

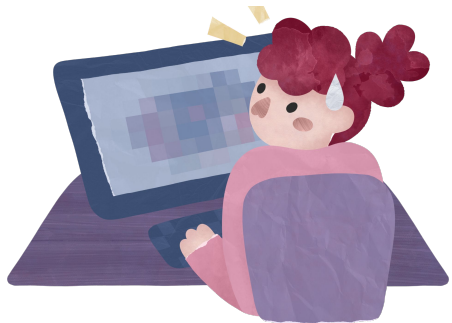
$\wedge[\backslash s\backslash u200c]+|[\backslash s\backslash u200c]+\$$

The Parsing and Analysis of [rR]eg(ular)* *([eE]x(pression(s)*))*

Will Beddow, Antonia Ritter, Shiyue Zhang,
Vicente Riquelme

The Challenges of Using Regexes

Why Aren't Regular Expressions a Lingua Franca? An Empirical Study on the Re-use and Portability of Regular Expressions



Malformed Stack Overflow Post Chokes Regex, Crashes Site

823
↑
↓



Blog/Article/Link [Cloudflare outage caused by deploying bad regular expression that caused 100% CPU usage worldwide, dropping up to 82% of traffic](#) (self.sysadmin)

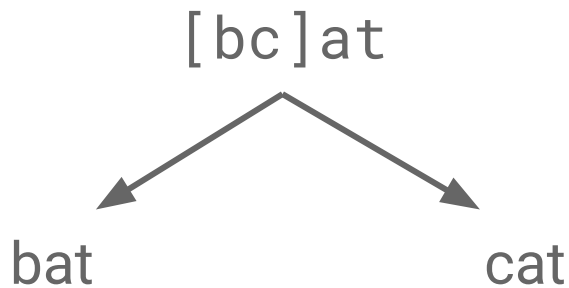
submitted 2 years ago by TyroPyro

[Cloudflare Blog](#)^[1]

275 comments [source](#) [share](#) [save](#) [hide](#) [give award](#) [report](#) [crosspost](#) [hide all child comments](#)

Regular Expressions

A regular expression (**regex**) is a pattern used to match particular strings



```
#?([\da-fA-F]{2})([\da-fA-F]{2})([\da-fA-F]{2})
```



Any hex code, for example
`#FF5733`

Project Overview

- Our tool evaluates:
 - **Security**
 - **Understandability**
 - **Generalizability**

Wh(
what |
why
)

(In | Out)put

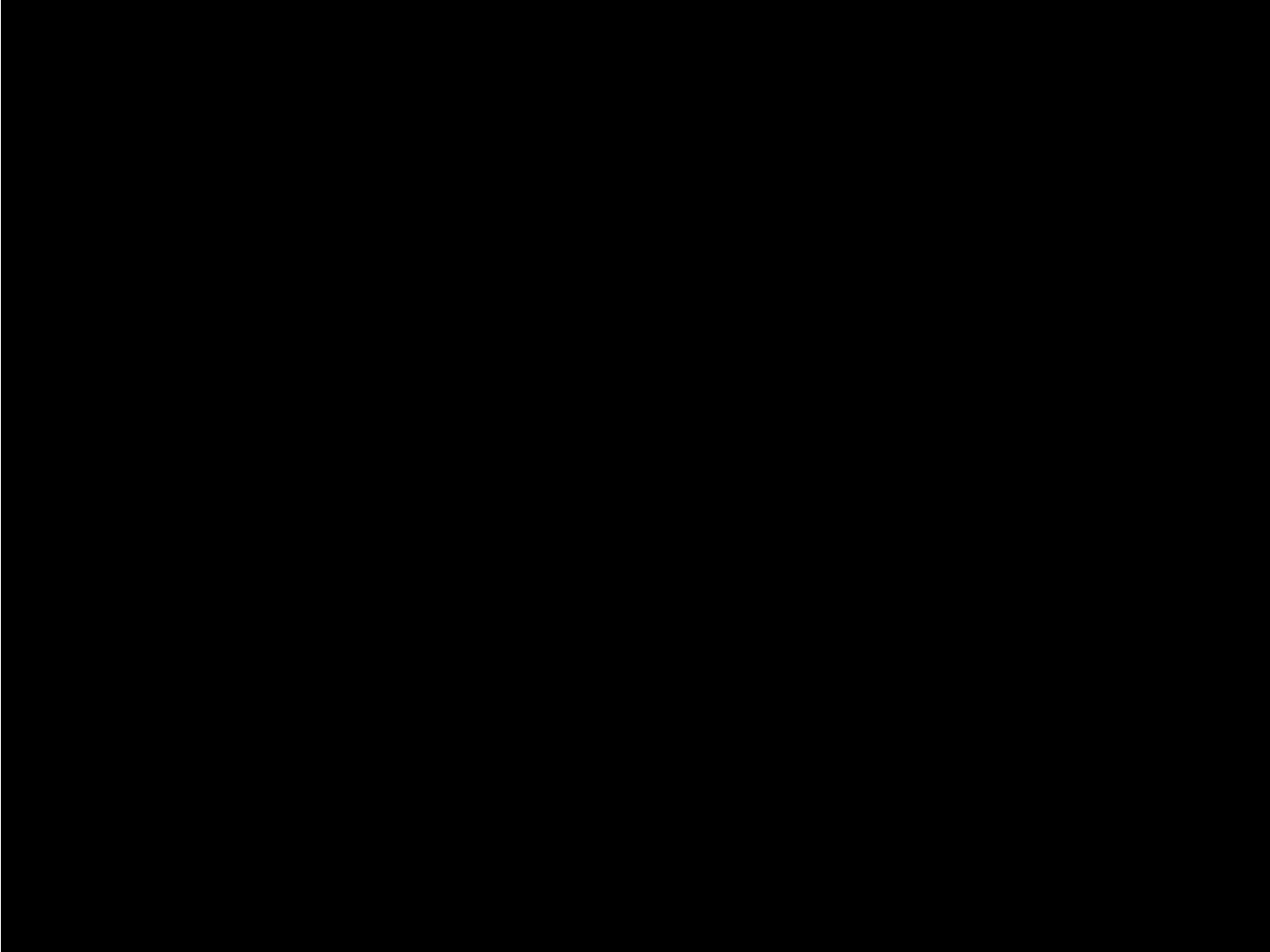
Directory containing:
 $\wedge[\backslash s\backslash u200c]^+|[\backslash s\backslash u200c]^+\$$



