

**An empirical study of regular expression use in practice, sampling from Python projects on Github, leading to new concepts for refactoring regular expressions for readability.**

by

Carl Allen Chapman

A thesis submitted to the graduate faculty  
in partial fulfillment of the requirements for the degree of  
MASTER OF SCIENCE

Major: Computer Science

Program of Study Committee:  
Kathryn Stolee, Major Professor  
Samik Basu  
Tien Nguyen

Iowa State University  
Ames, Iowa  
2016

Copyright © Carl Allen Chapman, 2016. All rights reserved.

## DEDICATION

I would like to dedicate this thesis to my mother, who believed in me and supported me through many years on a long winding road leading to a satisfying career. I'd also like to thank my wife Chien Wen Hung and our cat Siva.

## TABLE OF CONTENTS

<b>LIST OF TABLES</b> . . . . .	v
<b>LIST OF FIGURES</b> . . . . .	vi
<b>ACKNOWLEDGEMENTS</b> . . . . .	vii
<b>ABSTRACT</b> . . . . .	i
<b>CHAPTER 1. OVERVIEW</b> . . . . .	1
1.1 Background, terms, examples . . . . .	1
1.2 The 5 experiments conducted as part of this thesis, and the intentions behind them . . . . .	1
1.3 Introduction From Refactoring . . . . .	3
<b>CHAPTER 2. RELATED WORK</b> . . . . .	4
2.1 Regex in Languages . . . . .	4
2.2 Tools that depend on Regex . . . . .	4
2.3 Research on Regex . . . . .	4
2.4 Regex Analysis Tools . . . . .	4
2.5 Related Work From Features . . . . .	4
2.6 Related Work From Refactoring . . . . .	6
<b>CHAPTER 3. Feature Analysis</b> . . . . .	8
3.1 Study . . . . .	8
3.1.1 Research Questions . . . . .	9
3.1.2 Survey Design and Implementation . . . . .	10
3.1.3 Regex Corpus . . . . .	10

3.1.4	Analyzing Features . . . . .	13
3.1.5	Clustering and Behavioral Similarity . . . . .	13
3.2	Results . . . . .	15
3.2.1	RQ1: How do developers use regexes? . . . . .	16
3.2.2	RQ2: How is the <b>re</b> module used? . . . . .	18
3.2.3	RQ3: Regex language feature usage . . . . .	19
3.2.4	RQ4: Regex behavioral similarity . . . . .	23
<b>CHAPTER 4.</b>	<b>Conclusion . . . . .</b>	<b>31</b>
4.1	Discussion From Features . . . . .	31
4.1.1	Implications For Tool Designers . . . . .	31
4.1.2	Opportunities For Future Work . . . . .	32
4.2	Discussion From Refactoring . . . . .	35
4.2.1	Interpreting Results . . . . .	35
4.2.2	Opportunities For Future Work . . . . .	36
4.2.3	Threats to Validity . . . . .	39
4.3	Conclusion . . . . .	41
<b>APPENDIX A.</b>	<b>Patterns in Python projects from Github . . . . .</b>	<b>42</b>
<b>APPENDIX B.</b>	<b>Developer Survey . . . . .</b>	<b>43</b>
<b>APPENDIX C.</b>	<b>Mechanical Turk Study . . . . .</b>	<b>44</b>
<b>APPENDIX D.</b>	<b>Community Analysis . . . . .</b>	<b>45</b>
<b>BIBLIOGRAPHY</b>	<b>. . . . .</b>	<b>46</b>

## LIST OF TABLES

3.1	Survey results for number of regexes composed per year by technical environment (RQ1) . . . . .	16
3.2	Survey results for regex usage frequencies for activities, averaged using a 6-point likert scale: Very Frequently=6, Frequently=5, Occasionally=4, Rarely=3, Very Rarely=2, and Never=1 (RQ1) . . . . .	17
3.3	How saturated are projects with utilizations? (RQ2) . . . . .	18
3.4	How frequently do features appear in projects? . . . . .	28
3.5	Survey results for preferences between custom character and default character classes (RQ3) . . . . .	29
3.6	Survey results for regex usage frequencies, averaged using a 6-point likert scale: Very Frequently=6, Frequently=5, Occasionally=4, Rarely=3, Very Rarely=2, and Never=1 (RQ3) . . . . .	29
3.7	Sample from an example cluster (RQ4) . . . . .	29
3.8	Cluster categories and sizes (RQ4) . . . . .	30

## LIST OF FIGURES

3.1	Example of one regex utilization . . . . .	8
3.2	Two patterns parsed into feature vectors . . . . .	13
3.3	A similarity matrix created by counting strings matched . . . . .	14
3.4	Creating a similarity graph from a similarity matrix . . . . .	14
3.5	How often are <code>re</code> functions used? (RQ2) . . . . .	19
3.6	Which behavioral flags are used? (RQ2) . . . . .	19

## ACKNOWLEDGEMENTS

I would like to take this opportunity to express my thanks to those who helped me with various aspects of conducting research and the writing of this thesis. First and foremost, Dr. Kathryn Stolee for her guidance, patience and support throughout this research and the writing of this thesis. I would also like to thank my committee members for their efforts and contributions to this work: Dr. Samik Basu and Dr. Tien Nguyen.

# ABSTRACT

## Abstract

Though regular expressions (regex) provide a powerful search technique that is baked into every major language, is incorporated into a myriad of essential tools, and has been a fundamental aspect of Computer Science since Kleene in 1956, no one has ever formally studied how they are used in practice, or what can be done to make them easier to understand. This thesis presents the original work of studying a sample of regex taken from Python projects pulled from Github, determining what features are used most often, defining some categories that illuminate common use cases, and identifying areas of significance for tool builders. Furthermore, this thesis defines an equivalence class model used to explore comprehension of regex, identifying the most common and most understandable representations of semantically identical regex, suggesting several refactorings and preferred representations. Opportunities for future work include the novel and rich field of regex refactoring, semantic search of regexes, and further fundamental research into regex usage and understandability.



## CHAPTER 1. OVERVIEW

### 1.1 Background, terms, examples

### 1.2 The 5 experiments conducted as part of this thesis, and the intentions behind them

Regular expressions (regexes) are an abstraction of keyword search that enables the identification of text using a pattern instead of an exact keyword.

Regexes are commonly used for parsing text using a general purpose language like Python, validating content entered into web forms using Javascript, and searching text files for a particular pattern using tools like grep, vim or Eclipse.

Although regexes are powerful and versatile, they can be hard to understand, maintain, and debug, resulting in tens of thousands of bug reports Spishak et al. (2012).

Due in part to their common use across programming languages and how susceptible regexes are to error, many researchers and practitioners have developed tools to support more robust regex creation Spishak et al. (2012) or to allow visual debugging Beck et al. (2014). Other research has focused on learning regular expressions from text Babbar and Singh (2010); Li et al. (2008), avoiding human composition altogether. Researchers have also explored applying regexes to test case generation Ghosh et al. (2013); Galler and Aichernig (2014); Anand et al. (2013); Tillmann et al. (2014), as specifications for string constraint solvers Trinh et al. (2014); Kiezun et al. (2013) and using regexes as queries in a data mining framework Begel et al. (2010). Regexes are also employed in critical missions like MySQL injection prevention Yeole and Meshram (2011) and network intrusion detection network (2015), or in more diverse applications like DNA sequencing alignment Arslan (2005).

Regex researchers and tool designers must pick what features to include or exclude, which can be a difficult design decision. Supporting advanced features may be more expensive, taking more time and potentially making the project too complex and cumbersome to execute well. A selection of only the simplest of regex features limits the applicability or relevance of that work. Despite extensive research effort in the area of regex support, no research has been done about how regexes are used in practice and what features are essential for the most common use cases.

*The goal of this work is to explore 1) the context in which developers use regular expressions, and 2) the features and similarities of regular expressions found in Python<sup>1</sup> projects.*

First, we survey professional developers about how they use regexes and their pain points. Second, we gather a sample of regexes from Python projects and analyze the frequency of feature usage (e.g., kleene star: `*` and the end anchor: `$` are features). Third, we investigate what features are supported by four large projects that aim to support regex usage (brics Møller (2010), hampi Kiezun et al. (2013), Rex Veanes et al. (2010), and RE2 re2 (2015)), and which features are not supported, but are frequently used by developers. Finally, we cluster regular expressions that appear in multiple projects by behavior, investigating high-level behavioral themes in regex usage.

Our results indicate that regexes are most frequently used in command line tools and IDEs. Capturing the contents of brackets and searching for delimiter characters were some of the most apparent behavioral themes observed in our regex clusters, and developers frequently use regexes to parse source code. The contributions of this work are:

- A survey of 18 professional software developers about their experience with regular expressions,
- An empirical analysis of regex feature usage in nearly 14,000 regular expressions in 3,898 open-source Python projects, mapping of those features to those supported by common regex tools and survey results showing the impact of not supporting various features,

---

<sup>1</sup>*Python is the fourth most common language on GitHub (after Java, Javascript and Ruby) and Python's regex pattern language is close enough to other regex libraries that our conclusions are likely to generalize.*

- An approach for measuring behavioral similarity of regular expressions and qualitative analysis of the most common behaviorally similar clusters, and
- An evidence-based discussion of opportunities for future work in supporting programmers who use regular expressions, including refactoring regexes, developing regex similarity analyses, and providing migration support between languages.

### 1.3 Introduction From Refactoring

Regular expressions are used frequently by developers for many purposes, such as parsing files, validating user input, or querying a database. Regexes are also employed in MySQL injection prevention Yeole and Meshram (2011) and network intrusion detection network (2015). However, recent research has suggested that regular expressions (regexes) are hard to understand, hard to compose, and error prone Spishak et al. (2012). Given the difficulties with working with regular expressions and how often they appear in software projects and processes, it seems fitting that efforts should be made to ease the burden on developers.

Tools have been developed to make regexes easier to understand, and many are available online. Some tools will, for example, highlight parts of regex patterns that match parts of strings as a tool to aid in comprehension.<sup>2</sup> Others will automatically generate strings that are matched by the regular expressions Kiezun et al. (2013). Other tools will automatically generate regexes when given a list of strings to match Babbar and Singh (2010); Li et al. (2008). The commonality of such tools provides evidence that people need help with regex composition and understandability.

In software, code smells have been found to hinder understandability of source code Abbes et al. (2011); Du Bois et al. (2006). Once removed through refactoring, the code becomes more understandable, easing the burden on the programmer. In regular expressions, such code smells have not yet been defined, perhaps in part because it is not clear what makes a regex smelly.

---

<sup>2</sup><https://regex101.com/>

## CHAPTER 2. RELATED WORK

### 2.1 Regex in Languages

### 2.2 Tools that depend on Regex

### 2.3 Research on Regex

### 2.4 Regex Analysis Tools

**TODO.NOW: Create Intro For Related Section?** This is the opening paragraph to my thesis which explains in general terms the concepts and hypothesis which will be used in my thesis.

With more general information given here than really necessary.

