# Regex Feature ●Use In Practice

Carl Chapman and Kathryn T. Stolee
Department of Computer Science
Iowa State University
{carl1978, kstolee}@iastate.edu

*Abstract*—**Regular expressions are used frequently in many programming languages for form validation, ad-hoc file searches, and simple parsing. Many researchers have created solvers, parsers, tools and theoretical works which include a subset of the regular expression features used in practice. Yet, there does not exist an empirical study of regular expression feature usage in practice that could inform the choices that researchers are making about what features to include and exclude. In this paper, we explore feature usage in practice, focusing on how often certain excluded features are used, and what use cases are associated with what features. To do this, we analyzed about 4000 open source Python projects from GitHub and explored the regular expressions contained within. Our results indicate that TODO: high level results**

## I. INTRODUCTION

Regular expressions are used extensively in many programming languages, for example, to search text files [3], in form validation, and for XYZ.

## II. MOTIVATION

Bugs related to regular expressions are common, resulting in tens of thousands of bug reports [4].

## III. RELATED WORK

### A. Research on Regular Expressions

Visual debugging of regular expressions [1]

Static analysis to reduce errors in building regular expressions by using a type system to identify errors like `PatternSyntaxExceptions` and `IndexOutOfBoundsExceptions` at compile time [4].

### B. Research on Regular Expressions

Visual debugging of regular expressions [1]

### C. Research that Depends on Regular Expression Usage

Regular expressions are used as queries in a data mining framework [2]

## IV. STUDY

### A. RECORDING REGEX USAGES

Fun fact: while creating similarity matrix, row 5464 took 2 hours, or almost 1 second per cell avg, only suffering 18 timeouts (1.2 secs). What is this pesky pattern?

Using GHTorrent, we found the clone urls of -1 github projects that had Python listed as the main language. After consulting with github about the polite way to mine their service, we used 3 separate machines to clone projects in parallel. On each machine, one attempt was made to clone each project into a unique directory. -1 of these attempts failed, so a total of 9,644 projects were scanned overall. For each of these projects, the java program launched a Python process that used 'Astroid' to build the AST of each python file in the unique directory, recording uses of the 're' module. Here is an example Python code using the 're' module:

Placeholder for several examples of useage of the 're' module...image? Shows example of function, flags pattern.

TODO...mention rewinding, duplicate skipping, exact criteria for citing a usage...For each of the 151,025 regex usages observed, we recorded which 're' function was used, what flags (if any) were used, and what pattern was found.

### B. FILTERING REGEX USAGES

Because we want to let regex patterns alone define behavior later in our analysis, all regex usages where behavioral flags are used are discarded, along with all invalid pattern strings. TODO - change to percents 21,424 usages were discarded because they used behavior-altering flags, and 9,237 without flags were discarded because their pattern strings were invalid (could not be compiled by Python into a regex object).

120,345 regex usages remain (meow percent), which are then condensed to 24,605 distinct pattern strings. The number of distinct pattern strings is much smaller because the same pattern string is often used in many distinct files and projects by different functions with different flags and these all count as separate usages. The resulting set of patten strings were parsed using TODO - [1] PCRE parser. 0.5% of the strings caused the PCRE parser to raise an error. Another 0.2% used regex features that we have chosen to exclude in this study (most notably named capturing groups). The 24,325 distinct pattern strings that remain are given a weight based upon how many projects the pattern appeared in (FYI no forked projects were scanned).

TODO - clean this: So then we use the following 4? techniques to try and best represent the entire corpus using a few regexes: 1. weight (frequency of usage across projects) 2. features (matching properties of feature usage) 3. syntactic clustering 4. semantic clustering ...and also we give the topN clones by weight, and a list that tries to include information from all 4 lists? That will be all. Future work: mine more projects and compare that result with this result. Compare across languages (java, javascript,ruby,etc.).

---

[1] www.source

TODO - list research questions (metrics and how they are computed?)

RQ1: How is the `re module` used in python projects? RQ2: How frequently are regexes used in python projects? RQ3: Which regex language features are used most commonly in python? RQ4: How syntactically diverse are regexes used in python? RQ5: How semangically diverse are regexes used in python?

### 1) FEATURES

Here is a table showing all the features included in this study and which features are supported by four popular regex research projects/tools:

## V. RESULTS

### A. CONTEXT AND CORPUS ORIGIN

#### 1) SATURATION

Although 50.5% of the projects observed had at least one regex usage, only 9.1% of the files observed had at least one regex usage.

From the above figure/table, we see that on average each project had 2 files containing any regex usage, out of an average of 8 files. Each of the files that did have a regex usage had an average of 2 regex usages. Because we scanned 9,644 projects, we would expect to have seen 154,304 regex usages, which is lower than the actual 151,025 usages observed.
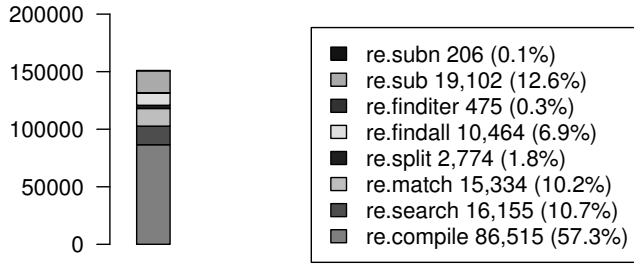


Fig. 1.   How often are the 8 re functions used?

