression using the <code>Test\_Message.find()</code> function, and since the two match each other, this program will print out the following message:

The pattern matches.

```
import java.util.regex.Pattern;
import java.util.regex.Matcher;

public class Main {
   public static void main(String[] args) {

   Pattern RegEx_Pattern =
        Pattern.compile("^[A-Za-z+.]+@[A-Za-z.]+[A-Za-z]+$");

   Matcher Test_Message =
        RegEx_Pattern.matcher("Damian.T.Gordon@mymail.com");

   if(Test_Message.find()) {
        System.out.println("The pattern matches");
    } else {
        System.out.println("The pattern does not match");
    }
}
```

The java.util.regex package has over 4000 classes in it, so we won't cover them all here, but we'll see some more of them over the coming weeks.

# **Strings and Things**

## Introduction

In computers we often use special terminology to represent fairly common things, the goal of this isn't to make things more complicated at all, but instead it is to be very precise about what we are talking about (because when dealing with computers we know that they want us to be very precise about what we say).

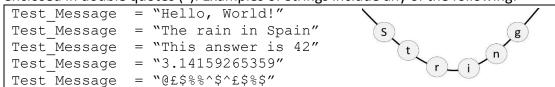
#### Character

A *character* is any single letter, number, or symbol. They are often enclosed in single quotation marks ('). Examples of characters include any of the following:

Test_Message	=	`H'																	
Test_Message	=	'X'					Fur	the	r E	xan	nple	es o	f C	har	act	ers			
Test_Message	=	'd'	ı		!	11	#	\$	%	&	·	(	)	*	+	,	-	•	/
Test Message	=	`j′	П	Θ	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
Test Message	=	14'	П	@	Α	В	С	D	E	F	G	Н	I	J	K	L	М	N	0
Test Message	=	181	H	P	Q	R	S	ď	U e	V	W	X	Y	Z	[	1	] m	n	 0
Test Message	=	<b>'</b> @'	l	р	q	r	s	t	u	v	W	х	у	z	{	_	}	~	7
Test_Message	=	181																	

#### String

The term *string* is used to describe text, because text is made up of a chain of characters (like a string of pearls, where each pearl is a character). They are often enclosed in double quotes ("). Examples of Strings include any of the following:



#### Substring

If we are referring to part of a specific string, we call it a *substring*. They are often enclosed in double quotes ("). If the String is "*Hello, World!*", then the following are examples of Substrings of that string:

```
Test_Message = "Hello"
Test_Message = "World!"
Test_Message = "llo"
Test_Message = "Wor"
Test_Message = "Wor"
Test_Message = "llo, Wor"
```

#### Reminder

We will remember that the key purpose of regular expressions is to create a pattern to locate a specific string, or to locate a collection of strings with a specific format (as we saw with the email example on Page 1). And we can use them for things like searching and replacing text in a Word Processing editor, and validating fields in an online form. We will also learn about other uses for regular expressions, including things such as Web Scraping and Data Mining, that we will discuss in the future.

# **Using the Period to Match**

#### Introduction

This will be our first use of special characters ("metacharacters") to match a String with a regular expression pattern, and the special character we are going to look at is the period (".") character, or the *Full Stop* character. <u>The period metacharacter</u> represents any other single character, which can be a letter, number, or symbol.

## A Simple Example

Sometimes I get emails from people who spell my name wrong, and they call me "Damien" instead of "Damian", which isn't a big deal, but if I had a big document full of emails to me, and I wanted to check how often my name is used in emails, I'd have to search for the string "Damian" and the string "Damien". Regular Expressions give us a nice, simple way of doing this using the period character. We could create a Regular Expression using the period metacharacter as follows:

#### RegEx Pattern = "Dami.n"

So, we can read this Regular Expression as to try to match Strings with the following pattern: "D", "a", "m", "I", any character, and "n".

#### So we would get the following matches:

Test_Message = "Damian"	MATCH ✓
Test_Message = "Damien"	MATCH ✓

#### And it won't match with names that don't fit the pattern:

Test_Message = "Damani"	NO MATCH 🗴
Test_Message = "Dameon"	NO MATCH 😕

#### But, it is worth noting, that the pattern will match many other Strings, including:

, ,	<u> </u>
<pre>Test_Message = "Damixn"</pre>	MATCH ✓
Test_Message = "Dami5n"	MATCH ✓
Test_Message = "Dami=n"	MATCH ✓
Test_Message = "Dami n"	MATCH ✓

Since the period can represent any character, including the space character ("").

We should also remember, our pattern is six characters long ("D", "a", "m", "I", any character, and "n"), so the following Strings will not be considered a match:

Test_Message = "Damin"	NO MATCH X
Test_Message = "Damiaan"	NO MATCH 🗴

#### As they don't follow the pattern.

# **More Details on the Period**

#### Introduction

We will remember that the period metacharacter represents any other single character, which can be a letter, number, or symbol. Now let's look at more advanced uses of the period when searching for patterns in a String (i.e. in a text).

#### **More than One Period**

Sometimes when people are emailing me, instead of "Damian" they type "Damain", so they swap around the "I" and the "a". An easy way to describe this is as follows:

#### RegEx Pattern = "Dam..n"

So, we can read this Regular Expression as to search for a String with the following pattern: "D", "a", "m", any character, any character, and "n". So we would get the following matches:

Test_Message = "Damian"	MATCH ✓
Test_Message = "Damain"	MATCH ✓

However, we know this will also match other Strings including any of the following:

Test_Message = "Damxxn"	MATCH ✓
Test_Message = "Dam35n"	MATCH ✓
Test_Message = "Dam&=n"	MATCH ✓
Test_Message = "Dam n"	MATCH ✓

So, we might match to more things than we want.