### 2.5 Related Work From Features

Regular expressions have been a focus point in a variety of research objectives. From the user perspective, tools have been developed to support more robust creation Spishak et al. (2012) or to allow visual debugging Beck et al. (2014). Building on the perspective that regexes are difficult to create, other research has focused on removing the human from the creation process by learning regular expressions from text Babbar and Singh (2010); Li et al. (2008).

Regarding applications, regular expressions have been used for test case generation Ghosh et al. (2013); Galler and Aichernig (2014); Anand et al. (2013); Tillmann et al. (2014), and as specifications for string constraint solvers Trinh et al. (2014); Kiezun et al. (2013). Regexes are also employed in MySQL injection prevention Yeole and Meshram (2011) and network intrusion

detection network (2015), or in more diverse applications like DNA sequencing alignment Arslan (2005) or querying RDF data Lee et al. (2010); Alkhateeb et al. (2009).

As a query language, lightweight regular expressions are pervasive in search. For example, some data mining frameworks use regular expressions as queries (e.g., Begel et al. (2010)). Efforts have also been made to expedite the processing of regular expressions on large bodies of text Baeza-Yates and Gonnet (1996).

Research tools like Hampi Kiezun et al. (2013), and Rex Veanes et al. (2010), and commercial tools like bricsMøller (2010) and RE2 re2 (2015), all support the use of regular expressions in various ways. Hampi was developed in academia and uses regular expressions as a specification language for a constraint solver. Rex was developed by Microsoft Research and generates strings for regular expressions that can be used in applications such as test case generation Anand et al. (2013); Tillmann et al. (2014). Brics is an open-source package that creates automata from regular expressions for manipulation and evaluation. RE2 is an open-source tool created by Google to power code search with an efficient regex engine.

Mining properties of open source repositories is a well-studied topic, focusing, for example, on API usage patterns Linares-Vásquez et al. (2014) and bug characterizations Chen et al. (2014). Exploring language feature usage by mining source code has been studied extensively for Smalltalk Callaú et al. (2011, 2013), JavaScript Richards et al. (2010), and Java Dyer et al. (2014); Grechanik et al. (2010); Parnin et al. (2013); Livshits et al. (2005), and more specifically, Java generics Parnin et al. (2013) and Java reflection Livshits et al. (2005). To our knowledge, this is the first work to mine and evaluate regular expression usages from existing software repositories. Related to mining work, regular expressions have been used to form queries in mining framework Begel et al. (2010), but have not been the focus of the mining activities. Surveys have been used to measure adoption of various programming languages Meyerovich and Rabkin (2013); Dattero and Galup (2004), and been combined with repository analysis Meyerovich and Rabkin (2013), but have not focused on regexes.

## 2.6 Related Work From Refactoring

Regular expression understandability has not been studied directly, though prior work has suggested that regexes are hard to read and understand since there are tens of thousands of bug reports related to regular expressions Spishak et al. (2012). To aid in regex creation and understanding, tools have been developed to support more robust creation Spishak et al. (2012) or to allow visual debugging Beck et al. (2014). Building on the perspective that regexes are difficult to create, other research has focused on removing the human from the creation process by learning regular expressions from text Babbar and Singh (2010); Li et al. (2008).

Regular expression refactoring has also not been studied directly, though refactoring literature abounds Mens and Tourwé (2004); Opdyke (1992); Griswold and Notkin (1993). The closest to regex refactoring comes from research toward expediting the processing of regular expressions on large bodies of text Baeza-Yates and Gonnet (1996), which could be thought of as refactoring for performance.

Code smells in object-oriented languages were introduced by Fowler (1999). Researchers have studied the impact of code smells on program comprehension Abbes et al. (2011); Du Bois et al. (2006), finding that the more smells in the code, the harder the comprehension. This is similar to our work, except we aim to identify which regex representations can be considered smelly. Code smells have been extended to other language paradigms including end-user programming languages Hermans et al. (2012, 2014); Stolee and Elbaum (2011, 2013). The code smells identified in this work are representations that are not common or not well understood by developers. This concept of using community standards to define smells has been used in other refactoring literature for end-user programmers Stolee and Elbaum (2011, 2013).

Exploring language feature usage by mining source code has been studied extensively for Smalltalk Callaú et al. (2011), JavaScript Richards et al. (2010), and Java Dyer et al. (2014); Grechanik et al. (2010); Parnin et al. (2013); Livshits et al. (2005), and more specifically, Java generics Parnin et al. (2013) and Java reflection Livshits et al. (2005). Our prior work (Chapman and Stolee (2016), under review) was the first to mine and evaluate regular expression

usages from existing software repositories. The intention of the prior work Chapman and Stolee (2016) was to explore regex language features usage and surveyed developers about regex usage. In this work, we define potential refactorings and use the mined corpus to find support for the presence of various regex representations in the wild. Beyond that, we measure regex understandability and suggest canonical representations for regexes to enhance conformance to community standards and understandability.

## CHAPTER 3. Feature Analysis

### 3.1 Study

To understand how programmers use regular expressions in Python projects, we scraped 3,898 Python projects from GitHub, and recorded regex usages for analysis. Throughout the rest of this paper, we employ the following terminology:

**Utilization:** A *utilization* occurs whenever a regex appears in source code. We detect utilizations by statically analyzing source code and recording calls to the `re` module in Python. Within a source code file, a utilization is composed of a function, a pattern, and 0 or more flags. Figure 3.1 presents an example of one regex utilization, with key components labeled. The function call is `re.compile`, `(0|-?[1-9][0-9]*)$` is the regex string, or pattern, and `re.MULTILINE` is an (optional) flag. When executed, this utilization will compile a regex object in the variable `r1` from the pattern `(0|-?[1-9][0-9]*)$`, with the `$` token matching at the end of each line because of the `re.MULTILINE` flag. Thought of another way, a regex utilization is one single invocation of the `re` library.

**Pattern:** A *pattern* is extracted from a utilization, as shown in Figure 3.1. In essence, it is a string, but more formally it is an ordered series of regular expression language feature tokens. The pattern in Figure 3.1 will match if it finds a zero at the end of a line, or a (possibly

	function	pattern	flags
<code>r1 =</code>	<code>re.compile</code>	<code>(0 -?[1-9][0-9]*)\$</code>	<code>re.MULTILINE</code>

Figure 3.1 Example of one regex utilization



negative) integer at the end of a line (i.e., due to the `-?` sequence denoting zero or one instance of the `-`).

Note that because the vast majority of regex features are shared across most general programming languages (e.g., Java, C, C#, or Ruby), a Python pattern will (almost always) behave the same when used in other languages, whereas a utilization is not universal in the same way (i.e., it may not compile in other languages, even with small modifications to function and flag names). As an example, the `re.MULTILINE` flag, or similar, is present in Python, Java, and C#, but the Python `re.DOTALL` flag is not present in C# though it has an equivalent flag in Java.

In this work, we primarily focus on patterns since they are cross-cutting across languages and are the primary way of specifying the matching behavior. Next, we describe the research questions, data set collection and analysis.

### 3.1.1 Research Questions

To understand the contexts in which regexes are used and feature usage, we perform a survey of developers and explore regular expressions found in Python projects on GitHub. We aim to answer the following research questions:

**RQ1:** In what contexts do professional developers use regular expressions?

We designed and deployed a survey about when, why, and how often they use regular expressions. This was completed by 18 professional developers at a small software company.

**RQ2:** How is the `re` module used in Python projects?

We explore invocations of the `re` module in 3,898 Python projects scraped from GitHub.

**RQ3:** Which regular expression language features are most commonly used in Python?

We consider regex language features to be tokens that specify the matching behavior of a regex pattern, for example, the `+` in `ab+`. All studied features are listed and described in Table 3.4 with examples. We then map the feature coverage for four common regex support

tools, brics, hampi, RE2 and Rex, and explore survey responses regarding feature usage for some of the less supported features.

**RQ4:** How behaviorally similar are regexes across projects?

As this is a first step in understanding behavioral overlap in regexes, we measure similarity between pairs of regexes by overlap in matching strings. For each regex, matching strings are generated and then evaluated against each other regex to compute pairwise similarity. Then we use clustering to form behaviorally similar groupings.

### 3.1.2 Survey Design and Implementation

To understand the context of when and how programmers use regular expressions, we designed a survey, implemented using Google Forms, with 40 questions. The questions asked about regex usage frequency, languages, purposes, pain points, and the use of various language features.<sup>1</sup> Participation was voluntary and participants were entered in a lottery for a \$50 gift card.

Our goal was to understand the practices of professional developers. Thus, we deployed the survey to 22 professional developers at Dwolla, a small software company that provides tools for online and mobile payment management. While this sample comes from a single company, we note anecdotally that Dwolla is a start-up and most of the developers worked previously for other software companies, and thus bring their past experiences with them. Surveyed developers have nine years of experience, on average, indicating the results may generalize beyond a single, small software company, but further study is needed.

### 3.1.3 Regex Corpus

Our goal was to collect regexes from a variety of projects to represent the breadth of how developers use the language features. Using the GitHub API, we scraped 3,898 projects containing Python code. We did so by dividing a range of about 8 million repo IDs into 32

---

<sup>1</sup>[https://github.com/softwarekitty/tour\\_de\\_source/blob/master/regex\\_usage\\_in\\_practice\\_survey.pdf](https://github.com/softwarekitty/tour_de_source/blob/master/regex_usage_in_practice_survey.pdf)

sections of equal size and scanning for Python projects from the beginning of those segments until we ran out of memory. At that point, we felt we had enough data to do an analysis without further perfecting our mining techniques. We built the AST of each Python file in each project to find utilizations of the `re` module functions. In most projects, almost all regex utilizations are present in the most recent version of a project, but to be more thorough, we also scanned up to 19 earlier versions. The number 20 was chosen to try and maximize returns on computing resources invested after observing the scanning process in many hours of trial scans. All regex utilizations were obtained, sans duplicates. Within a project, a duplicate utilization was marked when two versions of the same file have the same function, pattern and flags. In the end, we observed and recorded 53,894 non-duplicate regex utilizations in 3,898 projects.

In collecting the set of distinct patterns for analysis, we ignore the 12.7% of utilizations using flags, which can alter regex behavior. An additional 6.5% of utilizations contained patterns that could not be compiled because the pattern was non-static (e.g., used some runtime variable). The remaining 80.8% (43,525) of the utilizations were collapsed into 13,711 distinct pattern strings. Each of the pattern strings was pre-processed by removing Python quotes (`'\\W'` becomes `\\W`), unescaping escaped characters (`\\W` becomes `\W`) and parsing the resulting string using an ANTLR-based, open source PCRE parser<sup>2</sup>. This parser was unable to support 0.5% (73) of the patterns due to unsupported unicode characters. Another 0.2% (25) of the patterns used regex features that we chose to exclude because they appeared very rarely (e.g., reference conditions). An additional 0.1% (16) of the patterns were excluded because they were empty or otherwise malformed so as to cause a parsing error.

