

# Readability and Refactoring of Regular Expressions

Carl Chapman  
Department of Computer Science  
Iowa State University  
carl1976@iastate.edu

Kathryn T. Stolee  
Department of Computer Science  
North Carolina State University  
ktstolee@ncsu.edu

## ABSTRACT

Regexes are hard to understand. Let me tell you how.

## 1. INTRODUCTION

Our contributions are as follows:

- 
- Defined regex refactoring categories for understandability backed by empirical evidence
- Identified 3 or so other regex refactorings categories and specific instances that are worthy of further investigations
- Identified a few regex refactorings that can be eliminated because both options are equally readable

The rest of this paper is organized as follows:  
Related work, study, results, discussion, conclusion.

## 2. RELATED WORK

**TODO.MID: We are building on the survey that indicates regexes are hard to read, and the apparent lack of any regex readability refactoring attempts. Many papers have talked about refactoring, basically it is changing the form but not the behavior.**

Prior work that surveyed developers about regex usage found that in a small software company, the 18 surveyed developers compose an average of 172 regexes per year. This is 48% higher than the number of regexes composed annually by MTurk participants in this work, which may be due to the nature of the jobs performed by the two populations.

## 3. STUDY

First we define a 'Functional Regex'(FR) as some regex that performs in a specific way. For many FRs, there are several concrete ways to express a single FR. We define a

concrete regex(CR) as a regex expressed with a particular pattern String. Here is one illustration of these definitions:

**TODO.NOW: create some examples for these terms**

We identified 10 loose groups of FRs, described in this table:

**TODO.NOW: create a table explaining the 10 groups**

For each of these groups we created either two concrete versions of three FRs or three concrete versions of two FRs.

Each of the 10 categories had 6 concrete versions of some FR and so there are 60 CRs. For each CR, we selected 5 *example strings* designed to test the understanding of the CR. The idea is that different CRs may have different levels of readability, even when they are representing the same FR. We define readability as the ability to look at the CR and determine if an *example string* can be matched by it or not.

**TODO.NOW: create some illustration of one matching subtask**

### 3.1 Metrics

### 3.2 Design

In Mechanical Turk, we designed a 180 tasks composed of 10 matching subtasks, so that each of the 60 CRs had 30 separate observations (each an average of 5 *example string* problems). These 1800 observations are what the analysis will focus on.

**TODO.NOW: Were there qualification questions that required participants to demonstrate knowledge of regexes?**

### 3.3 Participants

In total, there were 180 different participants in the study. A majority were male (83%) with an average age of 31. Most had at least an Associates degree (72%) and most were at least somewhat familiar with regexes prior to the study (87%). On average, participants compose 67 regexes per year with a range of 0 to 1000. Fittingly, participants read more regexes than they write with an average of 116 and a range from 0 to 2000. Figure ?? summarizes the self-reported participant characteristics from the qualification survey.

## 4. RESULTS

(for rough draft...)

Looks like M6 and M8 are the best meta-refactorings according to ANOVA. If you peek at MTResults Processing.csv on google docs, M6 has the best refactoring...every OCTAL type should be converted to an OR or preferably a CCC.

	What is your gender?	n	%
1.	Male	149	83%
	Female	27	15%
	Prefer not to say	4	2%
2.	What is your age? $\mu = 31, \sigma = 9.3$		
	Education Level?	n	%
3.	High School	5	3%
	Some college, no degree	46	26%
	Associates degree	14	8%
	Bachelors degree	78	43%
	Graduate degree	37	21%
	Familiarity with regexes?	n	%
4.	Not familiar at all	5	3%
	Somewhat not familiar	16	9%
	Not sure	2	1%
	Somewhat familiar	121	67%
	Very familiar	36	20%
5.	How many regexes do you compose each year? $\mu = 67, \sigma = 173$		
6.	How many regexes (not written by you) do you read each year? $\mu = 116, \sigma = 275$		

**Figure 1: Participant Profiles,  $n = 180$**

M8 has a weak P value, but still ok in one case (0.12) and consistently says that 'aa\*' should be written as 'a+'.