But it will not match to these (as the pattern must be six characters long):

Test_Message = "Damn"	NO MATCH 🗴
Test_Message = "Dami an"	NO MATCH 😕

#### **Matching the Period Character**

Finding the actual period character (not the metacharacter) in a String is a bit tricky. So, if we wanted to find something easy like the letter "D" in a String, we could simply set the Regular Expression to be equal to the literal character "D" as follows:

#### RegEx Pattern = "D"

However, if we wanted to find the period character ".", the following won't work:

#### RegEx Pattern = "." 🗶 💢 💢

As it would match any character in the String. So, we use a specific code to represent the actual period character. Mathematically we can represent it as follows: \. However, most programming languages typically prefer we state it as follows:

#### RegEx Pattern = "\\."

And this will allow us to locate a period character (".") in a String, e.g. in an email.

# **Using the Vertical Bar to Match**

#### Introduction

Our next special character or "metacharacter" is the vertical bar ("|"). The vertical bar allows us to choose from a list of alternatives, and is thus also called the "alternation" operator, and also the "Logical OR" operator. <u>The vertical bar</u> metacharacter is used with a list of alternatives, from which any will be a match.

#### Why Use the Vertical Bar?

We will remember that sometimes when people are emailing me, instead of "Damian" they type "Damien", and that we can use the period metacharacter to match both of these, using the following Regular Expression:

#### RegEx Pattern = "Dami.n"

However, this is an *approximate solution*, as we know this will also match other variations, including: "Damisn", "Dami7n", "Dami@n", "Dami n", and many others. If we want a *precise solution*, that will match only to the Strings "Damian" or "Damien", we can use the vertical bar with the following Regular Expression:

#### RegEx Pattern = "Damian|Damien"

This Regular Expression matches with EITHER "Damian" OR "Damien".

My friend Rebecca has a lot more trouble than I do, in her emails people address her as "Rebecca", "Becca", "Becks", and "Becky", so if she asked us to search for her name, we could use the following Regular Expression:

## RegEx\_Pattern = "Rebecca|Becca|Becks|Becky"

This reads as EITHER "Rebecca" OR "Becca" OR "Becks" OR "Becky".

#### **Combining Metacharacters**

We notice that two of Rebecca's nicknames are very similar: "Becks" and "Becky". So we could use the Period metacharacter to represent these using the following Regular Expression: "Beck.", with the Vertical Bar for the other alternatives:

## RegEx Pattern = "Rebecca|Becca|Beck."

And this is a good approximate solution that will match any of the required names, but will match a few others as well, like "Beckr" or "Becko". We could use an even shorter Regular Expression: "Bec..." to cover three alternatives:

## RegEx\_Pattern = "Rebecca|Bec.."

And this will match any of the required names, but also match lots of other Strings like "Becss", "Bechg" and "Becdh", but it's a more compact Regular Expression.

#### Matching the Vertical Bar Character

If we are searching for the actual vertical bar character (not the metacharacter), then mathematically we represent the actual vertical bar as follows: \ | However, most programming languages typically prefer we state it as follows:

#### RegEx Pattern = "\\|"

And this will allow us to locate the vertical bar character ("|") in a String.

# **Using the Brackets to Group**

#### Introduction

Our next special characters or "metacharacters" are the brackets - () or []. We'll be mainly focusing on Round Brackets (also known as *Parentheses*), but we'll also look at the square brackets [] a little bit as well. The brackets allow us to group together operations, and is thus also called the "Grouping" operator. *The Brackets metacharacters are used to group together parts of a Regular Expression*.

#### Why Use the Round Brackets?

We will remember that if we are searching for "Damian" or "Damien", we can use the period metacharacter to give us an *approximate solution*:

RegEx Pattern = "Dami.n"

And for a *precise solution*, we can use the vertical bar:

RegEx\_Pattern = "Damian|Damien"

Using the brackets we can restate the *precise solution* more compactly, as follows:

#### RegEx Pattern = "Dami(a|e)n"

Which can be read as creating a Regular Expression to match to any String with the letters: "D", "a", "m", "i", ("a" or "e"), and "n". So, the brackets allow specific parts of the Regular Expression that can be grouped together.

In the example above we are creating a choice between two single characters ("a" and "e"), however, we will see in later examples that the brackets can also be used to create choices between multicharacter Strings as well.

#### **Character Class**

A Character Class (or Character Set) is a list of characters enclosed in square brackets, and any one of those characters will represent a match to the class. So for the moment we can say that the Regular Expression with the single characters above, and Character Class below are equivalent in terms of their functionality:

#### Character Class = "Dami[ae]n"

This can be read as creating a Character Class to match to any String with the letters "D", "a", "m", "I", ("a" or "e"), and "n".

#### Matching the Brackets Character

If we are searching for the Round Brackets characters (not the metacharacter), then mathematically we represent the Open Bracket as follows: \ ( and for the Close Bracket we can mathematically represent it as follows: \ ) However, most programming languages prefer we state the Open Bracket as:

RegEx Pattern = "\\("

and most programming languages prefer we state the Close Bracket as:

## RegEx Pattern = "\\)"

And it is the same for the Square Brackets: " $\setminus [$ " and " $\setminus ]$ ".

# **Examples with the Period**

#### Introduction

We will remember that the period metacharacter <u>represents any other single</u> <u>character, which can be a letter, number, or symbol</u>. And that we can match "Damian" and "Damien" using the following Regular Expression:

### RegEx Pattern = "Dami.n"

And, we note that it does match to other Strings as well including: "Damixn", "Dami5n", "Dami=n", and "Dami n", so it's an approximate solution.

#### **String Length**

If we were looking for a String that has to be at least 5 characters long, we could create a regular expression of 5 periods, and that will only match with Strings that are at least 5 characters long, as follows:

RegEx Pattern = "...."