The 13,597 distinct pattern strings that remain were each assigned a weight value equal to the number of distinct projects the pattern appeared in. We refer to this set of weighted, distinct pattern strings as the *corpus*.

**TODO.NOW: JOINT SEcTION FROM OTHER PAPER: To determine how common each regex representations is in the wild, we collected regexes from GitHub projects. We specifically targeted Python projects as it is a popular programming language with a strong presence on GitHub. Further, Python is the**

---

<sup>2</sup><https://github.com/bkiers/pcre-parser>

fourth most common language on GitHub (after Java, Javascript and Ruby) and Python’s regex pattern language is close enough to other regex libraries that our conclusions are likely to generalize.

**TODO.NOW:** We collected and analyzed static invocations to the Python `re` library. Figure 3.1 presents an example from Python with key components labeled. The *function* called is `re.compile`. **TODO.NOW:** The *pattern* defines what strings will be matched and the *flag* `re.MULTILINE` modifies the rules used by the regex engine when matching. When executed, a regex object `r1` is created and it will match if it finds a zero at the end of a line, or a (possibly negative) integer at the end of a line (i.e., due to the `-?` sequence denoting zero or one instance of the `-`).

**TODO.NOW:** Our goal was to collect regex patterns from a variety of projects to represent the breadth of how developers use regexes. We scraped 3,898 projects containing Python code using the GitHub API. This was done by systematically selecting repository IDs, checking the repository for Python files, and retaining the project if Python was found. After dividing eight million repository IDs into 32 groups, we scanned from the beginning until we had collected approximately four thousand Python projects. We did so by dividing a range of about 8 million repo IDs into 32 sections of equal size and scanning for Python projects from the beginning of those segments until we ran out of memory. At that point, we felt we had enough data to do an analysis without further perfecting our mining techniques.

**TODO.NOW:** Within a project, a duplicate utilization was marked when two versions of the same file have the same function, pattern and flags. In the end, we observed and recorded 16,088 non-duplicate patterns in 1,645 projects.

**TODO.NOW:** In collecting the set of distinct patterns for analysis, we ignore the 12.7% of `re` invocations using flags, which can alter regex behavior. An additional 6.5% of `re` invocations contained patterns that could not be compiled because the pattern was non-static (e.g., used some runtime variable). **TODO.NOW:** The remaining 80.8% (43,525) of the utilizations were collapsed into 13,711 dis-

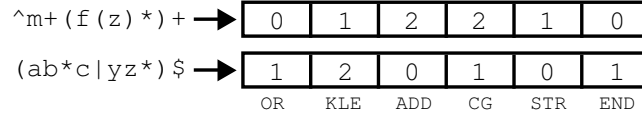


Figure 3.2 Two patterns parsed into feature vectors

tinct pattern strings. This parser was unable to support 0.8% (114) due to error. **TODO.NOW:** After removing all problematic patterns as described and collapsing on duplicates, we ended up with 13,597 distinct patterns from 1,544 projects remained to be used in this study.

### 3.1.4 Analyzing Features

For each escaped pattern, the PCRE-parser produces a tree of feature tokens, which is converted to a vector by counting the number of each token in the tree. For a simple example, consider the patterns in Figure 3.2. The pattern ‘^m+(f(z)\*)+’ contains four different types of tokens. It has the kleene star (KLE), which is specified using the asterisk ‘\*’ character, additional repetition (ADD), which is specified using the plus ‘+’ character, capture groups (CG), which are specified using pairs of parenthesis ‘(...)’ characters, and the start anchor (STR), which is specified using the caret ‘^’ character at the beginning of a pattern. A list of all features and abbreviations is provided in Table 3.4.

Once all patterns were transformed into vectors, we examined each feature independently for all patterns, tracking the number of patterns and projects that the each feature appears in at least once.

### 3.1.5 Clustering and Behavioral Similarity

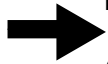
An ideal analysis of regex behavioral similarity would use subsumption or containment analysis. However, we struggled to find a tool that could facilitate such an analysis. Further, regular expressions in code libraries (e.g., for Python, Java) are not the same as regular languages in formal language theory. Some features of regular expression libraries, such as backreferences, make the libraries more expressive than regular languages. This allows a regular expression pattern to match, for example, repeat words, such as “cabcab”, using the pattern ([a-z]+\1.

Pattern A matches 100/100 of A's strings		
Pattern B matches 90/100 of A's strings		
Pattern A matches 50/100 of B's strings		
Pattern B matches 100/100 of B's strings		

	A	B
A	1.0	0.9
B	0.5	1.0

Figure 3.3 A similarity matrix created by counting strings matched

	A	B	C	D
A	1.0	0.0	0.9	0.0
B	0.2	1.0	0.8	0.7
C	0.6	0.8	1.0	0.2
D	0.0	0.6	0.1	1.0



	A	B	C	D
A	1.0			
B	0.1	1.0		
C	0.75	0.8	1.0	
D	0.0	0.65	0.15	1.0

Figure 3.4 Creating a similarity graph from a similarity matrix

However, building an automaton to recognize such a pattern and to facilitate containment analysis, is infeasible. For these reasons, we developed a similarity analysis based on string matching.

Our similarity analysis clusters regular expressions by their behavioral similarity on matched strings. Consider two unspecified patterns A and B, a set  $\mathbf{mA}$  of 100 strings that pattern A matches, and a set  $\mathbf{mB}$  of 100 strings that pattern B matches. If pattern B matches 90 of the 100 strings in the set  $\mathbf{mA}$ , then B is 90% similar to A. If pattern A only matches 50 of the strings in  $\mathbf{mB}$ , then A is 50% similar to B. We use similarity scores to create a similarity matrix as shown in Figure 3.3. In row A, column B we see that B is 90% similar to A. In row B, column A, we see that A is 50% similar to B. Each pattern is always 100% similar to itself, by definition.

Once the similarity matrix is built, the values of cells reflected across the diagonal of the matrix are averaged to create a half-matrix of undirected similarity edges, as illustrated in Figure 3.4. This facilitates clustering using the Markov Clustering (MCL) algorithm<sup>3</sup>. We chose MCL because it offers a fast and tunable way to cluster items by similarity and it is particularly useful when the number of clusters is not known *a priori*.

---

<sup>3</sup><http://micans.org/mcl/>

In the implementation, strings are generated for each pattern using Rex Veanes et al. (2010). Rex generates matching strings by representing the regular expression as an automaton, and then passing that automation to a constraint solver that generates members for it<sup>4</sup>. If the regex matches a finite set of strings smaller than 400, Rex will produce a list of all possible strings. Our goal is to generate 400 strings for each pattern to balance the runtime of the similarity analysis with the precision of the similarity calculations.

For clustering, we prune the similarity matrix to retain all similarity values greater than or equal to 0.75, setting the rest to zero, and then using MCL. This threshold was selected based on recommendations in the MCL manual. The impact of lowering the threshold would likely result in either the same number of more diverse clusters, or a larger number of clusters, but is unlikely to markedly change the largest clusters or their summaries, which are the focus of our analysis for RQ4 (Section 3.2.4), but further study is needed to substantiate this claim. We also note that MCL can also be tuned using many parameters, including inflation and filtering out all but the top-k edges for each node. After exploring the quality of the clusters using various tuning parameter combinations, the best clusters (by inspection) were found using an inflation value of 1.8 and k=83. The top 100 clusters are categorized by inspection into six categories of behavior.

The end result is clusters and categories of highly behaviorally similar regular expressions, though we note that this approach has a tendency to over-approximate the similarity of two regexes. We measure similarity based on a finite set of generated strings, but some regexes match an infinite set (e.g., `ab*c`), so measuring similarity based on the first 400 strings may lead to an artificially high similarity value. To mitigate this threat, we chose a large number of generated strings for each regex, but future work includes exploring other approaches to computing regex similarity.

## 3.2 Results

Next, we present the results of each research question.

---

<sup>4</sup><http://research.microsoft.com/en-us/projects/rex/>

Table 3.1 Survey results for number of regexes composed per year by technical environment (RQ1)

Language/Environment	0	1-5	6-10	11-20	21-50	51+
General (e.g., Java)	1	6	5	3	1	2
Scripting (e.g., Perl)	5	4	3	3	2	1
Query (e.g., SQL)	15	2	0	0	1	0
Command line (e.g., grep)	2	5	3	2	0	6
Text editor (e.g., IntelliJ)	2	5	0	5	1	5

### 3.2.1 RQ1: How do developers use regexes?

The survey was completed by 18 participants (82% response rate) that identified as software developer/maintainers. Respondents have an average of nine years of programming experience ( $\sigma = 4.28$ ). On average, survey participants report to compose 172 regexes per year ( $\sigma = 250$ ) and compose regexes on average once per month, with 28% composing multiple regexes in a week and an additional 22% composing regexes once per week. That is, 50% of respondents uses regexes at least weekly. Table 3.1 shows how frequently participants compose regexes using each of several languages and technical environments. Six (33%) of the survey participants report to compose regexes using general purpose programming languages (e.g., Java, C, C#) 1-5 times per year and five (28%) do this 6-10 times per year. For command line usage in tools such as grep, 6 (33%) participants use regexes 51+ times per year. Yet, regexes were rarely used in query languages like SQL. Upon further investigation, it turns out the surveyed developers were not on teams that dealt heavily with a database.

Table 3.2 shows how frequently, on average, the participants use regexes for various activities. Participants answered questions using a 6-point likert scale including very frequently (6), frequently (5), occasionally (4), rarely (3), very rarely (2), and never (1). Averaging across participants, among the most common usages are capturing parts of a string and locating content within a file, with both occurring somewhere between occasionally and frequently.

Using a similar 7-point likert scale that includes ‘always’ as a seventh point, developers indicated that they test their regexes with the same frequency as they test their code (average



Table 3.2 Survey results for regex usage frequencies for activities, averaged using a 6-point likert scale: Very Frequently=6, Frequently=5, Occasionally=4, Rarely=3, Very Rarely=2, and Never=1 (RQ1)

Activity	Frequency
Locating content within a file or files	4.4
Capturing parts of strings	4.3
Parsing user input	4.0
Counting lines that match a pattern	3.2
Counting substrings that match a pattern	3.2
Parsing generated text	3.0
Filtering collections (lists, tables, etc.)	3.0
Checking for a single character	1.7

response was 5.2, which is between frequently and very frequently). Half of the developers indicate that they use external tools to test their regexes, and the other half indicated that they only use tests that they write themselves. Of the nine developers using tools, six mentioned online composition aides such as [regex101.com](https://regex101.com) where a regex and input string are entered, and the input string is highlighted according to what is matched.

When asked an open ended question about pain points encountered with regular expressions, we observed three main categories. The most common, “hard to compose,” was represented in 61% (11) responses. Next, 39% (7) developers responded that regexes are “hard to read” and 17% (3) indicated difficulties with “inconsistency across implementations,” which manifest when using regexes in multiple languages. These responses do not sum to 18 as three developers provided multiple parts in their answers.