Regex-Library



Analysis and
Recommendations



```
antonia@ARXPS:~/Dropbox (Carleton College)/21_CS_Comps/regex-library$ python3 main.py ../demo-input/ --security
Expression ^<!|\-|-(.*)+(\|/){0,1}\-|-\>a(?:bc|b)c\141S{3}$ from ../demo-input/example_1.py line 4
|0import re
|1
|2def main():
|3     my_patt = re.compile(r"^<!|\-|-(.*)+(\|/){0,1}\-|-\>a(?:bc|b)c\141S{3}$")
|4     my_patt.match("<div>test</div>")
|5
|6if __name__ == "__main__":
|7     main()
|8
|9
|10
|11
^Understandability suggestion: a
[aaa]
^Understandability suggestion: repeated characters in []
S{3}
^Understandability suggestion: SSS
?=
^Java, JavaScript, and Golang may not support conditional
<!-->abcbcaaaa!---->abcbcaaaa!----->abcbc---->abcbcaaaS-->abcbc>abcb>abcbcaaaS"aS---->abcbcaa!---->abcbcaaaaS
^This string gave the expression a run time longer than 1 second
<!-->abcbcaaaa---->abcbcaaaaS>abcbc+rcae>abcbcs<!-->abcbP->abcaafu&<---->abcbcaaaaS@--
Bfi>abcbcaaaaS-caa<!-->abcbcaaaaSbcaaa<!-->abcbcaaaaS`-->abc>abcbca
^This string gave the expression a run time longer than 1 second
Expression hello from ../demo-input/folder_1/example_2.js line 4
|0
|1function test_function() {
|2     let text = "Hello";
|3     let n = text.search(/Hello/i);
|4}
|5
```