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)\\}' (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]' which has a lot of escape characters, so that the hex group was easier to understand than the CCC.

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\_]+);' (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'

**TODO.MID: more data from the composition problems**

## 5. DISCUSSION: REFACTORING OPPORTUNITIES

## 6. THREATS TO VALIDITY

### 6.1 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.

### 6.2 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 we used in the evaluation were inspired by those commonly 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.

## 7. CONCLUSION

### Acknowledgements

This work is supported in part by NSF SHF-EAGER-1446932.

Table 1: How frequently is each alternative expression style used?

code	R1	R2	R3	acc1	acc2	acc3	cmp1	cmp2	cmp2	Pacc	Pcmp
M0	'^[1-9][0-9]*\$'	'^[1-9][0-9]*\$'	NA	0.85	0.85	NA	27.0	26.3	NA	0.42	1.00
M1	'^[1-9][0-9]*\$'	'^[1-9][0-9]*\$'	NA	0.68	0.75	NA	21	23	NA	0.28	0.67
M2	'^[1-9][0-9]*\$'	'^[1-9][0-9]*\$'	NA	0.82	0.86	NA	23.3	25.3	NA	0.67	0.42
M3	'^[1-9][0-9]*\$'	'^[1-9][0-9]*\$'	'^[1-9][0-9]*\$'	0.81	0.86	0.84	15.5	27.5	19.5	0.62	0.19

Table 2: How frequently is each alternative expression style used?

name	description	example	nPatterns	% patterns	nProjects	% projects
C1	CCC and RNG	'^[1-9][0-9]*\$'	2,510	0.18	818	0.53
C2	CCC, no RNG or defaults	'[aeiouy]'	1,283	0.09	551	0.36
C3	OR of single chars	'\\(t n r " \ \\)'	156	0.01	174	0.11
C4	Contains NCCC	'[0-9]'	1,935	0.14	776	0.50
C5	CCC and defaults, no RNG	'[-+\\d.]'	675	0.05	382	0.25
C6	OR containing defaults	'(\\s _)+'	90	0.01	131	0.08
CCC REMAINDER	null	'^\\s*\$'	7,491	0.55	1,223	0.79
D1	DBB other than 0,1	'^x{1,4}\$'	279	0.02	212	0.14
D2	DBB exactly 0,1	'^[\\n\\r]{0,1}'	23	0.00	25	0.02
D3	Contains QST	'^https?:\\/\\S+\$'	1,871	0.14	646	0.42
D4	OR expressing repetition	'^(Q QQ)\\<(\\+)>\$'	10	0.00	27	0.02
DBB REMAINDER	null	'^[1-9][0-9]*\$'	11,491	0.85	1,473	0.95
LIT REMAINDER	null	'(\\w+)\\s+(\\d+)'	319	0.02	538	0.35
T1	no HEX, OCT or CCC-wrapped	'get_tag'	12,482	0.92	1,485	0.96
T2	Contains HEX literal	'(\\x1b ~{)'	479	0.04	243	0.16
T3	Contains CCC-wrapped	'[\$]([\\d+:([~}]+))[]]'	307	0.02	268	0.17
T4	Contains OCT literal	'[\\041-\\176]+:\$'	14	0.00	37	0.02
L1	LWB like M,, M>0	'^_{2,}\\.[^_]+_?\$'	91	0.01	166	0.11
L2	Contains KLE	'^(\\.\\*?)''\\s*(#\\.\\*)?\$',	6,017	0.44	1,097	0.71
L3	Contains ADD	'[A-Z][a-z]+'	6,003	0.44	1,207	0.78
LWB REMAINDER	null	'-py([123]\\.[0-9])\$'	3,979	0.29	944	0.61
OR REMAINDER	null	'^012TF\\*{9}\$'	11,495	0.85	1,489	0.96
R1	OR with prefix or suffix	'^(def class)\\s+'	1,529	0.11	603	0.39
R2	top-level OR	'([ ]+_ ([ ]+) ([ ]+))'	573	0.04	396	0.26
S1	Contains SNG	'[^ ]*\\. [A-Z]{3}'	581	0.04	340	0.22
S2	Repeated non-literals	'ff:ff:ff:ff:ff:ff'	927	0.07	497	0.32
S3	DBB like M,M	'U\\dA-F]{5,5}'	27	0.00	32	0.02
SNG REMAINDER	null	'@@BINARYDIR@@'	12,209	0.90	1,498	0.97