**Summary - RQ1:** Common uses of regexes include locating content within a file, capturing parts of strings, and parsing user input. The fact that all the surveyed developers compose regexes, and half of the developers use tools to test their regexes indicates the importance of tool development for regex. Developers complain about regexes being hard to read and hard to write.

Table 3.3 How saturated are projects with utilizations? (RQ2)

source	Q1	Avg	Med	Q3	Max
utilizations per project	2	32	5	19	1,427
files per project	2	53	6	21	5,963
utilizing files per project	1	11	2	6	541
utilizations per file	1	2	1	3	207

### 3.2.2 RQ2: How is the `re` module used?

We explore regex utilizations and flags used in the scraped Python projects. Out of the 3,898 projects scanned, 42.2% (1,645) contained at least one regex utilization. To illustrate how saturated projects are with regexes, we measure utilizations per project, files scanned per project, files contained utilizations, and utilizations per file, as shown in Table 3.3.

Of projects containing at least one utilization, the average utilizations per project was 32 and the maximum was 1,427. The project with the most utilizations is a C# project<sup>5</sup> that maintains a collection of source code for 20 Python libraries, including larger libraries like `pip`, `celery` and `ipython`. These larger Python libraries contain many utilizations. From Table 3.3, we also see that each project had an average of 11 files containing any utilization, and each of these files had an average of 2 utilizations.

The number of projects that use each of the `re` functions is shown in Figure 3.5. The y-axis denotes the total utilizations, with a maximum of 53,894. The `re.compile` function encompasses 57.6% of all utilizations. Note that compiled objects can also be used to call functions of the `re` module, ie `compiledObject.findall(...)`, but we ignore these utilizations so that our analysis is easier to automate, and because we are primarily interested in extracting the patterns which these 8 functions contain.

Of all utilizations, 87.3% had no flag, or explicitly specified the default flag. The debug flag, which causes the `re` regex engine to display extra information about its parsing, was never observed. This may be because developers use it for debugging and choose not to commit it to their repositories. **Summary - RQ2:** Only about half of the Python projects sampled

<sup>5</sup><https://github.com/Uuroboros/Arianrhod>

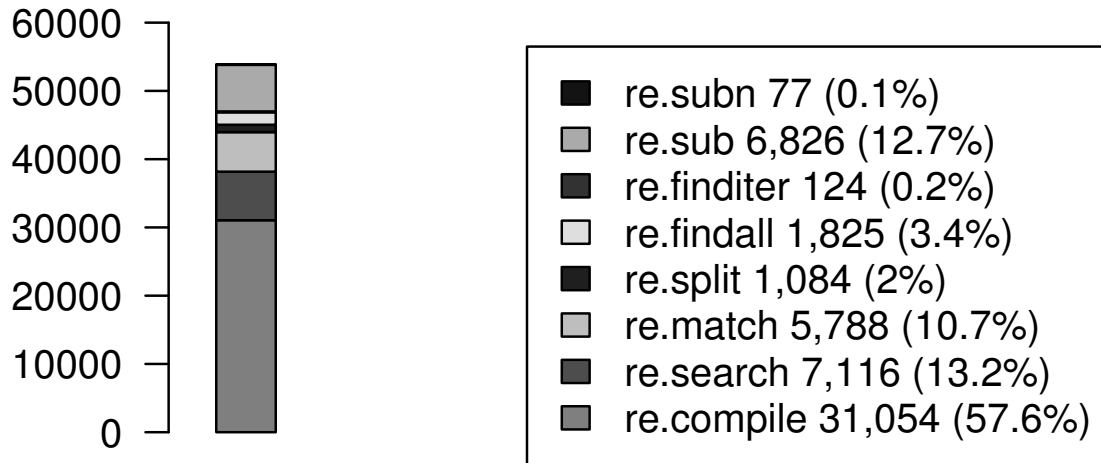
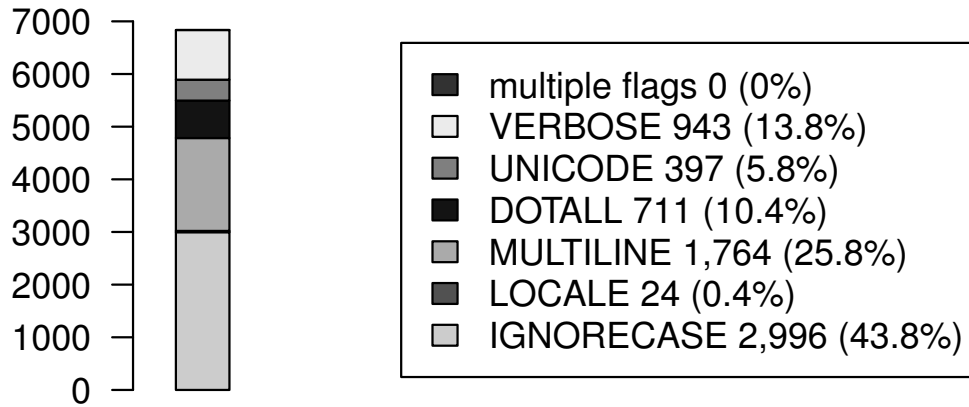
Figure 3.5 How often are `re` functions used? (RQ2)

Figure 3.6 Which behavioral flags are used? (RQ2)

contained any utilizations. Most utilizations used the `re.compile` function to compile a regex object before actually using the regex to find a match. Most utilizations did not use a flag to modify matching behavior.

### 3.2.3 RQ3: Regex language feature usage

We count the usages of each feature per project and as compared to all distinct regex patterns in the corpus.

### 3.2.3.1 Feature Usage

Table 3.4 displays feature usage from the corpus and relates it to four major regex related projects. Only features appearing in at least 10 projects are listed. The first column, *rank*, lists the rank of a feature (relative to other features) in terms of the number of projects in which it appears. The next column, *code*, gives a succinct reference string for the feature, and is followed by a *description* column that provides a brief comment on what the feature does. The *example* column provides a short example of how the feature can be used.

The next four columns, (i.e., *brics*, *hampi*, *Rex*, and *RE2*), map to the four major research projects chosen for our investigation (see Section 3.2.3.2). We indicate that a project supports a feature with the ‘●’ symbol, and indicate that a project does not support the feature with the ‘○’ symbol. The final four columns contain two pairs of usage statistics. The first pair contains the number and percent of *patterns* that a feature appears in, out of the 13,597 patterns that make up the corpus. The second pair of columns contain the number and percent of *projects* that a feature appears in out of the 1,645 projects scanned that contain at least one utilization.

One notable omission from Table 3.4 is the literal feature, which is used to specify matching any specific character. An example pattern that contains only one literal token is the pattern ‘a’. This pattern only matches the lowercase letter ‘a’. The literal feature was found in 97.7% of patterns.

We consider the literal feature to be necessary for any regex related tool to support, and so exclude it from Table 3.4 and the rest of the feature analysis.

The eight most commonly used features, ADD, CG, KLE, CCC, ANY, RNG, STR and END, appear in over half the projects. CG is more commonly used in patterns than the highest ranked feature (ADD) by a wide margin (over 8%), even though they appear in similar numbers of projects.

### 3.2.3.2 Feature Support in Regex Tools

While there are many regex tools available, in this work, we focus on the feature support for four tools, brics, hampi, Rex and RE2, which offer diversity across developers (i.e., Microsoft,

Google, open source, and academia) and applications. Further, as we wanted to perform a feature analysis, these four tools and their features are well-documented, allowing for easy comparison.

To create the tool mappings, we consulted documentation for each tool. For brics, we collected the set of supported features using the formal grammar<sup>6</sup>. For hampi, we manually inspected the set of regexes included in the `lib/regex-hampi/sampleRegex` file within the hampi repository<sup>7</sup> (this may have been an overestimation, as this included more features than specified by the formal grammar<sup>8</sup>). For RE2, we used the supported feature documentation<sup>9</sup>. For Rex, we collected the feature set empirically because we tried to parse all scraped patterns with Rex for the behavioral analysis (Section 3.2.4), and Rex provides comprehensive error feedback for unsupported features.

Of the four projects selected for this analysis, RE2 supports the most studied features (28 features) followed by hampi (25 features), Rex (21 features), and brics (12 features). All projects support the 8 most commonly used features except brics, which does not support STR or END.

No projects support the four look-around features LKA, NLKA, LKB and NLKB. RE2 and hampi support the LZY, NCG, PNG and OPT features, whereas brics and Rex do not.

### 3.2.3.3 Survey Results for Feature Usage

The pattern language for Python, which is used to compose regexes, supports default character classes like the ANY or dot character class: `.` meaning, ‘any character except newline’. It also supports three other default character classes: `\d`, `\w`, `\s` (and their negations). All of these default character classes can be simulated using the custom character class (CCC) feature, which can create semantically equivalent regexes. For example the decimal character class: `\d` is equivalent to a CCC containing all 10 digits: `\d ≡ [0123456789] ≡ [0-9]`.

---

<sup>6</sup><http://www.brics.dk/automaton/doc/index.html?dk/brics/automaton/RegExp.html>

<sup>7</sup><https://code.google.com/p/hampi/downloads/list>

<sup>8</sup><http://people.csail.mit.edu/akiezun/hampi/Grammar.html>

<sup>9</sup><https://re2.googlecode.com/hg/doc/syntax.html>

Other default character classes such as the word character class: `\w` may not be as intuitive to encode in a CCC: `[a-zA-Z0-9_]`.

Survey participants were asked if they use only CCC, use CCC more than default, use both equally, use default more than CCC or use only default. Results for this question are shown in Table 3.5, with 67% (12) indicating that they use default the most.

Participants who favored CCC indicated that “it is more explicit,” whereas the participants who favored default character classes said, “it is less verbose” and “I like using built-in code.”

To further explore how participants use various regex features, participants were asked five questions about how frequently they use specific related groups of features:

- endpoint anchors (STR, END): `^` and `$`
- capture groups(CG): (capture me)
- word boundaries (WNW): `word\b`
- (negative) look-ahead/behinds (LKA, NLKA, LKB, NLKB): `a(?!bc)`, `(?<x)yz!`, `(?<=a)`, `a(?yz)!`
- lazy repetition (LZY): `ab+?`, `xy{2,3}?`

These features were chosen based on the tool feature support explored in Section 3.2.3.2. Results are shown in Table 3.6, indicating that lazy repetition and look-ahead features are rarely used and capture groups and endpoint anchors are occasionally to frequently used.

**Summary - RQ3:** The eight most common features are found in over 50% of the projects. Shown in Table 3.4, the STR and END features are present in over half of the scanned projects containing utilizations. In our survey, over half (56%) of the respondents answered that they use endpoint anchors frequently or very frequently, and none of them claimed to never use them.

The LZY feature is present in over 36% of scanned projects with utilizations, and yet was not supported by two of the four major regex projects we explored, brics and RE2. In our developer survey, 11% (2) of participants use this feature frequently and 6 (33%) use it occasionally, showing a modest impact on potential users.

When survey participants were asked if they prefer to always use numbered (BKR) or named (BKRN) back references, 66% (12) of survey participants said that they always use BKR, and the remaining 33% (6) said “it depends.” No participants preferred named capture groups. BKR is present in 5% of scanned projects, while BKRN is present in only 1.7%, which corroborates our findings that numbered are generally preferred over named capture groups.