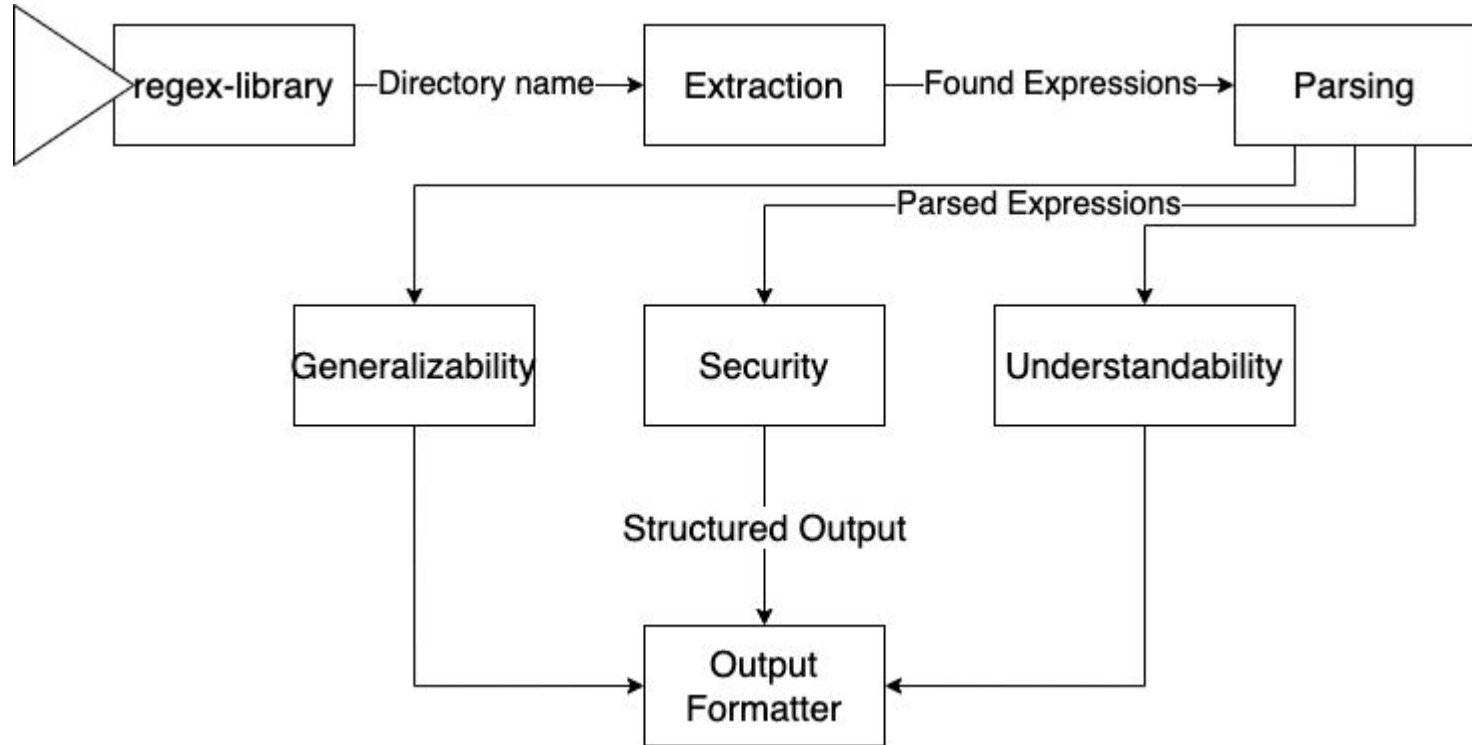
The raw expression, and where it was found

Excerpt of file where expression was found

Targeted tips and analysis

Output

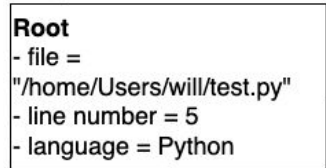
Structure of the library



Data Format

Anatomy of an Expression Representation

Example of how a simple expression, `"\d+"` might be represented in the system



The Root expression object tells the system what file the expression came from, where in the file it can be found, and what language it's in

```
Expression
raw = "\d+"
```

The raw string of the expression can be used to reference the plain regex, without any modifications

```
Tokens
```

```
Token
raw = "\d"
type = "digit"
```

```
Token
raw = "+"
type = "Quantifier Modifier Plus"
```

The token types, generated by the parser, give the different components of the analyzer insights into what the Regex is doing, and how it's trying to do it

Brew



```
brew install zhangshyue/regexanalyzer/regexAnalyzer  
regex-library [directory_path]
```

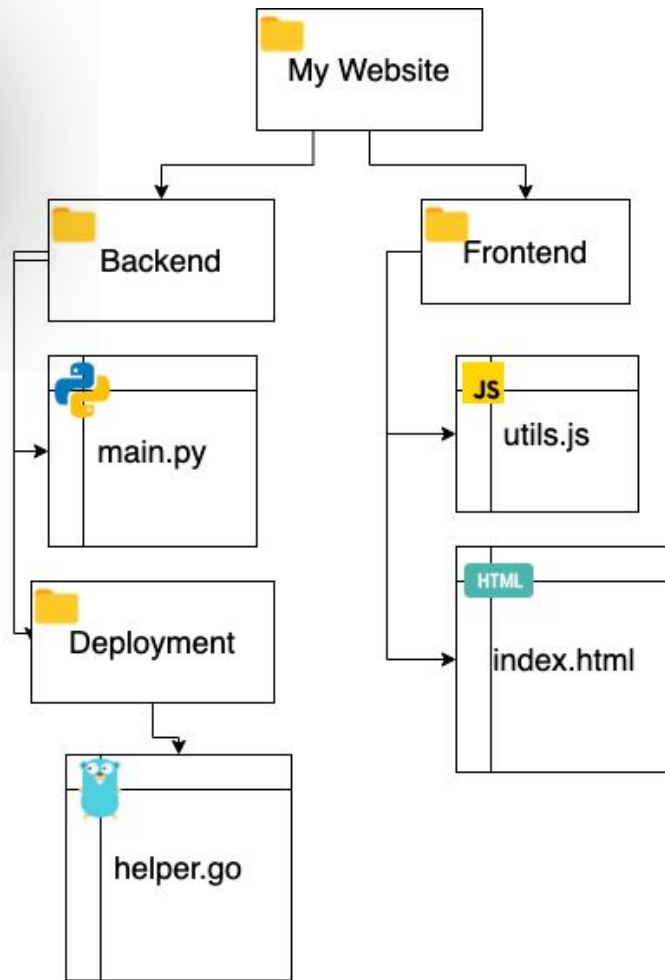


Homebrew

The missing package manager
for macOS (or Linux)



```
python3 regex-library/main.py "My Website"
```