So we would get the following matches (and many others):

Test_Message = "abcde"	MATCH ✓
<pre>Test_Message = "H3llo"</pre>	MATCH ✓
Test_Message = "a z"	MATCH ✓

And it wouldn't match with a String that was 4 characters (or less) long:

Test_Message = "wxyz'	•	NO MATCH 😕
Test_Message = "1234"	1	NO MATCH 🗴

#### **String Format**

An ISBN number (International Standard Book Number) is a code that is used to uniquely identify books. It takes the format of a series of digits separated by dashes, with a total of 13 digits (3 digits followed by 2 digits, followed by 5 digits, followed by 2 digits, followed by 1 digit). For example:

- 978-02-41953-53-2
- 978-03-74537-14-2
- 978-04-25054-71-0

The digits in the ISBN can be represented by the period, but the dash character is a special character in using Regular Expressions, so it needs to be proceeded with two slashes (as we've seen before with other special characters), as follows: \\-

So bringing it all together, a potential Regular Expression for ISBNs is as follows:

```
RegEx_Pattern = "...\\-..\\-...\\-.."
```

This will only match a collection of digits (or numbers or symbols) in the ISBN format. #RegExThursday © Damian Gordon

# **Another Example with the Period**

#### Introduction

We remember that the period metacharacter <u>represents any other single character</u>, which can be a letter, number, or symbol.

#### **String Format**

Mobile telephone numbers in Ireland have a prefix of either 085, 086 or 087, and then are followed by seven digits. For example:

- (085) 1234567
- (086) 2345678
- (087) 3456789

## So let's look at each individual element of that format:

Element	Regular Expression	
Open Bracket	We can express this as: \\ (	
Prefix	It's either 085, 086 or 087, so let's express it as: 08.	
Close Bracket	We can express this as: \\)	
Space	Let's use the space character itself to match with space	
Seven Digits	We can use seven periods to match the seven digits	

So bringing it all together, the Regular Expression is the following:

#### RegEx Pattern = "\\(08.\\) ....."

And, this will match to any of the phone numbers presented above, as follows:

Test_Message	= "(085) 1234567"	MATCH ✓
Test_Message	= <b>"</b> (086) 2345678 <b>"</b>	MATCH ✓
Test_Message	= "(087) 3456789"	MATCH ✓

## And it wouldn't match with a String that doesn't match that format:

Test_Message = "(085)1234567"	NO MATCH 😕
Test_Message = "(08)5 1234567"	NO MATCH 🗴
Test_Message = "(085) 123456"	NO MATCH 🗴

But we will remember that this Regular Expression also matches many other Strings that might be invalid phone numbers, for example:

Test_Message = "(08A)	ABCDEFG"	MATCH ✓
Test_Message = "(080)	!£\$%^&*"	MATCH ✓
Test_Message = "(08)	11	MATCH ✓

So the Period will match letters and symbols as well as numbers.

# **Examples with the Vertical Bar**

#### Introduction

We will remember that the vertical bar metacharacter <u>is used to choose from a list</u> <u>of alternatives, where any will be a match</u>. And we saw that if we want to match with the Strings "Damian" or "Damien", we can use the following:

#### RegEx Pattern = "Damian|Damien"

This Regular Expression matches with EITHER "Damian" OR "Damien". And we saw that the same thing can be done as follows:

#### RegEx Pattern = "Dami(a|e)n"

And we noted that sometimes when people are emailing me, instead of "Damian" they type "Damain"; they swap around the "I" and "a", and we used this RegEx:

#### RegEx Pattern = "Dam..n"

However, we could also do the same as follows (and get a more *precise solution*):

#### RegEx\_Pattern = "Dam(ia|ai)n"

And if we wanted to take in account both the people who confuse the spelling of the name, and those who confuse the order of letters, we could do the following:

#### RegEx\_Pattern = "Dam(ia|ai|ie|ei)n"

This Regular Expression matches with any of the following four Strings, EITHER "Damian" OR "Damain" OR "Damien" OR "Damein".

#### **Clear Regexes**

It is really important that we make our Regular Expressions as clear as possible, so that if we have to modify them in a few months, or someone else has to, it will be easy to understand the purpose of the RegEx. This is not an exact science because something that is very clear to you might be less so to others, but it is important to be aware that we should try to make our RegExes as easy to understand as possible. If we can add comments in to explain to purpose of the RegEx, that will help a lot, but we should also think about how we present the RegEx. As an example, we'll look at searching for a Period or a Vertical Bar character.

So, if we are looking for the Period character, we put two slashes in front of it:

#### RegEx Pattern = "\\."

And it's the same for a Vertical Bar character:

## RegEx\_Pattern = "\\|"

So to search for one or the other of them, we can express it in one of two ways, either as Vertical Bar character or Period character ("|" OR "."):

## RegEx\_Pattern = "\\||\\."

Or as Period character or Vertical Bar character ("." OR "|"):

#### RegEx Pattern = "\\.|\\|"

I, personally, find the second of these Regexes easier to read, the first one (slash, slash, bar, <u>bar</u>, slash, slash, period) is harder for me to find the OR symbol, than the second one (slash, slash, period, <u>bar</u>, slash, slash, bar), but it really is a matter of perspective, and it is always worth getting someone else to check if your Regexes are clear to them.

# **Examples with Brackets**

#### Introduction

We will remember that the brackets metacharacters <u>are used to group together</u> <u>parts of a Regular Expression</u>. And we saw that if we want to match with the Strings "Damian" or "Damien", we can do it as follows

#### RegEx Pattern = "Dami(a|e)n"

Which can be read as creating a Regular Expression to match to any String with the letters: "D", "a", "m", "i", ("a" or "e"), and "n".