### 3.2.4 RQ4: Regex behavioral similarity

In clustering the regular expressions, we are most interested in observing behavior of regexes found in multiple projects. Starting with the 13,597 patterns of the corpus, we discarded 10,015 (74%) patterns that were not found in multiple projects. Then we excluded an additional 711 (5%) patterns that contain features not supported by Rex. We studied the remaining 2,871 (21%) patterns using our similarity analysis technique. The impact is that 923 projects were excluded from the data set for the similarity analysis. Omitted features are indicated in Table 3.4 for Rex.

From 2,871 distinct patterns, MCL clustering identified 186 clusters with 2 or more patterns, and 2,042 clusters of size 1. The average size of clusters larger than size one was 4.5. Each pattern belongs to exactly one cluster.

Three example strings generated by Rex for the first pattern are: ‘-()’, ‘\*’8(5)’, ‘Oe()’. For the third pattern, Rex generated these three strings: ‘()’, ‘(q)F’, ‘(n)M’. The pattern: `\(.*\)\$` is very similar, but will not match the string ‘(n)M’, and so was placed in a different cluster.

Table 3.7 provides an example of a behavioral cluster containing 12 patterns (four longer patterns omitted for brevity). Patterns from this cluster are present in 31 different projects. All patterns in this cluster share the literal ‘:’ character. The smallest pattern, ‘:+', matches one or more colons.

Another pattern from this cluster, `([^\:]+):(.*)`, requires at least one non-colon character to occur before a colon character. Our similarity value between these two regexes was below the minimum of 0.75 because Rex generated many strings for ‘:+' that start with one or more colons. We observe that the smallest pattern in a cluster provides insight about key characteristic that all the patterns in the cluster have in common. A shorter pattern will tend

to have less extraneous behavior because it is specifying less behavior, yet, in order for the smallest pattern to be clustered, it had to match most of the strings created by Rex from many other patterns within the cluster, and so we observe that the smallest pattern is useful as a representative of the cluster.

For the rest of this paper, a cluster will be represented by one of the shortest patterns it contains, followed by the number of projects any member of the cluster appears in, so the cluster in Table 3.7 will be represented as ‘:+(31)’. This representation is not an attempt to express all notable behavior of patterns within a cluster, but is a useful and meaningful abbreviation. Other regexes in the cluster may exhibit more diverse behavior, for example the pattern ‘([^\s: ]+):(.\*)’ requires a non-colon character to appear before a colon character.

We manually mapped the top 100 largest clusters based on the number of projects into 6 behavioral categories (determined by inspection). The largest cluster was left out, as it was composed of patterns that trivially matched almost any string, like ‘b\*’ and ‘^’. The remaining 99 clusters were all categorized. These clusters are briefly summarized in Table 3.8, showing the name of the category and the number of clusters it represents, patterns in those clusters, and projects. The most common category is *Multi Matches*, which contains clusters that have alternate behaviors (e.g., matching a comma or a semicolon, as in ‘,|;(18)’). Each cluster was mapped to exactly one category. Next, we describe the categories, ordered by the number of projects the regex patterns map to.

#### 3.2.4.1 Multiple Matching Alternatives

The patterns in these clusters match under a variety of conditions by using a character class or a disjunctive |. For example: ‘(\W)’ (89) matches any alphanumeric character, ‘(\s)’ (89) matches any whitespace character, ‘\d’ (58) matches any numeric character, and ‘,|;(18) matches a comma or semicolon. Most of these clusters are represented by patterns that use default character classes, as opposed to custom character classes. This provides further support for our survey results to the question, *Do you prefer to use custom character classes or default character classes more often?*, in which a majority of participants indicated they use the default



classes more than custom. This category contains 21 clusters, each appearing in an average of 33 projects.

#### 3.2.4.2 Specific Character Must Match

Each cluster in this category requires one specific character to match, for example: `'\n\s*'` (42) matches only if a newline is found, `':+'` (31) matches only if a colon is found, `'%'` (22), matches only if a percent sign is found and `'}'` (14) matches only if a right curly brace is found. Table 3.7 presents a cluster that falls under this category. While the cluster is centered on the presence of the `:` character, the other regexes in the cluster also exhibit more diverse behavior. The commonality of this cluster category contrasts with the survey in Section 3.2.1 in which participants reported to very rarely or never use regexes to check for a single character (Table 3.2). This category contains 17 clusters, each appearing in an average of 17.1 projects. These clusters have a combined total of 103 patterns, with at least one pattern present in 184 projects.

#### 3.2.4.3 Anchored Patterns

Each of the clusters uses at least one endpoint anchor to require matches to be absolutely positioned, for example: `'(\w+)$'` (35) captures the word characters at the end of the input, `'^\s'` (16) matches a whitespace at the beginning of the input, and `'^-?\d+$'` (17) requires that the entire input is an (optionally negative) integer. These anchors are the only way in regexes to guarantee that a character does (or does not) appear at a particular location by specifying what is allowed. As an example, `^[-_A-Za-z0-9]+$` says that from beginning to end, only `[-_A-Za-z0-9]` characters are allowed, so it will fail to match if undesirable characters, such as `?`, appear anywhere in the string. This category contains 20 clusters, each appearing in an average of 15.4 projects. These clusters have a combined total of 85 patterns, with at least one pattern present in 141 projects.

**TODO.MID: The thing I want to mention about anchored patterns (but have struggled to say in the past) is that they are the only way to guarantee that a character does not appear in a particular location by specifying what is allowed.**

Consider the regex `^[-_A-Za-z0-9]+$` **TODO.MID:** which will fail to match if an undesirable character like ‘?’ appears anywhere in the input. In logic, there is a similar phenomenon. That is, ‘Always’ is true iff ‘Not Exists’ of the negation is true, and by requiring an entire input to always maintain some abstraction, you can indirectly specify the negation of another (inverse) abstraction. Even with only one anchor point, a regex like `.*[0-9]$` **TODO.MID:** is creating an ultimatum about the end being a digit. Without the endpoint anchors, I don’t see how one could specify absolutes about an input.

#### 3.2.4.4 Content of Brackets and Parenthesis

The clusters in this category center around finding a pair of characters that surround content, often also capturing that content. For example, `\(.*\)` (29) matches when content is surrounded by parentheses and `".*"` (25) matches when content is surrounded by double quotes. The cluster `<(.+)>` (23) matches and captures content surrounded by angled brackets. This category contains 10 clusters, each appearing in an average of 18.4 projects. These clusters have a combined total of 46 patterns, with at least one pattern present in 111 projects. **TODO.MID: include this?, and** `\[.*\]` (22) **TODO.MID: matches when content is surrounded by square brackets**

#### 3.2.4.5 Two or More Characters in Sequence

These clusters require several characters in a row to match some pattern, for example: `\d+\.\d+` (30) requires one or more digits followed by a period character, followed by one or more digits. The cluster  (17) requires two spaces in a row, `([A-Z][a-z]+[A-Z][^ ]+)` (11), and `@[a-z]+'` (9) requires the at symbol followed by two or more lowercase characters, as in a twitter handle. This category contains 16 clusters, each appearing in an average of 13 projects. These clusters have a combined total of 40 patterns, with at least one pattern present in 120 projects.

**TODO.MID:** Again, it might be interesting to look at what particular sequences are looking like. I think I mention this again in the discussion, but should we put it here instead?

#### 3.2.4.6 Code Search and Variable Capturing

These clusters show a recognizable effort to parse source code or URLs. For example, ‘`^https?:/`’(23) matches a web address, and ‘`(.+)=(.+)`’(9) matches an assignment statement, capturing both the variable name and value. The cluster ‘`\$\{([\w\-\+])\}`’(11) matches an evaluated string interpolation and captures the code to evaluate. This category contains 15 clusters, each appearing in an average of 11.7 projects. These clusters have a combined total of 27 patterns, with at least one pattern present in 92 projects.

**Summary - RQ4:** When tool designers are considering what features to include, data about usage in practice is valuable. Behavioral similarity clustering helps to discern these behaviors by looking beyond the structural details of specific patterns and seeing trends in matching behavior. We are also able to find out what features are being used in these behavioral trends so that we can make assertions about why certain features are important. We used the behavior of individual patterns to form clusters, and identified six main categories for the clusters. Overall, we see that many clusters are defined by the presence of particular tokens, such as the colon for the cluster in Table 3.7. We identified six main categories that define regex behavior at a high level: matching with alternatives, matching literal characters, matching with sequences, matching with endpoint anchors, parsing contents of brackets or braces, or searching and capturing code, and can be considered in conjunction with the self-described regex activities from the survey in Table 3.2 to be representative of common uses for regexes. One of the six common cluster categories, *Code Search and Variable Capturing*, has a very specific purpose of parsing source code files. This shows a very specific and common use of regular expressions in practice.

Table 3.4 How frequently do features appear in projects?

rank	code	description	example	brics	hampi	Rex	RE2	nPatterns	% patterns	nProjects
1	ADD	one-or-more repetition	<code>z+</code>	●	●	●	●	6,003	44.1	1,204
2	CG	a capture group	<code>(caught)</code>	●	●	●	●	7,130	52.4	1,194
3	KLE	zero-or-more repetition	<code>.*</code>	●	●	●	●	6,017	44.3	1,099
4	CCC	custom character class	<code>[aeiou]</code>	●	●	●	●	4,468	32.9	1,026
5	ANY	any non-newline char	<code>.</code>	●	●	●	●	4,657	34.3	1,005
6	RNG	chars within a range	<code>[a-z]</code>	●	●	●	●	2,631	19.3	848
7	STR	start-of-line	<code>^</code>	○	●	●	●	3,563	26.2	846
8	END	end-of-line	<code>\$</code>	○	●	●	●	3,169	23.3	827
9	NCCC	negated CCC	<code>[^qwxzf]</code>	●	●	●	●	1,935	14.2	776
10	WSP	<code>\t \n \r \v \f</code> or space	<code>\s</code>	○	●	●	●	2,846	20.9	762
11	OR	logical or	<code>a b</code>	●	●	●	●	2,102	15.5	708
12	DEC	any of: 0123456789	<code>\d</code>	○	●	●	●	2,297	16.9	692
13	WRD	<code>[a-zA-Z0-9_]</code>	<code>\w</code>	○	●	●	●	1,430	10.5	650
14	QST	zero-or-one repetition	<code>z?</code>	●	●	●	●	1,871	13.8	645
15	LZY	as few reps as possible	<code>z+?</code>	○	●	○	●	1,300	9.6	605
16	NCG	group without capturing	<code>a(?:b)c</code>	○	●	○	●	791	5.8	404
17	PNG	named capture group	<code>(?P&lt;name&gt;x)○</code>	○	●	○	●	915	6.7	354
18	SNG	exactly n repetition	<code>z{8}</code>	●	●	●	●	581	4.3	340
19	NWSP	any non-whitespace	<code>\S</code>	○	●	●	●	484	3.6	270
20	DBB	$n \leq x \leq m$ repetition	<code>z{3,8}</code>	●	●	●	●	367	2.7	238
21	NLKA	sequence doesn't follow	<code>a(?:!yz)</code>	○	○	○	○	131	1	183
22	WNW	word/non-word boundary	<code>\b</code>	○	○	○	●	248	1.8	166
23	NWRD	non-word chars	<code>\W</code>	○	●	●	●	94	0.7	165
24	LWB	at least n repetition	<code>z{15,}</code>	●	●	●	●	91	0.7	158
25	LKA	matching sequence follows	<code>a(?=bc)</code>	○	○	○	○	112	0.8	158
26	OPT	options wrapper	<code>(?i)CasE</code>	○	●	○	●	231	1.7	154
27	NLKB	sequence doesn't precede	<code>(?&lt;!x)yz</code>	○	○	○	○	94	0.7	137
28	LKB	matching sequence precedes	<code>(?&lt;=a)bc</code>	○	○	○	○	80	0.6	120
29	ENDZ	absolute end of string	<code>\Z</code>	○	○	○	●	89	0.7	90
30	BKR	match the $i^{th}$ CG	<code>\1</code>	○	○	○	○	60	0.4	84
31	NDEC	any non-decimal	<code>\D</code>	○	●	●	●	36	0.3	58
32	BKRN	references PNG	<code>\g&lt;name&gt;</code>	○	●	○	○	17	0.1	28
33	VWSP	matches U+000B	<code>\v</code>	○	○	●	●	13	0.1	15
34	NWNW	negated WNW	<code>\B</code>	○	○	○	●	4	0	11