Finding Regular Expressions in Files

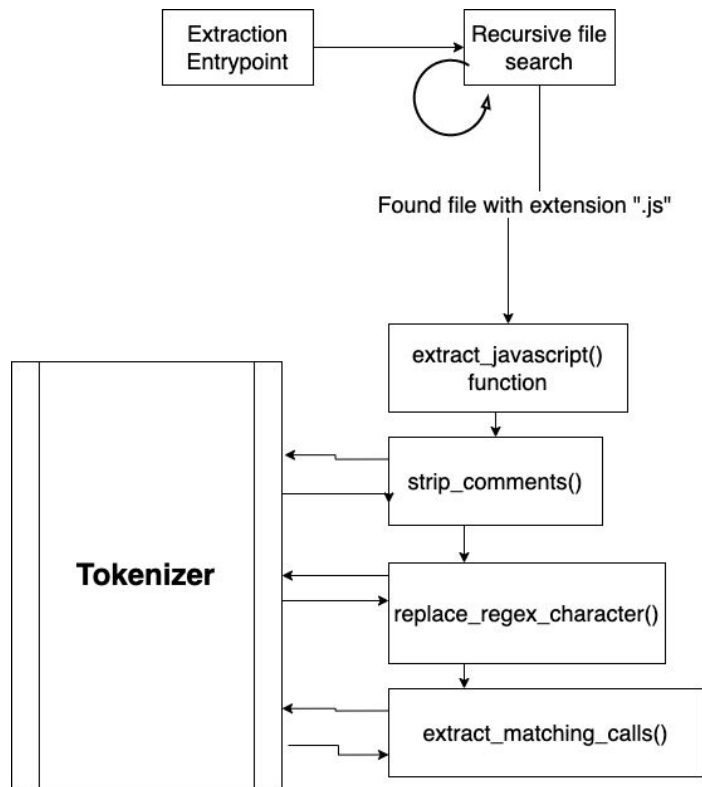
```
match language {  
  SupportedLanguage::Python => res = extract_python( contents_raw: file_contents),  
  SupportedLanguage::JavaScript => res = extract_js( contents_raw: file_contents),  
  SupportedLanguage::GoLang => res = extract_go( contents_raw: file_contents),  
  SupportedLanguage::Java =  
  SupportedLanguage::Rust =  
  SupportedLanguage::PHP =  
  SupportedLanguage::Ruby =  
}
```

```
async function get_page() {  
  let internet_res = await fetch("https://example.com"); // Download the result from the internet  
  return await internet_res.text();  
}  
  
/**  
 * The main function  
 */  
function main() {  
  let my_regex = /(.*)\W+/  
  get_page.then((page_text) => {  
    my_regex.exec(get_page());  
    do_other_stuff();  
  })  
}
```

JS

utils.js

Extraction Process Overview

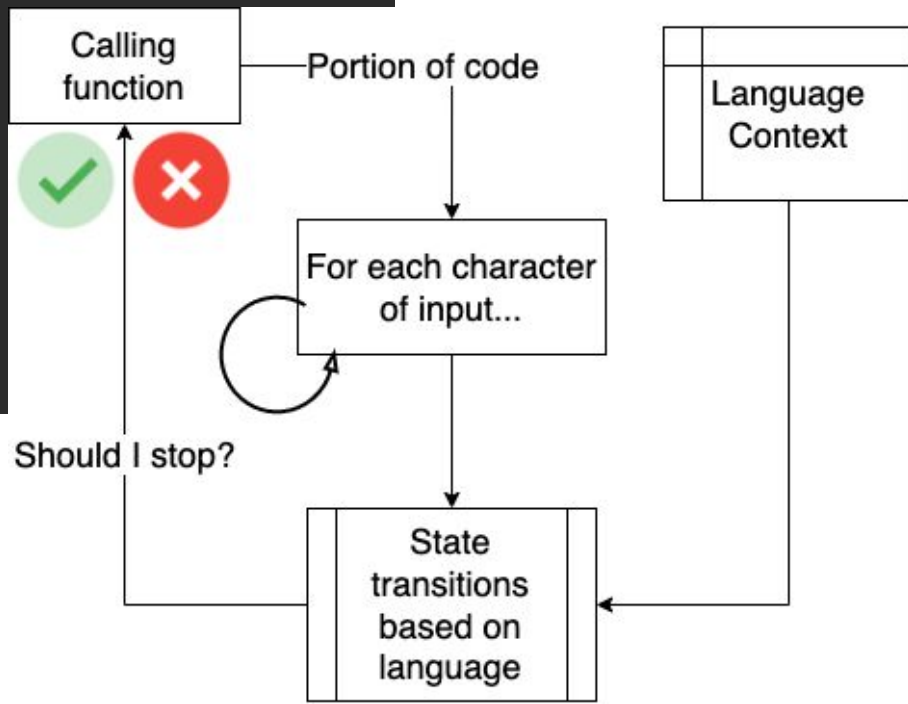


```
async function get_page() {
  let internet_res = await fetch("https://example.com");
  return await internet_res.text();
}

function main() {
  let my_regex = new RegExp("(.*?)" + "\\W+");
  get_page.then((page_text) => {
    my_regex.exec(get_page());
    do_other_stuff();
  })
}
```