#### **Telephone Numbers**

We saw that mobile telephone numbers in Ireland have a prefix of either 085, 086 or 087, and then are followed by seven digits. For example:

- (085) 1234567
- (086) 2345678
- (087) 3456789

And we suggested a Regular Expression as follows:

```
RegEx_Pattern = "\\(08.\\) ...."
```

However, given that there are only three choices for the prefix, we can do refine this Regular Expression using the Vertical Bar and Brackets as follows:

#### RegEx Pattern = "\\(08(5|6|7)\\) ....."

Which searches for either 085, 086 or 087. And as we noted before, a Character Class (or Character Set) uses Square Brackets to achieve the same thing as the above selection, but without the need for the Vertical Bar notation, so the above RegEx can equivalently be expressed as follows:

RegEx Pattern = "\\(08[567]\\) ....."

#### **Movie Titles**

The "Jaws" series of movies (1975-1987) concerns the saga of a great white shark and its attacks on people, the titles of the first four movies are:

- Jaws
- Jaws II
- Jaws 3D
- Jaws The Revenge

If we wanted to search for any of these titles in a text, we could do so, as follows:

#### RegEx Pattern = "Jaws(\$| II| 3D| The Revenge)"

In this case we are looking for a title starting with the String "Jaws" followed by either: Nothing OR " II" or " 3D" or " The Revenge". So the interesting thing here is that to represent nothing, we use the Dollar symbol (\$), and this means "End of String", so, in other words, that there is no further text after the word "Jaws".

# **Examples with Letter Cases**

#### Introduction

When we are working with Strings, it's important to note that Capital Letters (Uppercase Letters) are treated as being different from Lowercase Letters. So, for example, if the Regular Expression is the Uppercase Letter "D", as follows:

#### RegEx Pattern = "D"

This will match with the String "D", but not to the String "d". We say that Regular Expressions are *Case Sensitive* to describe this property of Regular Expressions.

#### **Uppercase and Lowercase Letters**

So, the following three Strings should be considered as three different Strings:

- hello
- Hello
- HELLO

And to detect these three Strings, we could use the following Regular Expression:

## RegEx Pattern = "hello|Hello|HELLO"

And we can reduce the length of that Regular Expression by doing the following:

RegEx\_Pattern = "(h|H)ello|HELLO"

#### **Lowercase Letters**

If we had a Notepad file full of text and we wanted to read all of that text into a computer program, and then only print out the letters that are in Lowercase, there is a really easy way to do that using Regular Expressions. We just read in the file one character at a time, and then check if that is a Lowercase character (or not) by comparing it to the following Regular Expression:

# RegEx\_Pattern = "(a|b|c|d|e|f|g|h|i|j|k|l|m|n|o|p|q|r|s|t|u|v|w|x|y|z)"

So, if the character read in is a Lowercase letter, it will match the RegEx pattern above, and we can print it out, and if it is anything other than a Lowercase letter, it will not match and therefore we should not print it out. So, for example:

Test_Message	= "a"	MATCH ✓
Test_Message	= "z"	MATCH ✓

#### And any Uppercase letters, Numbers and Symbols will not match, for example:

Test_Message = "A"	NO MATCH X
Test_Message = "4"	NO MATCH 😕
Test_Message = "@"	NO MATCH 😕

We also note that the same Pattern can be created by using a Character Class:

#### RegEx Pattern = "[abcdefghijklmnopqrstuvwxyz]"

And this will work exactly the same way, but the Pattern is more compact.

# **Errors in Regular Expressions**

#### Introduction

Creating Regular Expressions is like every other skill we learn in life, the more we practice doing them, the better we will get at them. When we start writing them, we might construct them wrong sometimes, so we need to check whether or not the Regular Expression is doing what we want it to do, or not.

#### A Simple Example

Let's imagine we were doing our example of searching for "Damian" or "Damien", but we put the period in the wrong place, so instead of "Dami.n", we did this:

## RegEx Pattern = "Dam.an"

we can easily see there is an error (or "bug") in the RegEx, because it's a simple example. However, if we tested it on the word "Damian", it would work, but if we tested it on the word "Damien", it wouldn't work. So, we need to test our Regular Expressions on things we expect it to match, and things it shouldn't match with.

#### How long will it really take?

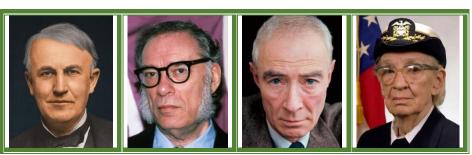
When people are learning to write Regular Expressions by themselves, a lot of them experience a similar cycle; in the beginning, it takes them about 5 minutes to write a Regular Expression, and 25 minutes to fix it (to get it working exactly the way they want). After a few weeks, one of two things happen; either they give up and decide Regular Expressions are too hard for them, or they hang in there, and after a few more weeks they become <u>Regular Expression Wizards</u>, which doesn't mean that they get everything right first time, it just means it takes a much shorter period of time to fix any issues, so, maybe 5 minutes to write them, and 5 minutes to fix them.

The lessons we are doing here will make it a lot easier to understand how to write Regular Expressions, and because we are doing it one page at a time, we are learning in a slow and steady manner, so we should never get to a place where the content is too hard for us, but it is also important that we aren't too hard on ourselves. So, we can't expect ourselves to be brilliant at Regular Expressions immediately, and even once we become proficient using them, we should still treat them like writing an essay or an e-mail, we should have the expectation that we will have to draft and redraft a few times to make the content better quality. So, this is my golden rule:



So as long as we don't expect that it will all work out first time, everything is easier. #RegExThursday © Damian Gordon

# What is Debugging?



Thomas Edison, Isaac Asimov, J. Robert Oppenheimer, Grace Hopper

## Debugging is ...

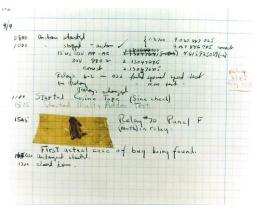
Sometimes errors in programs are called "bugs", so we have a special name for finding and fixing errors in computer programs, we call it "debugging" (in other words, taking the bugs out). These terms are not exclusive to computers, as far back as the 1870s, Thomas Edison uses the term in a letter, where he says: "then that "Bugs"—as such little faults and difficulties are called—show themselves".