#### 2) FUNCTIONS AND FLAGS

As seen in Figure 1 The 'compile' function encompasses 57.3% of all usages, even though every compiled regex object can only be used by calling other functions. (TODO-Why?)
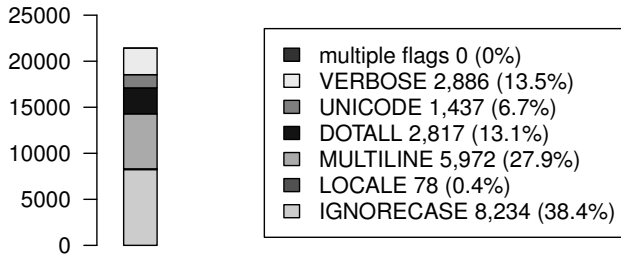


Fig. 2.   Which behavioral flags are used?

85.8% of all regex usages did not use a flag or specified a non-behavioral flag (default or debug). Of all behavioral flags used, ignorecase (38.4%) and multiline (27.9%) were the most frequently used. It is also worth noting that although multiple flags can be combined using a bitwise or, this was never observed. (remove this last part if it is observed later)

#### 3) GENERAL CHARACTERISTICS OF REGEXES FOUND :
...TODO

#### 4) Top 10 Regex Patterns by weight :

#### 5) All Features :

Literal tokens were found in (TODO) 101% of patterns, and accounted for 75% of all tokens. Excluding literal tokens and features that were not present in any pattern, the following stats...make a sentence, these are some stats about the features:

some more text, IDK

| pair | example from corpus | nTimes |
|------|---------------------|--------|
| CG::ADD | `'(:+)'` | 7564 |
| CG::KLE | `'(:)*'` | 7147 |
| ANY::KLE | `'.*'` | 6741 |
| CG::ANY | `'(.)'` | 5782 |
| CCC::CG | `"([']")"` | 4779 |
| CCC::ADD | `'[ ]+'` | 4699 |
| RNG::CCC | `'[A-Z]'` | 4625 |
| ADD::KLE | `'-*(.+)'` | 4376 |
| ANY::ADD | `'.+'` | 3953 |
| WSP::KLE | `'\\s*'` | 3774 |

OK now that is all for section 2. Now in section 3 I want to look at clustering by string similarity using mcl clustering algorithm. Here are the top 6 clusters using various string similarity metrics:

TODO - multiple boxplots for all 5-6 demonstrating cluster size and then also have # of clusters, pick smallest number of clusters and then use that.

## VI. DISCUSSION

...only 9.1% of the files observed had at least one regex usage. This indicates that regex usage may usually be concentrated in just a few files.

## VII. CONCLUSION

### ACKNOWLEDGMENT

### REFERENCES

[1] F. Beck, S. Gulan, B. Biegel, S. Baltes, and D. Weiskopf. 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, 2014. ACM.

| pattern | weight |
|---|---|
| '.*?\n' | 482 |
| '[0-9]+' | 475 |
| '^([\\w., ]+=)?\\s*[\\w\\.]+\\\(.*\\))$' | 471 |
| '\\s+' | 467 |
| '^[^(]*' | 464 |
| '   \\.\\.\\.\\.+:' | 461 |
| '(In \\[[0-9]+\\]: )|(   \\.\\.\\.\\.+:)' | 461 |
| '(Out\\[[0-9]+\\]: )|(   \\.\\.\\.\\.+:)' | 461 |
| '\\-+' | 461 |
| '^\\s*[#*=]{4,}\\n[a-z0-9 -]+\\n[#*=]{4,}\\s$' | 454 |