General Purpose Tokenizer: A Finite State Machine

```
pub static ref JS_CONTEXT: LanguageContext = LanguageContext {  
  single_line_comment_starts: vec!["//"],  
  multi_line_comment_starts: vec!["/*", "*/"],  
  single_line_string_starts: vec!["'", "\""],  
  multi_line_string_starts: vec!["`"],  
  call_open_sym: b'(',  
  call_close_sym: b')',  
  special_regex_char: Some(b'/')  
};
```



Parsing

```
message Token {  
  string token = 8;  
  TokenType type = 9;  
  
  oneof sub_type {  
    FlagType flag = 10;  
    SubstitutionType substitution = 11;  
    QuantifierModifierType quantifiermodifier = 12;  
    AnchorType anchor = 13;  
    string character = 14;  
    LookaroundType lookaround = 15;  
    EscapeType escape = 16;  
    GroupReferenceType groupref = 17;  
    CharacterClassType characterclass = 18;  
  }  
}
```

```
{  
  token: "g"  
  type: Flag  
  flag: Global  
}  
  
{  
  token: "\t"  
  type: Escape  
  escape: Tab  
}
```

Parsing



```
expression {
  raw: "(.*)\W+"
  tokens {
    token: "("
    type: GroupReference
    groupref: OpenCapture
  }
  tokens {
    token: "."
    type: CharacterClass
    characterclass: Dot
  }
  tokens {
    token: "*"
    type: QuantifierModifier
    quantifiermodifier: Star
  }
  tokens {
    token: ")"
    type: GroupReference
    groupref: CloseCapture
  }
  tokens {
    token: "\\W"
    type: CharacterClass
    characterclass: NotWord
  }
  tokens {
    token: "+"
    type: QuantifierModifier
    quantifiermodifier: Plus
  }
}
```

4. token: ")"
type: GroupReference
groupref: CloseCapture
5. token: "\W"
type: CharacterClass
characterclass: NotWord
6. token: "+"
type: QuantifierModifier
quantifiermodifier: Plus

Understandability Component

- Single bounded class:
 $S\{3\} = \mathbf{SSS} = S\{3,3\}$
- Lower bounded class:
 $A\{2,\} = AAA^* = \mathbf{AA+}$
- Custom Character class:
 $\mathbf{[0-9a]} = [\backslashda] = [0123456789a]$
- Literal class:
 $\mathbf{\backslasha\}$} = \backslasha[\$] = \backslashx61\backslashx24 = \backslash141\backslash044$

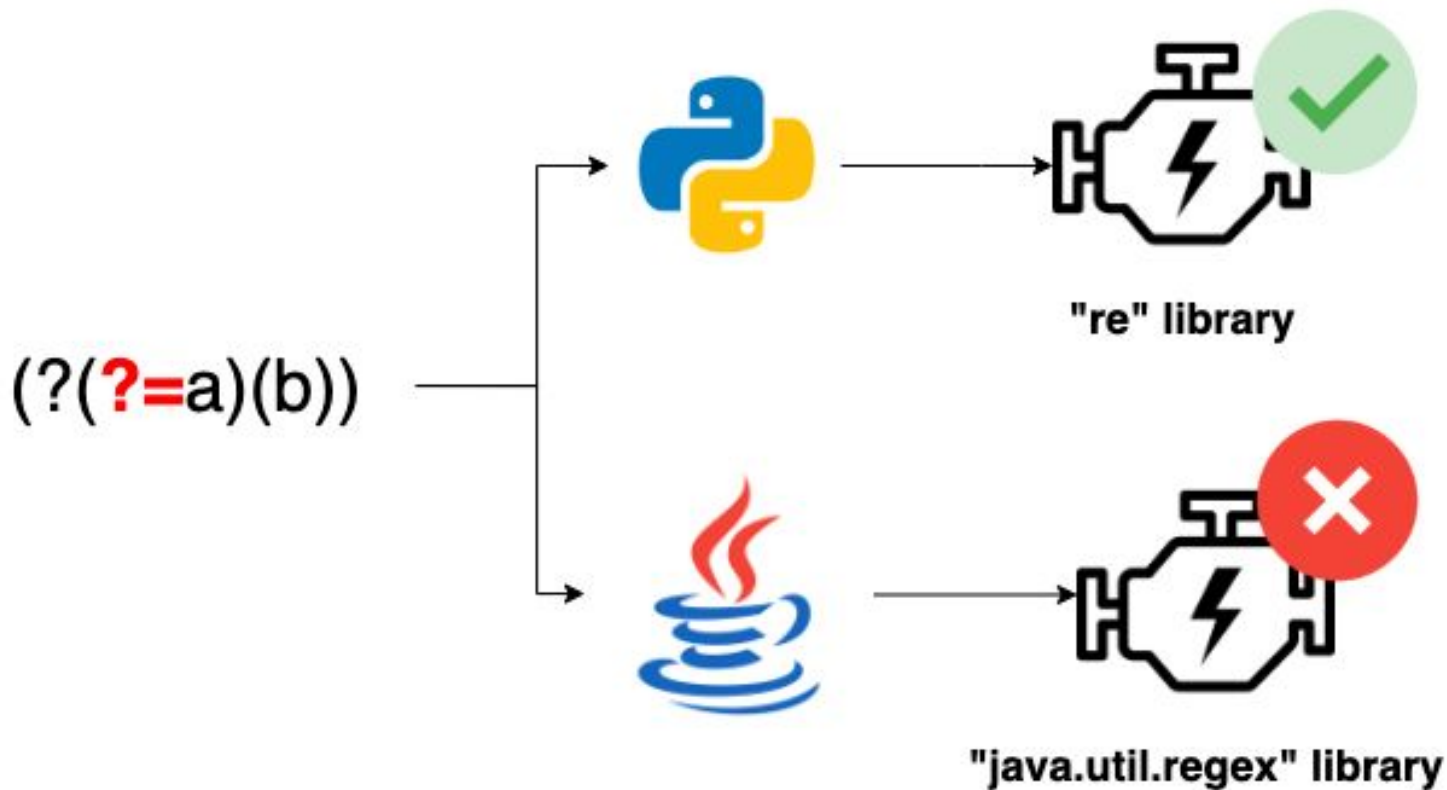
Carl Chapman, Peipei Wang, and Kathryn T. Stolee. 2017. Exploring regular expression comprehension.