This term was used extensively in the 1930s and 1940s to describe *flaws* or *glitches* in mechanical devices. In 1944, the writer Isaac Asimov used the term in a fictional context describing potential errors in robots, in his short story "*Catch That Rabbit*", published in 1944: "*U. S. Robots had to get the bugs out of the multiple robots, and there were plenty of bugs, and there are always at least half a dozen bugs left for the field-testing.*"

That same year, on October 27<sup>th</sup>, in a letter from theoretical physicist, J. Robert Oppenheimer, when he was discussing the building of the first atomic bomb, he mentions, when discussing staff recruitment, that most of the existing staff are "occupied in getting into operation and debugging" the bomb.

On September 9th, 1947, computer developer Grace Hopper was tracing an error on the Harvard Mark II electromechanical computer.

Along with one of the operators, Bill Burke, they found a moth trapped in a relay that was the cause of the error, so they taped the moth into the logbook, and recorded it as the first actual bug.



The notation reads:

"First actual case of bug being found."

# **Patient Debugging**

#### Fixing Regular Expressions Sometimes Isn't Easy

Now that we have seen some Regular Expressions, it's important to note that writing the Regular Expression often isn't the hard bit, it's when we try to run it, and it gives us lots and lots of wrong matches, and sometimes it can be hard to figure out exactly what the problem is. So we have to have the patience and persistence to review each part of the Regular Expression (as well as the rest of the program that the RegEx is in) to see if we can identify the problem. And often we have to stare at the RegEx for a few minutes before we will, in a flash, figure out what's wrong. This is not easy, and it requires a lot of determination, because often when we have figured out (and fixed) one bug, another one follows. So, we have to type in our Regexes very carefully, and review each part as we write it. And when we are finished writing a RegEx, and go to run it, we need to accept that it may not work correctly the first time, and when we fix the initial bug, there may be another, and another.

#### **Breathing**

Fixing RegExes can be a bit stressful, and the longer we are working on a single one, the more frustrated we can get. This can lead to short-term thinking, where we move parts of the RegEx around at random in hopes that it will fix itself. Try to avoid doing this, and try to avoid getting stressed, by breathing. There are a variety of breathing techniques that can be used to calm down, including the following:

- **Left Nostril Breathing**: As the name suggests, just close your right nostril off, and breath in and out through your left nostril slowly, with your eyes closed. This creates a calming effect in your nervous systems within minutes.
- **7-2-11** Breathing: Breath in through your nose for 7 seconds, hold the breath for 2 seconds, and exhale through your mouth for 11 seconds. This takes a bit of practice, but after a few days of 4-8 sessions a day, you will master it.

#### Take a Break

One important trick to know is when to take a break; so if I am staring at an bug and I can't figure it out, my rule of thumb is after 7 minutes I walk away from the computer and get a glass of water and stop thinking about it for 2-3 minutes. More often than not as soon as I return to the computer, I know exactly the bug is, because I gave my unconscious mind time to work on it, and can fix it in no time.

#### **Cardboard Helper**

Sometimes the easiest way to fix a bug is to ask someone for help, and you will find that as soon as you say to someone "Excuse me, can you help me with this problem...", even before you have even outlined the problem, you know what the solution is, because you got a chance to think about it in a different way. In fact, you don't really need another person, just get a cardboard cutout and ask them for help.

# How to Find a Bug

#### **Debugging Approaches**

There are two parts to debugging, the locating, and fixing, of bugs. We'll look at locating the bugs first, so if the RegEx (or RegExes) runs but doesn't give us the results we are expecting, there is some bug in the RegEx (or RegExes) that we need to find (we should also check the rest of the program). The computer is only doing what it is told, so there must be a wrong instruction somewhere. There are a number of debugging approaches that can be taken to find that instruction:

- Brute Force Approach: The is probably the most common approach to
  debugging, and it typically involves splitting the Regular Expression into parts
  from the start, and running each part separately to make sure we understand
  what it is doing. There are also tools that can be used in this approach, these
  include both tracing tools and debugging tools, e.g. <a href="https://regex101.com/">https://regex101.com/</a>
- Backtracking Approach: The backtracking approach is exactly what it sounds like, you start at the end of the RegEx, manually reviewing each part of the RegEx, or each metacharacter in the RegEx to see if it is correctly written, until the incorrect instruction is found.
- Cause Elimination Approach: This approach involves creating a list of
  possible causes (or hypothesis) for the bug, and initial tests are carried out to
  eliminate each hypothesis. Of the ones that cannot be eliminated in the
  initial testing, further tests are carried out to eliminate more and more
  hypotheses, until there is only one cause left. The bug is then located.

#### One More Thing....

This may be just me, but when I'm trying to debug a program, and I don't feel like I'm making progress; sometimes if I recite a verse of poetry, or part of a song, and I do that a couple of times, and it gives me the fortitude to continue and succeed. The two verses below are the ones I most commonly use.

There's nothing you can do that can't be done;
Nothing you can sing that can't be sung;
Nothing you can say, but you can learn how to play the game;
It's easy.

From the song "All You Need Is Love" by John Lennon and Paul McCartney (1967)

Great bugs have little bugs upon their backs to bite 'em,
And little bugs have lesser bugs, and so ad infinitum.
And the great bugs themselves, in turn, have greater bugs to go on;
While these again have greater still, and greater still, and so on.