| rank | code | description | example | brics | hampi | Rex | RE2 | nPatterns | % patterns | nFiles | %files | nProjects | % projects |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | CG | a capture group | (caught) | ● | ● | ● | ● | 12,836 | 52.8 | 26,241 | 52.1 | 3,772 | 77.4 |
| 2 | ADD | one-or-more repetition | z+ | ● | ● | ● | ● | 10,930 | 44.9 | 26,502 | 52.6 | 3,764 | 77.2 |
| 3 | KLE | zero-or-more repetition | .* | ● | ● | ● | ● | 10,831 | 44.5 | 22,912 | 45.5 | 3,585 | 73.6 |
| 4 | CCC | custom character class | [aeiou] | ● | ● | ● | ● | 8,099 | 33.3 | 22,440 | 44.6 | 3,338 | 68.5 |
| 5 | ANY | any non-newline char | . | ● | ● | ● | ● | 8,525 | 35 | 17,959 | 35.7 | 3,315 | 68 |
| 6 | RNG | chars within a range | [a-z] | ● | ● | ● | ● | 4,850 | 19.9 | 15,399 | 30.6 | 2,985 | 61.3 |
| 7 | STR | start-of-line | ^ | ○ | ● | ● | ● | 6,305 | 25.9 | 15,582 | 30.9 | 2,937 | 60.3 |
| 8 | END | end-of-line | $ | ○ | ● | ● | ● | 5,844 | 24 | 15,037 | 29.9 | 2,872 | 58.9 |
| 9 | WSP | \t \n \r \b \f or space | \s | ○ | ● | ● | ● | 4,979 | 20.5 | 13,887 | 27.6 | 2,796 | 57.4 |
| 10 | NCCC | negated CCC | [^qwxf] | ● | ● | ● | ● | 3,571 | 14.7 | 10,872 | 21.6 | 2,706 | 55.5 |
| 11 | OR | logical or | a|b | ● | ● | ● | ● | 3,613 | 14.9 | 11,215 | 22.3 | 2,604 | 53.4 |
| 12 | WRD | [a-zA-Z0-9_] | \w | ○ | ● | ● | ● | 2,525 | 10.4 | 8,242 | 16.4 | 2,365 | 48.5 |
| 13 | QST | zero-or-one repetition | z? | ● | ● | ● | ● | 3,347 | 13.8 | 10,614 | 21.1 | 2,215 | 45.5 |
| 14 | LZY | as few reps as possible | z+? | ○ | ● | ○ | ● | 2,604 | 10.7 | 6,581 | 13.1 | 2,155 | 44.2 |
| 15 | DEC | any of: 0123456789 | \d | ○ | ● | ● | ● | 4,237 | 17.4 | 11,298 | 22.4 | 1,864 | 38.3 |
| 16 | NCG | group without capturing | a(?:b)c | ○ | ● | ○ | ● | 1,441 | 5.9 | 4,756 | 9.4 | 1,522 | 31.2 |
| 17 | NCG | named capture group | (?P<name>x) | ○ | ● | ○ | ● | 1,802 | 7.4 | 3,627 | 7.2 | 877 | 18 |
| 18 | NLKA | sequence doesn't follow | a(?!yz) | ○ | ○ | ● | ○ | 217 | 0.9 | 1,638 | 3.3 | 836 | 17.2 |
| 19 | SNG | exactly n repetition | z{8} | ● | ● | ● | ● | 1,123 | 4.6 | 3,524 | 7 | 835 | 17.1 |
| 20 | LWB | at least n repetition | z{15,} | ● | ● | ● | ● | 199 | 0.8 | 1,294 | 2.6 | 828 | 17 |
| 21 | NWSP | any non-whitespace | \S | ○ | ● | ● | ● | 835 | 3.4 | 1,882 | 3.7 | 671 | 13.8 |
| 22 | DBB | $n \le x \le m$ repetition | z{3,8} | ● | ● | ● | ● | 683 | 2.8 | 1,800 | 3.6 | 612 | 12.6 |
| 23 | WNW | word/non-word boundary | \b | ○ | ○ | ○ | ● | 405 | 1.7 | 1,217 | 2.4 | 551 | 11.3 |
| 24 | NWRD | non-word chars | \W | ○ | ● | ● | ● | 173 | 0.7 | 839 | 1.7 | 503 | 10.3 |
| 25 | NLKB | sequence doesn't precede | (?<!x)yz | ○ | ○ | ○ | ○ | 165 | 0.7 | 919 | 1.8 | 487 | 10 |
| 26 | NDEC | any non-decimal | \D | ○ | ● | ● | ● | 60 | 0.2 | 632 | 1.3 | 407 | 8.4 |
| 27 | LKA | matching sequence follows | a(?=bc) | ○ | ○ | ○ | ○ | 255 | 1 | 871 | 1.7 | 389 | 8 |
| 28 | OPT | options wrapper | (?i)CasE | ○ | ● | ○ | ● | 430 | 1.8 | 932 | 1.9 | 380 | 7.8 |
| 29 | LKB | matching sequence precedes | (?<=a)bc | ○ | ○ | ○ | ○ | 183 | 0.8 | 663 | 1.3 | 315 | 6.5 |
| 30 | ENDZ | absolute end of string | \Z | ○ | ○ | ○ | ● | 156 | 0.6 | 351 | 0.7 | 214 | 4.4 |
| 31 | BKR | match the $i^{th}$ CG | \1 | ○ | ○ | ○ | ○ | 98 | 0.4 | 334 | 0.7 | 211 | 4.3 |
| 32 | BKRN | references NCG | \g<name> | ○ | ● | ○ | ○ | 30 | 0.1 | 108 | 0.2 | 84 | 1.7 |
| 33 | NWNW | negated WNW | \B | ○ | ○ | ○ | ● | 9 | 0 | 33 | 0.1 | 29 | 0.6 |
| 34 | VWSP | matches U+000B | \v | ○ | ○ | ● | ● | 13 | 0.1 | 22 | 0 | 20 | 0.4 |

[2] A. Begel, Y. P. Khoo, and T. Zimmermann. 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, 2010. ACM.

[3] C. L. A. Clarke and G. V. Cormack. On the use of regular expressions for searching text. *ACM Trans. Program. Lang. Syst.*, 19(3):413–426, May 1997.

[4] E. Spishak, W. Dietl, and M. D. Ernst. 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, 2012. ACM.