Table 3.5 Survey results for preferences between custom character and default character classes (RQ3)

Preference	Frequency
use only CCC	1
use CCC more than default	5
use both equally	2
use default more than CCC	10
use only default	2

Table 3.6 Survey results for regex usage frequencies, averaged using a 6-point likert scale: Very Frequently=6, Frequently=5, Occasionally=4, Rarely=3, Very Rarely=2, and Never=1 (RQ3)

Group	Code	Frequency
endpoint anchors	(STR, END)	4.4
capture groups	(CG)	4.2
word boundaries	(WNW)	3.5
lazy repetition	(LZY)	2.9
(neg) look-ahead/behind	(LKA, NLKA, LKB, NLKB)	2.5

Table 3.7 Sample from an example cluster (RQ4)

index	pattern	nProjects	index	pattern	nProjects
1	'[:+']	8	5	'[::]'	6
2	'(:)'	8	6	'([[:+]]+:.*)'	6
3	'(:+)'	8	7	'\s*:\s*'	4
4	'(:)(:*)'	8	8	'\:'	2

Table 3.8 Cluster categories and sizes (RQ4)

<b>Category</b>	<b>Clusters</b>	<b>Patterns</b>	<b>Projects</b>
Multi Matches	21	237	295
Specific Char	17	103	184
Anchored Patterns	20	85	141
Content of Parens	10	46	111
Two or More Chars	16	40	120
Code Search	15	27	92

## CHAPTER 4. Conclusion

### 4.1 Discussion From Features

In this section, we discuss the implications of these empirical findings and opportunities for future work.

#### 4.1.1 Implications For Tool Designers

The results have implications for regex tool designers.

##### 4.1.1.1 Finding Specific Content

Two categorical clusters, *Specific Characters Must Match* (Section 3.2.4.2) and *Two or More Characters in Sequence* (Section 3.2.4.5), deal with identifying the presence of specific character(s). While multiple character matching subsumes single character matching, the overarching theme is that these regexes are looking to validate strings based on the presence of very specific content, as would be done for many common activities listed in Table 3.2, such as, “Locating content within a file or files.” More study is needed into what content is most frequently searched for, but from our cluster analysis we found that version numbers, twitter or user handles, hex values, decimal numbers, capitalized words, and particular combinations of whitespace, slashes and other delimiters were discernible targets.

##### 4.1.1.2 Capturing Specific Content Near A Delimiter

The survey results from Section 3.2.1 indicate that capturing parts of strings is among the most frequent activities for which developers use regexes. From a feature perspective, the capture group (CG) is the most frequently used in terms of patterns (Table 3.4). This feature

has two functions: 1) logical grouping as would be expected by parenthesis, and 2) retrieval of information in one logical grouping. As mentioned in Section 3.2.4, capturing content was a primary goal evident in several cluster categories. The fourth-largest category is based entirely on capturing the content between brackets or parentheses (Section 3.2.4.4).

Many uses of CG also use the ANY and KLE features, eg. `(.*){(.*)}(.*)` and `\\s*([~: ]*)\\s*:(.*)`. This type of usage frequently revolves around an important delimiter character such as `:` or `\`. This use case is well supported by existing tools for ASCII characters, but future tools should consider the centrality of this use case and its implications for non-English users of regex tools. For example, Unicode characters like ‘U+060D’ the Arabic Date Separator, or ‘U+1806’ the Mongolian Todo Soft Hyphen may be used to locate segments of text that a user would want to capture.

#### 4.1.1.3 Counting Lines

Text files containing one unit of information per line are common in a wide variety of applications (for example .log and .csv files). Out of the 13,597 patterns in the corpus, 3,410 (25%) contained ANY followed by KLE (i.e., ‘.\*’), often at the end of the pattern. One reasonable explanation for this tendency to put ‘.\*’ at the end of a pattern is that users want to disregard all matches after the first match on a single line in order to count how many distinct lines the match occurs on. Survey participants indicated an average frequency of “Counting lines that match a pattern” and “Counting substrings that match a pattern” at 3.2 or rarely/occasionally. It may be valuable for tool builders to include support for common activities such as line counting.

### 4.1.2 Opportunities For Future Work

There are many opportunities for future work.

#### 4.1.2.1 Refactoring Regexes

The survey showed that users want readability and find the lack of readable regexes to be a major pain point. This provides an opportunity to introduce refactoring transformations



to enhance readability or comprehension. As one opportunity, certain character classes are logically equivalent and can be expressed differently, for example, `\d`  $\equiv$  `[0123456789]`  $\equiv$  `[0-9]`. While `\d` is more succinct, `[0-9]` may be easier to read, so a refactoring for *default to custom character classes* could be introduced. Human studies are needed to evaluate the readability and comprehension of various regex features in order to define and support appropriate regex refactorings.

Another avenue of refactoring could be for performance. Various implementations of regex libraries may perform more efficiently with some features than others. An evaluation of regex feature implementation speeds would facilitate semantic transformations based on performance, similar to performance refactorings for LabVIEW Chambers and Scaffidi (2013, 2015).

Additionally, some developers may *find* specific content with a regex and then subsequently *capture* it with string parsing, which may be more error prone than using a capture group and indicates a missed opportunity to use the full extent of regex libraries. Future work will explore source code to identify the frequency of such occurrences and design refactorings to better utilize regex library features.

#### 4.1.2.2 Migration Support for Developers

Within standard programming languages, regular expressions libraries are very common, yet there are subtle differences between language libraries in the supported features. For example, Java supports possessive quantifiers like `'ab*+c'` (here the `'+'` is modifying the `'*'` to make it possessive) whereas Python does not. Differences among programming language implementations was identified as a pain point for using regular expressions by 17% of the survey participants. This provides a future opportunity for tools that translate between regex utilizations in various languages.

#### 4.1.2.3 Similarity Beyond String Matching

There are various ways to compute similarity between regexes, each with different trade-offs. While the similarity analysis we employ over-approximates similarity when compared to containment analysis, it may under-approximate similarity in another sense.

For example, two regexes that have dissimilar matching behavior could be very similar in purpose and in the eyes of the developer. For example, `commit:\[(\d+)\] - (.*)` and `push:\[(\d+)\] - (.*)` could both be used to capture the id and command from a versioning system, but match very different sets of strings. Future work would apply abstractions to the regex strings, such as removing or relaxing literals, prior to similarity analysis to capture and cluster such similarities.

From another perspective, our regex similarity measure, and even containment analysis, could treat behaviorally identical regexes as the same, when their usage in practice is completely different. For example, in Table 3.7, the regexes `:+` and `(:+)` are behaviorally identical in that they match the same strings, except the latter uses a capture group. In practice, these may be used very differently, where the former may be used for validation and the latter for extraction. This usage difference could be observed by code analysis, and is left for future work.

#### 4.1.2.4 Automated Regex Repair

Regular expression errors are common and have produced thousands of bug reports Spishak et al. (2012). This provides an opportunity to introduce automated repair techniques for regular expressions. Recent approaches to automated program repair rely on mutation operators to make small changes to source code and then re-run the test suite (e.g., Weimer et al. (2010); Le Goues et al. (2012)). In regular expressions, it is likely that the broken regex is close, semantically, to the desired regex. Syntax changes through mutation operators could lead to big changes in behavior, so we hypothesize that using the semantic clusters identified in Section 3.2.4 to identify potential repair candidates could efficiently and effectively converge on a repair candidate.

#### 4.1.2.5 Developer Awareness of Best Practices

One category of clusters, *Content of Brackets and Parenthesis*, parses the contents of angle brackets, which may indicate developers are using regexes to parse HTML or XML. As the contents of angle brackets are usually unconstrained, regexes are a poor replacement for XML or HTML parsers. This may be a missed opportunity for the regex users to take advantage of

more robust tools. More research is needed into how regex users discover best practices and how aware they are of how regexes should and should not be used.

#### 4.1.2.6 Tool-Specific Regex Exploration

In some environments, such as command line or text editor, regexes are used extensively by the surveyed developers (Section 3.2.1), but these regular expressions do not persist. Thus, using a repository analysis for feature usage only illustrates part of how regexes are used in practice. Exploring how the feature usage differs between environments would help inform tool developers about how to best support regex usage in context, and is left for future work.

## 4.2 Discussion From Refactoring

Based on our analyses of source code and our empirical study on the understandability of regex representations, we have identified preferred regex representations that may make regexes easier to understand and thus maintain. In this section, we describe the implications of these results.

### 4.2.1 Interpreting Results

In the CCC equivalence class, C1 (e.g., `[0-9a]`) is more commonly found in the patterns and projects. Representations C2 (e.g., `[0123456789a]`) and C3 (e.g., `[\x00-/:-‘b-\x7F]`) appear in similar percentages of patterns and projects but there is no significant difference in understandability considering two pairs of regexes tested as part of E13 (Table ??). However, a small preference is shown for C1 over C2 (E7), leading this to be the winner of both the community support and understandability analyses. Regex length is probably important for understandability, though we did not test for this.

**TODO.NOW: the longest regex in the corpus is X characters long...**

In the DBB group, D3 (e.g., `pBs|pBBs|pBBBs`) merits further exploration because it is the most understandable but least common node in DBB group. This may be because explicitly listing the possibilities with an OR is easy to grasp, but if the number of items in the OR is too large, the understandability may go down. Further analysis is needed to determine the

optimal thresholds for representing a regex as D3 compared to D1 (e.g., `pB{1,3}s`) or D2 (e.g., `pBB?B?s`).