"Siphonaptera" from Augustus De Morgan's A Budget of Paradoxes (1872)

# How to Fix a Bug

#### **Read the Error Message**

If the bug is producing an error message, that's great, read it carefully. The error messages generally can give us a clue as to what the issue might be, or at least where to start looking from. It can also be incredibly useful for us to put the error message into Google, and search for it; we will often find that someone else has experienced this same bug and has found a solution.

## **Beware of Side-Effects**

Fixing a bug should be done very carefully, if the fix is done poorly, it can introduce other issues into the RegEx, and do more harm than good, so it's important to be careful when fixing bugs. American software engineer Tom Van Vleck outlined three simple questions that we should ask ourselves before making a fix:

- Is this bug (or a similar bug) likely to appear in another RegEx?
   In many cases a programmer will use the same type of logic throughout multiple RegExes, so if an error is found in one RegEx, it may be worth reflecting on whether or not there are other RegExes that has similar functionality.
- 2. What new issue might be introduced into the RegEx when fixing the bug?
  Before the bug is fixed, it is a good idea to explore the design of the RegEx
  (and the program around it), to check the context in which the RegEx is being used in. If it depends on other programs or RegExes, then it is important to carefully check what the consequences could be of any changes made.
- 3. What can be done to prevent this same bug from happening again?
  This is the first step in creating a good Software Quality Assurance Process going forward. If there was a problem in the development process that caused the bug to occur, fixing that issue with the process will prevent similar bugs from occurring in the future.

#### **Other Terms for Bugs**

Depending on who is discussing bugs, they may use different terms to describe them, so for example, someone in IT Sales might call them *features*, whereas an IT Tester might call them *issues*. Here are some other terms for bugs:

Defects	Faults	Problems	Incidents
Anomalies	Inconsistencies	Variances	Failures
Mistakes	Exceptions	Errors	Side Effects

This is a small sampling of the range of terms used for bugs.

# **Character Classes: Ranges**

#### Introduction

We are already seen the idea of a Character Class (or Character Set), it's <u>a list of</u> <u>characters enclosed in square brackets, and any one of those characters will</u> <u>represent a match to the class.</u> So, for example, the following Regular Expression:

## RegEx Pattern = "Dami(a|e)n"

Can be more compactly represented as follows:

Character\_Class = "Dami[ae]n"

#### **Character Ranges**

If we wanted to match only with Lowercase letters, we could do either this:

RegEx Pattern =

||(a|b|c|d|e|f|g|h|i|j|k|l|m|n|o|p|q|r|s|t|u|v|w|x|y|z)||

or by using a Character Class:

RegEx Pattern = "[abcdefghijklmnopqrstuvwxyz]"

However, we can also represent the same range of letters as follows:

#### RegEx Pattern = "[a-z]"

Which means to match to any character in the range "a" to "z". So, in other words, the Regular Expression will match with any single Lowercase letter.

We can create a similar Regular Expression for Uppercase letter, as follows:

#### RegEx Pattern = "[A-Z]"

Which means to match to any character in the range "A" to "Z".

And we can create a similar Regular Expression for numbers, as follows:

#### RegEx Pattern = "[0-9]"

Which means to match to any number in the range "0" to "9".

#### **Combining Ranges**

So, if we wanted to search for a String of five characters long, with the first one being Uppercase and the rest being Lowercase, we could do that as follows:

#### RegEx Pattern = "[A-Z][a-z][a-z][a-z]"

And that Pattern will match to Upper, Lower, Lower, Lower, Lower, for example:

Test_Message = "Hello"	MATCH ✓
Test_Message = "Aaaaa"	MATCH ✓
Test_Message = "Abcde"	MATCH ✓

And it wouldn't match with a String that doesn't match that format:

Test_Message	= "AAaaa"	NO MATCH 😕
Test_Message	= "aaaaa"	NO MATCH 😕

And we will see that Character Ranges are very common and very useful in RegExes.

# **Using the Question Mark**

#### Introduction

We are now moving onto looking at "qualification metacharacters" which are metacharacters that allow us to specific the number of instances of a single character or a grouping of characters in a Regular Expression. Our first qualification metacharacters is the question mark ("?"). <u>The question mark metacharacter is matched if there is zero or one instances of the preceding character (or grouping) in a Regular Expression</u>.

#### Why Use the Question Mark?

If we were searching for the word "colour", but we are dealing with a collection of documents some of which are from the USA and some from Europe; we know that in the USA we will be looking for the word "color", so we can create a Regular Expression to describe this as follows:

#### RegEx Pattern = "colour|color"

But a more compact way of saying the same thing is to state that every time we find the word "colour", the "u" may appear either zero or one times in that word. So, if it appears zero times, then the word is "color" and if it appears one time, then the word is "colour". This can be stated using the Question Mark as follows:

## RegEx\_Pattern = "colou?r"

So it's a more compact way of matching with either "colour" and "color", as the character preceding the Question Mark ("u") can appear zero or one times.

#### **Matching with Multiple Characters**

If we were looking for documents about Regular Expressions, and we were looking for the phrase "Regular Expression" or "RegEx" in them, we could do it as follows:

## RegEx\_Pattern = "Regular Expression|RegEx"

We can achieve the same result using the Question Mark metacharacter in conjunction with the round brackets as follows:

## RegEx\_Pattern = "Reg(ular )?Ex(pression)?"

So it will match with either "Regular Expression" or "RegEx", as the two strings preceding the question marks ("ular" and "pression") can appear zero or one times.

#### **Matching the Question Mark Character**

If we are searching for the actual question mark character (not the metacharacter), then mathematically we represent it as follows:  $\$ ?

However, most programming languages prefer we state the Question Mark as:

## RegEx\_Pattern = "\\?"

And this will allow us to locate the question mark character in a String.