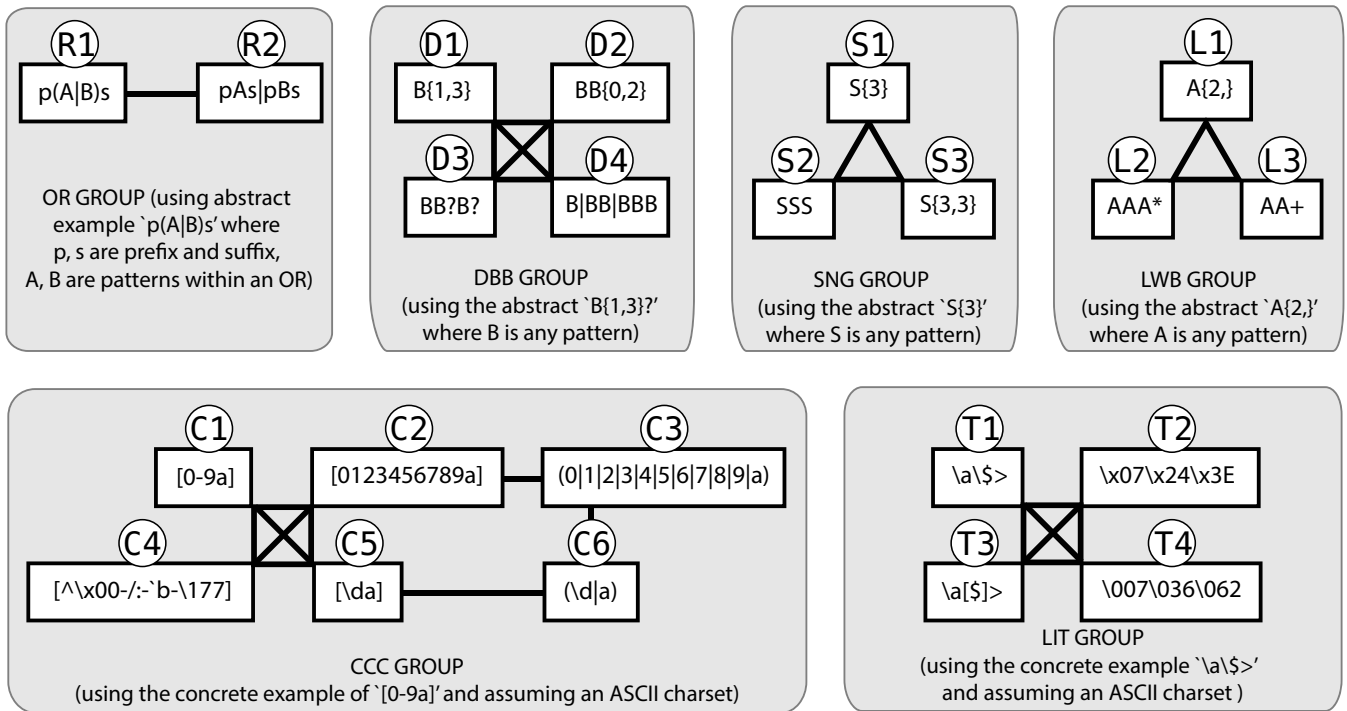


Figure 2: Some possible refactorings