**TODO.LAST:** Intuitively, it seems that D2 may be more common because 0,1 is just a more common use case than an arbitrary range like 4, 25.

In the SNG group, S1 is a compact representation (e.g., `S{3}`), but S2 was preferred (e.g., `SSS`). Similar to the DBB group, this may be do to the particular examples chosen in the analysis, as a large number of explicit repetitions may not be as preferred.

In the LWB group, L1 (e.g., `A{2,}`) is rare, appearing in  $< 1\%$  of the patterns. Representations L2 (e.g., `AAA*`) and L3 (e.g., `AA+`) appear in similar numbers of patterns and projects, but there is a significant difference in their understandability, favoring L3. **TODO.LAST:** it's clear that this is a rare use case, and also that L3 is the most common use case. Patterns using star are secondary, helper patterns because they will trivially match anything, so they are less common. But anyway...

**TODO.LAST:** S2 is over-weighted because of double-characters in regular words like `foot`. In the LIT group, T1 (e.g., `\a\>`) is the typical way to list literals, but the reason to use hex (T2) or oct (T4) types is because some characters cannot be represented any other way, like invisible chars. One main result of our work is that T4 (e.g., `\007\036\062`) is less understandable than T2 (e.g., `\x07\x24\x3E`), so if invisible chars are required, hex is the more understandable representation. Regarding T3 (e.g., `\a[$]>`), initially we thought the square brackets would be more understandable than using an escape character, but we found the opposite. Given a choice between T1 and T3, the escape character was more understandable.

#### 4.2.2 Opportunities For Future Work

There are several directions for future work related to regex study and refactoring.

**Equivalence Class Models** We looked at five equivalence classes, each with three to five nodes. Future work could consider richer models with more or different classes and nodes. **TODO.LAST:** For example, we have looked at all ranges as equivalent, all defaults as equivalent, and relied on many such generalizations. However, the range `[a-f]`

**TODO.LAST:** is likely to be more understandable for most people than a range like `[:-' ]`.

Additional equivalence groups to consider may include:

**Single line option** `'''(.|\n)+'''`  $\equiv$  `(?s)'''(.)+'''`

**Multi line option** `(?m)G\n`  $\equiv$  `(?m)G$`

**Case insensitive** `(?i)[a-z]`  $\equiv$  `[A-Za-z]`

**Backreferences** `(X)q\1`  $\equiv$  `(?P<name>X)q\g<name>`

**Word Boundaries** `\bZ`  $\equiv$  `((?<=\w)(?=\W)|(?<=\W)(?=\w))Z`

It might also be the case that there exist critical comprehension differences within a representation. For example, between C1 (e.g., `[0-9a]`) and C4 (e.g., `[\da]`), it could be the case that `[0-9]` is preferred to `[\d]`, but `[A-Za-z0-9_]` is not be preferred to `[\w]`). By creating a more granular model of equivalence classes, and making sure to carefully evaluate alternative representations of the most frequently used specific patterns, additional useful refactorings could be identified.

**TODO.LAST:** One of the most straightforward ways to address understandability is to directly ask software professionals which from a list of equivalent regexes they prefer and why. If understandability measurements used regexes sampled from the codebase of a specific community (most frequently observed regexes, most buggy regexes, regexes on the hottest execution paths, etc.), and measured the understanding of programming professionals working in that community, then the measurements and the refactorings they imply would be more likely to have a direct and certain positive impact. In another study, we did a survey where software professionals indicated that understandability of regexes they find in source code is a major pain point. In this study, our participants indicated that they read about twice as many regexes as they compose. What is the impact on maintainers, developers and contributors to open-source projects of not being able to understand a regex that they find in the code they are working with? Presumably this is a frustrating experience - how much does a confusing regex slow down a

software professional? What bugs or other negative factors can be attributed to or associated with regexes that are difficult to understand? How often does this happen and in what settings? Future work could tailor an in-depth exploration of the overall costs of confusing regexes and the potential benefits of refactoring or other treatments for confusing regexes.

**Regex Migration Libraries** We have identified opportunities to improve the understandability of regexes in existing code bases by looking for some of the less understandable regex representations, which can be thought of as antipatterns, and refactoring to the more common or understandable representations. Building migration libraries is a promising direction of future work to ease the manual burden of this process, similar in spirit to prior work on class library migration Balaban et al. (2005).

**Regex Refactoring Applications** Maintainers of code that is intentionally obfuscated for security purposes may want to develop regexes that they understand and then automatically transform them into the least understandable regex possible.

One fundamental concept that many users of regex struggle to learn is when to use regexes for simple parsing, and when to write a full-fledged parser (for example, when parsing HTML). Regexes that are trying to parse HTML, XML or similar languages could be refactored not into a better regex, but into some code with an equivalent intention that does parsing much better.

**Regex Programming Standards** Many organizations enforce coding standards in their repositories to ease understandability. Presently, we are not aware of coding standards for regular expressions, but this work suggests that enforcing standard representations for various regex constructs could ease comprehension.

**Regex Refactoring for Performance** The representation of regexes may have a strong impact on the runtime performance of a chosen regex engine. Prior work has sought to expedite the processing of regexes over large bodies of text Baeza-Yates and Gonnet (1996). Refactoring

regexes for performance would complement those efforts. Further study is needed to determine which representations are most efficient, leading to a whole new area of study on regex refactoring for performance, a topic already explored for Depending on the efficiency of an organization’s chosen regex engine, an organization may want to enforce standards for efficiency, , or for compatibility with a regex analysis tool like Z3, HAMPI, BRICS or REX.

#### 4.2.3 Threats to Validity

**Internal** We measure understandability of regexes using two metrics, matching and composition. However, these measures may not reflect actual understanding of the regex behavior. For this reason, we chose to use two metrics and present the analysis in the context of reading and writing regexes, but the threat remains.

Participants evaluated regular expressions during tasks on MTurk, which may not be representative enough of the context in which programmers would encounter regexes in practice. Further study is needed to determine the impact of the experimentation context on the results.

Some regex representations from the equivalence classes were not involved in the understandability analysis and that may have biased the results against those nodes. Repetition of the analysis with more complete coverage of the edges in the equivalence classes is needed.

We treated unsure responses as omissions that did not count against the matching scores. Thus, if a participant answered two strings correctly and marked the other three strings as unsure, then this was 2/2 correct, not 2/5. This may have inflated the matching scores, however, less than 5% of the matching scores were impacted by such responses.

**TODO.LAST:** In our analyses, we measure understandability using matching and composition metrics. However, there may be other ways to approach regex understandability, such as deciding which regexes in a set are equivalent, finding the minimum modification to some text so that a given regex will match it. It may also be meaningful to provide some code that exists around a regex as context, since that would better represent a scenario in which programmers would encounter regexes in practice. Further study is needed to determine if the cho-

sen metrics and experimentation context have resulted in a reasonable measure of understandability.

**External** Participants in our survey came from MTurk, which may not be representative of people who read and write regexes on a regular basis.

The regexes used in the evaluation were inspired by those found in Python code, which is just one language that has library support for regexes. Thus, we may have missed opportunities for other refactorings based on how programmers use regexes in other programming languages.

The results of the understandability analysis may be closely tied to the particular regexes chosen for the experiment. For many of the representations, we had several comparisons. Still, replication with more regex patterns is needed.

**TODO.MID:** what about the threat of too few examples per node? Didn't cover every edge. Regex set is randomly collected online, not focused on any specific target audience.

TODO.LAST: Our community analysis only focuses on the Python language, but as the vast majority of regex features are shared across most general programming languages (e.g., Java, C, C#, or Ruby), a Python pattern will (almost always) behave the same when used in other languages and our results are likely to generalize. , whereas a utilization is not universal in the same way (i.e., it may not compile in other languages, even with small modifications to function and flag names). As an example, the `re.MULTILINE` flag, or similar, is present in Python, Java, and C#, but the Python `re.DOTALL` flag is not present in C# though it has an equivalent flag in Java.

**TODO.MID:** Looks like M0, M1, M2, M3 and M9 are very dependent on the regex chosen, so regex-specific refactorings like:

0.1401 &d([aeiou][aeiou])z', &d([aeiou]{2})z', 0.075 [\t\r\f\n ], [\s]', 0.1024  
[a-f]([0-9]+)[a-f]', [a-f](\d+)[a-f]', 0.1271 [\{ \} \[ \] (\d+.[\d]) [\] ]', \\{\\\$(\d+.[\d])\\}'

**TODO.MID:** (from M0,M1,M2,M9 respectively) have okay P-values and may indicate regex-specific refactorings, but do not indicate an overall trend for that type of refactoring. Notice that M3 does not even have a strong p-value candidate,



but this may be thrown off because of the very confusing regex chosen for CCC: **0.78 0.79** `xyz[_\[\]‘\^\]\’ xyz[\x5b-\x5f]\’` **TODO.MID:** which has a lot of escape characters, so that the hex group was easier to understand than the CCC.

**TODO.MID:** Meanwhile M4,M5 and M7 have both ambiguous p-values and anova results. But this is still a finding: that no refactoring is needed between things like: `(q4fab|ab)\’ ((q4f){0,1}ab)\’ tri[abcdef]3\’ tri(a|b|c|d|e|f)3\’ &(\w+);\’ &([A-Za-z0-9_]+);\’` **TODO.MID: (from M4,M5,M7 respectively)** Although one refactoring from M5 might be of slight interest: `0.1196 FALSE tri[a-f]3\’ tri(a|b|c|d|e|f)3\’`

### 4.3 Conclusion

In an effort to find refactorings that improve the understandability of regexes, we created five equivalence class models and used these models to investigate the most common representations and most comprehensible representations per class. We found the most common representations per class by both number of patterns and number of projects to be C1, D2, T1 and S2 (L3 has the most patterns, L2 has the most projects). We also identified three strongly preferred transformations between representations (i.e.,  $\overrightarrow{T4T1}$ ,  $\overrightarrow{D2D3}$ , and  $\overrightarrow{L2L3}$ ) according to the results of our comprehension tests. We combined the results of these two investigations using a version of Kahn’s topological sorting algorithm to produce a total ordering of representations within each model. The agreement between Community Standards and Understandability in this analysis validates our results and suggests that indeed one particular representation can be preferred over others in most cases. We can also recommend using hex to represent invisible characters in regexes instead of octal, and to escape special characters with slashes instead of wrapping them in brackets to avoid escaping them. Further research is needed into more granular models that treat common specific cases separately, and that address the effect of length on readability when transforming from one representation to another.

## **APPENDIX A. Patterns in Python projects from Github**

This is now the same as any other chapter except that all sectioning levels below the chapter level must begin with the \*-form of a sectioning command.

**top 100 clusters used**

**top 10 by feature group**

**13,579 patterns: the corpus(1510 per page=9 pages)**

Supplemental material.

## **APPENDIX B. Developer Survey**

This is now the same as any other chapter except that all sectioning levels below the chapter level must begin with the \*-form of a sectioning command.