# **Using the Wildcard Star**

#### Introduction

Another "qualification metacharacter" is the Wildcard Star ("\*"), which is commonly also known as the Kleene Star, or the Asterisk. it is similar to the Question Mark metacharacter but looks for zero or more instances of a character or group of characters. The Wildcard Star metacharacter is matched if there is zero or more instances of the preceding character (or grouping) in a Regular Expression.

#### Why Use the Wildcard?

If we are trying to match with either "colour" and "color", just for fun, we could do by placing the wildcard after the letter 'u', as follows:

#### RegEx Pattern = "colou\*r"

So, this would match with either "colour" and "color", but it would also match with "colouur", "colouuur", "colouuuur", etc.

If I get an email with good news, I will reply with the phrase "Yippeeeee", the number of e's varies, depending on how good the news is, so if I had a collection of emails that we needed to search for these, we could do so as follows:

# RegEx Pattern = "Yippee\*"

And this would match with "Y", "I", "p", "p", "e", and zero or more "e"'s.

If we had an online form, and one of the fields in the form could be left blank, or it could be filled in with any number of characters, we could that as follows:

#### RegEx Pattern = ".\*"

Which means any character (period), any number of times (wildcard star).

If we wanted the field to be either blank or any amount of lowercase characters:

#### RegEx Pattern = "[a-z]\*"

Which means any lowercase character ([a-z]), any number of times (wildcard star). If we wanted the field to be either blank or any amount of uppercase characters:

#### RegEx Pattern = "[A-Z]\*"

Which means any uppercase character ([A-Z]), any number of times (wildcard star). And if we wanted the field to be either blank or any amount of numerical characters:

## RegEx Pattern = "[0-9]\*"

Which means any numerical character ([0-9]), any number of times (wildcard star).

#### Matching the Wildcard Star Character

If we are searching for the actual star character (not the metacharacter), then mathematically we represent it as follows: \\*

However, most programming languages prefer we state the star as:

#### RegEx Pattern = "\\\*"

And this will allow us to locate the star character in a String.

# **Using the Plus Sign**

#### Introduction

Another "qualification metacharacter" is the Plus Sign ("+"). It is also similar to the Question Mark and the Wildcard Star but looks for one or more instances of a character or group of characters. <u>The Plus metacharacter is matched if there is one</u> or more instances of the preceding character (or grouping) in a Regular Expression.

## Why Use the Plus?

If we are trying to match with either "colour" and "color", and we tried to do it by placing the plus after the letter 'u' (it would be wrong), as follows:

#### RegEx Pattern = "colou+r"

This is wrong because it would match with "colour", but it would not match with "color", as the "u" must appear at least once. And we remember that It would also match with "colouur", "colouuur", "colouuuur", "colouuuur", etc.

However, if we wanted to search a text file for all instances of "Yippee", including "Yippe" and "Yippeeeeeee" (so the number of e's varies), we could do so as follows:

## RegEx Pattern = "Yippe+"

And this would match with "Y", "I", "p", "p", and one or more "e"'s. This is slightly more compact than the equivalent Wildcard version we have seen ("Yippee\*").

If we had an online form, and one of the fields in the form had to have at least one character (or more) in it, we could that as follows:

#### RegEx Pattern = ".+"

Which means any character (period), at least one times or more times (plus sign).

If we wanted the field to have one or more lowercase characters:

## RegEx\_Pattern = "[a-z]+"

Which means any lowercase character ([a-z]), one or more times (plus sign). If we wanted the field to have one or more uppercase characters:

## RegEx\_Pattern = "[A-Z]+"

Which means any uppercase character ([A-Z]), one or more times (plus sign). And if we wanted the field to one or more numerical characters:

# RegEx Pattern = "[0-9]+"

Which means any numerical character ([0-9]), one or more times (plus sign).

#### **Matching the Plus Character**

If we are searching for the actual plus character (not the metacharacter), then mathematically we represent it as follows: \+

However, most programming languages prefer we state the plus sign as:

#### RegEx Pattern = "\\+"

And this will allow us to locate the plus character in a String.

# **Using Curly Braces**

#### Introduction

Another "qualification metacharacter" is the curly braces { } characters, which are commonly also known as braces, or curly brackets. They are used to specify the number of instances of a character or grouping of characters. The Curly Braces metacharacters are a repetition qualification that specifies the number of instances of the preceding character (or grouping) in a Regular Expression.

#### Why Use the Curly Braces?

If we know the number of times a character, or a group of characters is going to repeat, or even if we know a minimum and/or maximum number of occurrences of a character or grouping, we should use the Curley Braces. We can use it in four ways, to indicate that the preceding character (or grouping) appears the following times:

{num}	Appears exactly num times.
{min,}	Appears at least min times.
{,max}	Appears up to max times.
{min,max}	Appears between min and max times.

For example, if we are looking for the phrase "aaaaa", we can do:

RegEx Pattern = "a{5}"

If we are trying to match no less than 4 "a" characters in a row, we can do:

## RegEx Pattern = "a{4,}"

And this will match with "aaaa", "aaaaaa", etc., but not "a", "aa" or "aaa".

If we are trying to match no more than 4 "a" characters in a row, we can do:

# RegEx Pattern = "a{,4}"

And this will match with "a", "aa", "aaa", and "aaaa".

If we aren't sure exactly how many times "a" will appear in a row, but we know it will be at least 2 times, but no more than 5 times, we can do:

#### RegEx Pattern = "a{2,5}"

And this will match with "aa", "aaaa", "aaaa" and "aaaaa", but nothing else.