Understandability Component

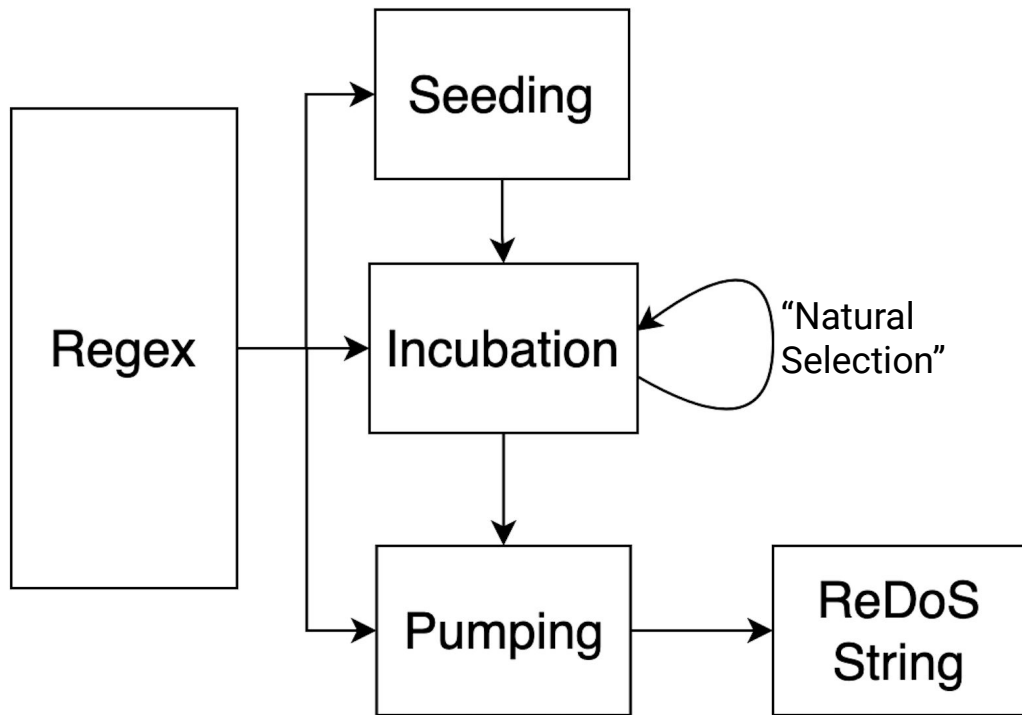
~~\x61[0-9]S{2}~~

a\dSS

Generalizability Component: Regex Engines



ReScue Algorithm



Input: Regex

Output: ReDoS String

- Genetic algorithm for incubation.
- Most desirable strings kept each generation

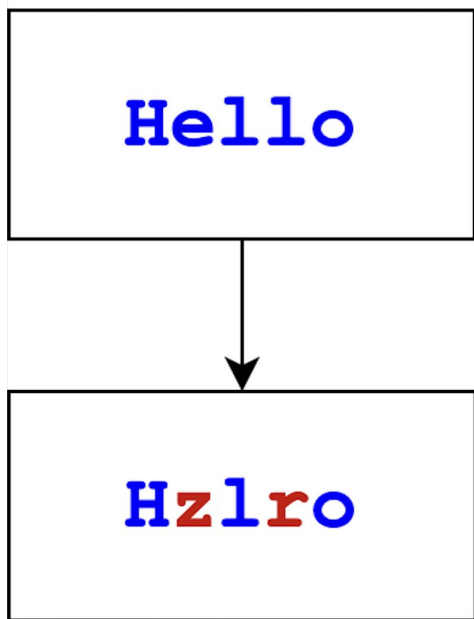
“ReScue: Crafting Regular Expression DoS Attacks”
Yuju Shen et al., Nanjing University