### **Survey Questions**

### **Survey Responses**

### **Survey Statistics**

More stuff.

## **APPENDIX C. Mechanical Turk Study**

This is now the same as any other chapter except that all sectioning levels below the chapter level must begin with the \*-form of a sectioning command.

### **Qualifying Test**

#### **Template**

#### **MT Data input**

#### **MT All Results**

#### **MT Summary Statistics**

More stuff.

## **APPENDIX D. Community Analysis**

This is now the same as any other chapter except that all sectioning levels below the chapter level must begin with the \*-form of a sectioning command.

### **Filter Criteria**

### **Summary Statistics**

More stuff.

## BIBLIOGRAPHY

- Abbes, M., Khomh, F., Gueheneuc, Y.-G., and Antoniol, G. (2011). An empirical study of the impact of two antipatterns, blob and spaghetti code, on program comprehension. In *Software Maintenance and Reengineering (CSMR), 2011 15th European Conference on*, pages 181–190. IEEE.
- Alkhateeb, F., Baget, J.-F., and Euzenat, J. (2009). Extending sparql with regular expression patterns (for querying rdf). *Web Semant.*, 7(2):57–73.
- Anand, S., Burke, E. K., Chen, T. Y., Clark, J., Cohen, M. B., Grieskamp, W., Harman, M., Harrold, M. J., and Mcminn, P. (2013). An orchestrated survey of methodologies for automated software test case generation. *J. Syst. Softw.*, 86(8):1978–2001.
- Arslan, A. (2005). Multiple sequence alignment containing a sequence of regular expressions. In *Computational Intelligence in Bioinformatics and Computational Biology, 2005. CIBCB '05. Proceedings of the 2005 IEEE Symposium on*, pages 1–7.
- Babbar, R. and Singh, N. (2010). Clustering based approach to learning regular expressions over large alphabet for noisy unstructured text. In *Proceedings of the Fourth Workshop on Analytics for Noisy Unstructured Text Data, AND '10*, pages 43–50, New York, NY, USA. ACM.
- Baeza-Yates, R. A. and Gonnet, G. H. (1996). Fast text searching for regular expressions or automaton searching on tries. *J. ACM*, 43(6):915–936.
- Balaban, I., Tip, F., and Fuhrer, R. (2005). Refactoring support for class library migration. *SIGPLAN Not.*, 40(10):265–279.

- Beck, F., Gulan, S., Biegel, B., Baltes, S., and Weiskopf, D. (2014). Regviz: Visual debugging of regular expressions. In *Companion Proceedings of the 36th International Conference on Software Engineering*, ICSE Companion 2014, pages 504–507, New York, NY, USA. ACM.
- Begel, A., Khoo, Y. P., and Zimmermann, T. (2010). Codebook: Discovering and exploiting relationships in software repositories. In *Proceedings of the 32Nd ACM/IEEE International Conference on Software Engineering - Volume 1*, ICSE '10, pages 125–134, New York, NY, USA. ACM.
- Callaú, O., Robbes, R., Tanter, E., and Röthlisberger, D. (2011). How developers use the dynamic features of programming languages: The case of smalltalk. In *Proceedings of the 8th Working Conference on Mining Software Repositories*, MSR '11, pages 23–32, New York, NY, USA. ACM.
- Callaú, O., Robbes, R., Tanter, E., and Röthlisberger, D. (2013). How (and why) developers use the dynamic features of programming languages: The case of smalltalk. *Empirical Software Engineering*, 18(6):1156–1194.
- Chambers, C. and Scaffidi, C. (2013). Smell-driven performance analysis for end-user programmers. In *Proc. of VLH/CC '13*, pages 159–166.
- Chambers, C. and Scaffidi, C. (2015). Impact and utility of smell-driven performance tuning for end-user programmers. *Journal of Visual Languages & Computing*, 28:176–194. to appear.
- Chapman, C. and Stolee, K. T. (2016). Exploring regular expression usage and context in python. under review.
- Chen, T.-H., Nagappan, M., Shihab, E., and Hassan, A. E. (2014). An empirical study of dormant bugs. In *Proceedings of the 11th Working Conference on Mining Software Repositories*, MSR 2014, pages 82–91, New York, NY, USA. ACM.
- Dattero, R. and Galup, S. D. (2004). Programming languages and gender. *Commun. ACM*, 47(1):99–102.

- Du Bois, B., Demeyer, S., Verelst, J., Mens, T., and Temmerman, M. (2006). Does god class decomposition affect comprehensibility? In *IASTED Conf. on Software Engineering*, pages 346–355.
- Dyer, R., Rajan, H., Nguyen, H. A., and Nguyen, T. N. (2014). Mining billions of ast nodes to study actual and potential usage of java language features. In *Proceedings of the 36th International Conference on Software Engineering, ICSE 2014*, pages 779–790, New York, NY, USA. ACM.
- Fowler, M. (1999). *Refactoring: improving the design of existing code*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA.
- Galler, S. J. and Aichernig, B. K. (2014). Survey on test data generation tools. *Int. J. Softw. Tools Technol. Transf.*, 16(6):727–751.
- Ghosh, I., Shafiei, N., Li, G., and Chiang, W.-F. (2013). Jst: An automatic test generation tool for industrial java applications with strings. In *Proceedings of the 2013 International Conference on Software Engineering, ICSE ’13*, pages 992–1001, Piscataway, NJ, USA. IEEE Press.
- Grechanik, M., McMillan, C., DeFerrari, L., Comi, M., Crespi, S., Poshyvanyk, D., Fu, C., Xie, Q., and Ghezzi, C. (2010). An empirical investigation into a large-scale java open source code repository. In *Proceedings of the 2010 ACM-IEEE International Symposium on Empirical Software Engineering and Measurement, ESEM ’10*, pages 11:1–11:10, New York, NY, USA. ACM.
- Griswold, W. G. and Notkin, D. (1993). Automated assistance for program restructuring. *ACM Trans. Softw. Eng. Methodol.*, 2(3):228–269.
- Hermans, F., Pinzger, M., and van Deursen, A. (2012). Detecting code smells in spreadsheet formulas. In *Proc. of ICSM ’12*, pages 409–418.
- Hermans, F., Pinzger, M., and van Deursen, A. (2014). Detecting and refactoring code smells in spreadsheet formulas. *Empirical Software Engineering*, pages 1–27.



- Kiezun, A., Ganesh, V., Artzi, S., Guo, P. J., Hooimeijer, P., and Ernst, M. D. (2013). Hampi: A solver for word equations over strings, regular expressions, and context-free grammars. *ACM Trans. Softw. Eng. Methodol.*, 21(4):25:1–25:28.
- Le Goues, C., Nguyen, T., Forrest, S., and Weimer, W. (2012). GenProg: A generic method for automated software repair. *Transactions on Software Engineering*, 38(1):54–72.
- Lee, J., Pham, M.-D., Lee, J., Han, W.-S., Cho, H., Yu, H., and Lee, J.-H. (2010). Processing sparql queries with regular expressions in rdf databases. In *Proceedings of the ACM Fourth International Workshop on Data and Text Mining in Biomedical Informatics*, DTMBIO '10, pages 23–30, New York, NY, USA. ACM.
- Li, Y., Krishnamurthy, R., Raghavan, S., Vaithyanathan, S., and Jagadish, H. V. (2008). Regular expression learning for information extraction. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, EMNLP '08, pages 21–30, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Linares-Vásquez, M., Bavota, G., Bernal-Cárdenas, C., Oliveto, R., Di Penta, M., and Poshyvanyk, D. (2014). Mining energy-greedy api usage patterns in android apps: An empirical study. In *Proceedings of the 11th Working Conference on Mining Software Repositories*, MSR 2014, pages 2–11, New York, NY, USA. ACM.
- Livshits, B., Whaley, J., and Lam, M. S. (2005). Reflection analysis for java. In *Proceedings of the Third Asian Conference on Programming Languages and Systems*, APLAS'05, pages 139–160, Berlin, Heidelberg. Springer-Verlag.
- Mens, T. and Tourwé, T. (2004). A survey of software refactoring. *IEEE Trans. Soft. Eng.*, 30(2):126–139.
- Meyerovich, L. A. and Rabkin, A. S. (2013). Empirical analysis of programming language adoption. In *Proceedings of the 2013 ACM SIGPLAN International Conference on Object Oriented Programming Systems Languages & Applications*, OOPSLA '13, pages 1–18, New York, NY, USA. ACM.

- Møller, A. (2010). dk.brics.automaton – finite-state automata and regular expressions for Java. <http://www.brics.dk/automaton/>.
- network (2015). The Bro Network Security Monitor. <https://www.bro.org/>.
- Opdyke, W. F. (1992). *Refactoring Object-oriented Frameworks*. PhD thesis, Champaign, IL, USA. UMI Order No. GAX93-05645.
- Parnin, C., Bird, C., and Murphy-Hill, E. (2013). Adoption and use of java generics. *Empirical Softw. Engg.*, 18(6):1047–1089.
- re2 (2015). RE2. <https://github.com/google/re2>.
- Richards, G., Lebresne, S., Burg, B., and Vitek, J. (2010). An analysis of the dynamic behavior of javascript programs. *SIGPLAN Not.*, 45(6):1–12.
- Spishak, E., Dietl, W., and Ernst, M. D. (2012). A type system for regular expressions. In *Proceedings of the 14th Workshop on Formal Techniques for Java-like Programs, FTfJP '12*, pages 20–26, New York, NY, USA. ACM.
- Stolee, K. T. and Elbaum, S. (2011). Refactoring pipe-like mashups for end-user programmers. In *International Conference on Software Engineering*.
- Stolee, K. T. and Elbaum, S. (2013). Identification, impact, and refactoring of smells in pipe-like web mashups. *IEEE Trans. Softw. Eng.*, 39(12):1654–1679.
- Tillmann, N., de Halleux, J., and Xie, T. (2014). Transferring an automated test generation tool to practice: From pex to fakes and code digger. In *Proceedings of the 29th ACM/IEEE International Conference on Automated Software Engineering, ASE '14*, pages 385–396, New York, NY, USA. ACM.
- Trinh, M.-T., Chu, D.-H., and Jaffar, J. (2014). S3: A symbolic string solver for vulnerability detection in web applications. In *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security, CCS '14*, pages 1232–1243, New York, NY, USA. ACM.

- Veanes, M., Halleux, P. d., and Tillmann, N. (2010). Rex: Symbolic regular expression explorer. In *Proceedings of the 2010 Third International Conference on Software Testing, Verification and Validation*, ICST '10, pages 498–507, Washington, DC, USA. IEEE Computer Society.
- Weimer, W., Forrest, S., Le Goues, C., and Nguyen, T. (2010). Automatic program repair with evolutionary computation. *Communications of the ACM Research Highlight*, 53(5):109–116.
- Yeole, A. S. and Meshram, B. B. (2011). Analysis of different technique for detection of sql injection. In *Proceedings of the International Conference & Workshop on Emerging Trends in Technology*, ICWET '11, pages 963–966, New York, NY, USA. ACM.