## **Matching the Brackets Character**

If we are searching for the Curly Braces characters (not the metacharacter), then mathematically we represent the Open Curly Brace as follows: \ { and for the Close Bracket we can mathematically represent it as follows: \ } However, most programming languages prefer we state the Open Curly Brace as:

# RegEx Pattern = "\\{"

and most programming languages prefer we state the Close Curly Brace as:

#### RegEx Pattern = "\\}"

And this will allow us to locate Curly Braces in a String.



# GLOSSARY OF TERMS FOR REGULAR EXPRESSIONS



# **Regular Expressions Glossary**

- **ANN**, see artificial neural network.
- **Approximate Solution**, this is when someone is searching for a particular set of *Strings*, and if the *Regular Expression* successfully matches that set of *Strings*, but can also match to additional *Strings* that are not specifically being searched for. The opposite is a *Precise Solution*.
- **Artificial Neural Network**, a computer program that can learn and recognise patterns. They are inspired by biological brains.
- **AWK** is a scripting language that can be used to extract data from files, developed by Alfred Aho, Peter Weinberger, and Brian Kernighan in the 1970s.
- C, a programming language developed by Dennis Ritchie in the 1970s.
- C++, a programming language developed by Bjarne Stroustrup in the 1980s.
- Case Sensitive, this refers to a property of some computer programs that differentiates between Uppercase and Lowercase letters. So, for example, if the program treats the letter "D" as being different from the letter "d", it is Case Sensitive, but if they are treated as being the same, it's Case Insensitive.
- Character, a single letter, number or symbol.
- **Church–Turing Thesis**, a collection of mathematical approaches that explore which types of problems have a solution, and which ones do not.
- **Compatible Time-Sharing System**, an operating system developed by IBM and the Massachusetts Institute of Technology in the 1960s.
- **CTSS,** see Compatible Time-Sharing System.
- **ECMAScript**, a script programming language developed by Brendan Eich in the 1990s.
- ed, a text editor developed by Ken Thompson in the 1970s.
- **Emacs**, a text editor developed by David A. Moon, Guy L. Steele Jr., and Richard Stallman in the 1970s.

# **Regular Expressions Glossary**

- **Empty String**, see *Null String*.
- **Field-Programmable Gate Arrays**, a special form of integrated circuit that can be configured and after it has been manufactured.
- **FPGA**, see *Field-Programmable Gate Array*.
- **Graphics Processing Units**, a specialized electronic circuit designed to accelerate the display of images and video on a display device.
- **grep**, a command-line tool for searching files for text that matches a pattern on *Unix*-type operating systems, developed by Ken Thompson in the 1970s.
- **GPU**, see *Graphics Processing Unit*.
- International Organization for Standardization, an international standards organization made up of representatives from the national standards organizations of member countries. It was founded in 1947.
- **ISO**, see International Organization for Standardization.
- **ISO SGML**, an ISO standard for a general way to add categorizing codes to text files called the "Standard Generalized Markup Language". It was based on work developed by Charles Goldfarb, Edward Mosher, and Raymond Lorie in the 1960s.
- Java, a programming language developed by James Gosling in the 1990s.
- Lambda calculus ( $\lambda$ -calculus), a formal mathematical system of logic that models problems that have a solution. It was developed by Alonzo Church in the 1930s.
- **Library**, see *Software Library*.
- **Literal Character**, a *character* that is being searched for in a *string*.
- **Metacharacter**, a special *character* that can represent other *characters*.
- McCulloch-Pitts neural network, an artificial neural network that is a simple model ("caricature model") of how biological brains work.

# **Regular Expressions Glossary**

- **Null String**, A *String* with no characters in it (often expressed as "a *String* of zero length"). This is also called an *Empty String*.
- Package, see Software Library.
- Perl, a script programming language developed by Larry Wall in the 1980s.
- **Portable Operating System Interface**, a family of standards specified by the IEEE Computer Society for maintaining compatibility between operating systems. It was released in 1988.
- **POSIX**, see *Portable Operating System Interface*.
- **PostgreSQL**, a *Relational Database Management System* developed by the PostgreSQL Global Development Group in the 1990s.
- **Precise Solution**, this is when someone is searching for a particular set of *Strings*, and if the *Regular Expression* successfully matches only to that exact set of *Strings*. The opposite is an *Approximate Solution*.
- Python, a programming language developed by Guido van Rossum in the 1980s.
- **QED**, a text editor developed by Butler Lampson, L. Peter Deutsch, and Dana Angluin in the 1960s.
- **RDBMS**, see *Relational Database Management System*.
- **Recursive Function**, a computer program where the number of times each part of the program will execute is known before the program runs.
- **RegEx,** another name for a *regular expression*.
- RegExes, a collection of regular expressions.
- **Regular Expression**, a special search pattern that can combine *literal characters* and *metacharacters*.
- Relational Database Management System, a system that stores and accesses data in tables that are highly structured. It is based on a structure defined by E. F. Codd in 1970.

# **Regular Expressions Glossary**

- **sed**, a text editor developed by Lee E. McMahon in the 1970s.
- **Software Library**, a collection of functions that can be used by computer programs, but aren't built directly into the programming language.
- **SQL (Structured Query Language)**, a database design language used in managing databases. Developed by Donald Chamberlin and Raymond Boyce in the 1970s.
- **String**, a piece of text made up of a sequence of *characters*. It is analogous to a string of pearls, where each pearl is a *character*.
- **Substring**, part of another *string*, or more formally, a contiguous sequence of *characters* within a *string*.
- **Theory of Computation**, a branch of theoretical computer science and mathematics that deals with which problems can be solved, and how efficiently.
- **Turing Machine**, a simple mathematical concept that models problems that have a solution. It was developed by Alan Turning in the 1930s.
- **Unix**, an operating system developed by Ken Thompson, Dennis Ritchie, Brian Kernighan, Douglas McIlroy, and Joe Ossanna in the 1960s.
- vi, a text editor developed by Bill Joy in the 1970s.