Seeding (Phase 1)

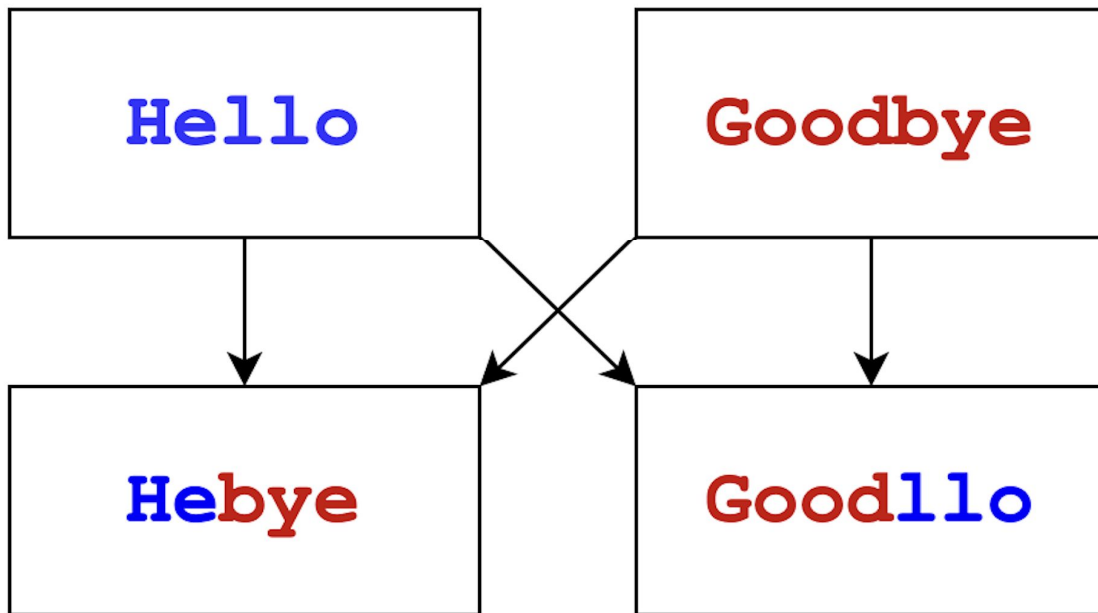
- Regex = "hello.{1,10}goodbye"
- Extract tokens from expression:
 - "hello"
 - "goodbye"
- Make seeds from the tokens, mixing in random characters:
 - "goodbyeYY7goodbye6qL<GWhellogoodbye2>hello"
 - "hellogoodbyehellogoodbyeQhellogoodbye#K(goodbye
\\\\fXhellohello4hellogoodbye"

Genetic Algorithm Mutations (Phase 2)

Mutation:

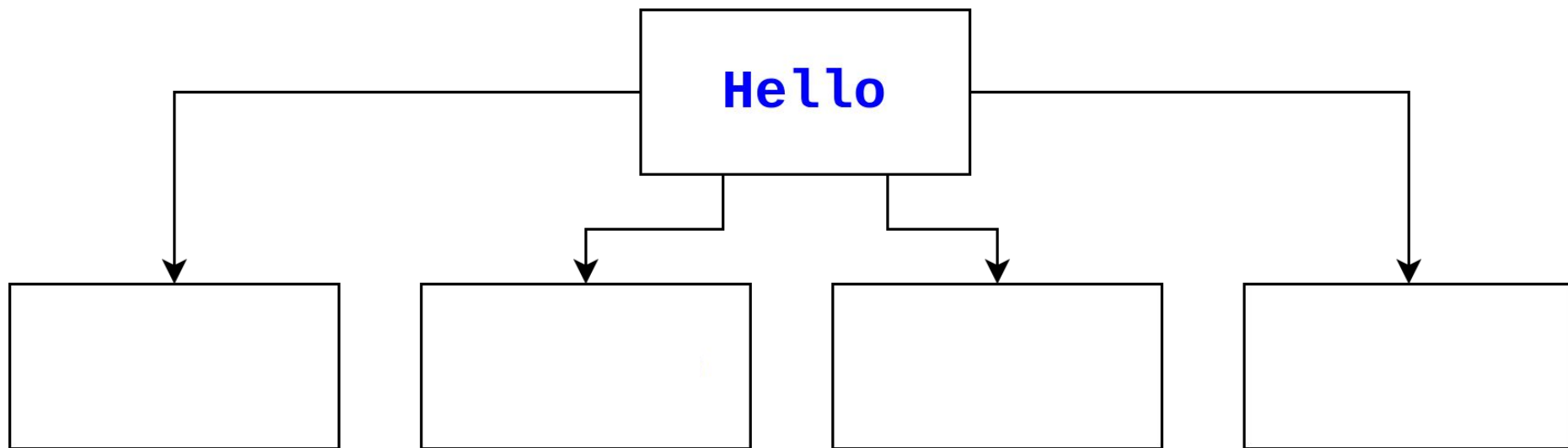


Crossover:

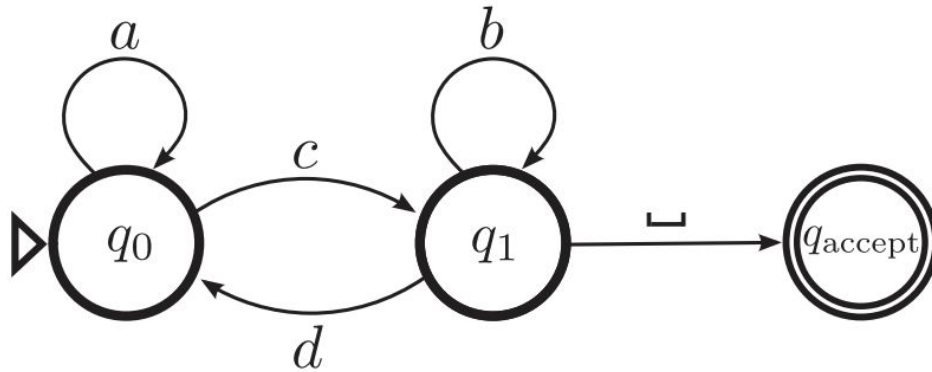


Genetic Algorithm Mutations (Phase 2)

Randomized Mutation:



Regex as Finite State Machines



State diagram for regex:
 $(a^*cb^*d)^* a^*cb^*$

- Any regex can be described as a Finite State Machine
- Match steps are graph dependent
- Implications for fitness functions and timing

Timing as a Fitness Function

- Slower strings = fitter
- Kept fittest strings in each generation
- RESCUE finds fittest strings by calculating matching steps per character
 - Essentially $O(\text{match}(\text{regex}, \text{input}))$
- This requires a regex engine



Timing as a Fitness Function



- Instead we use **time**(match(regex, input))

$$f_{fitness}(s) = \frac{|\tau(s)|}{|s|} \rightarrow f_{fitness}(s) = t(s)$$

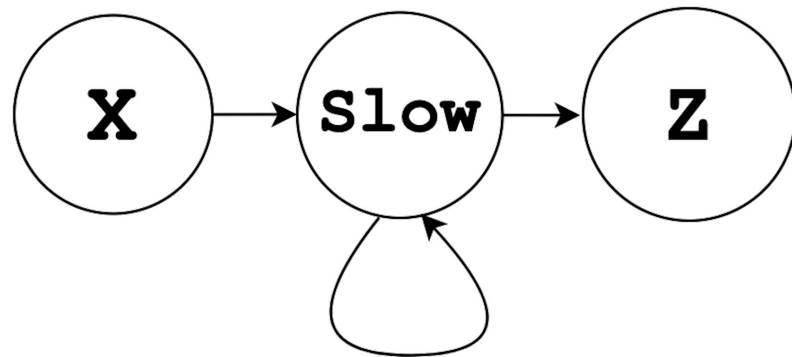
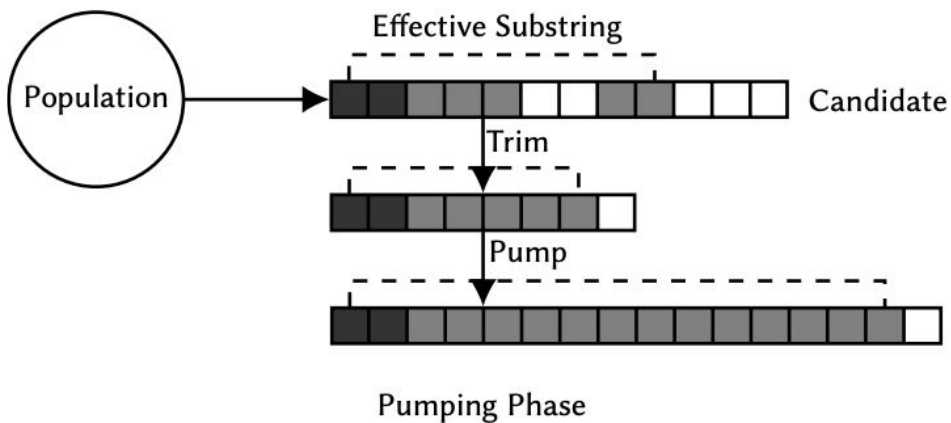
- Matching steps $\tau(s)$ \triangleright process time $t(s)$
- Multiprocessing to allow for timeouts

s = input string

$\tau(s)$ = matching steps of s

$t(s)$ = time to match s

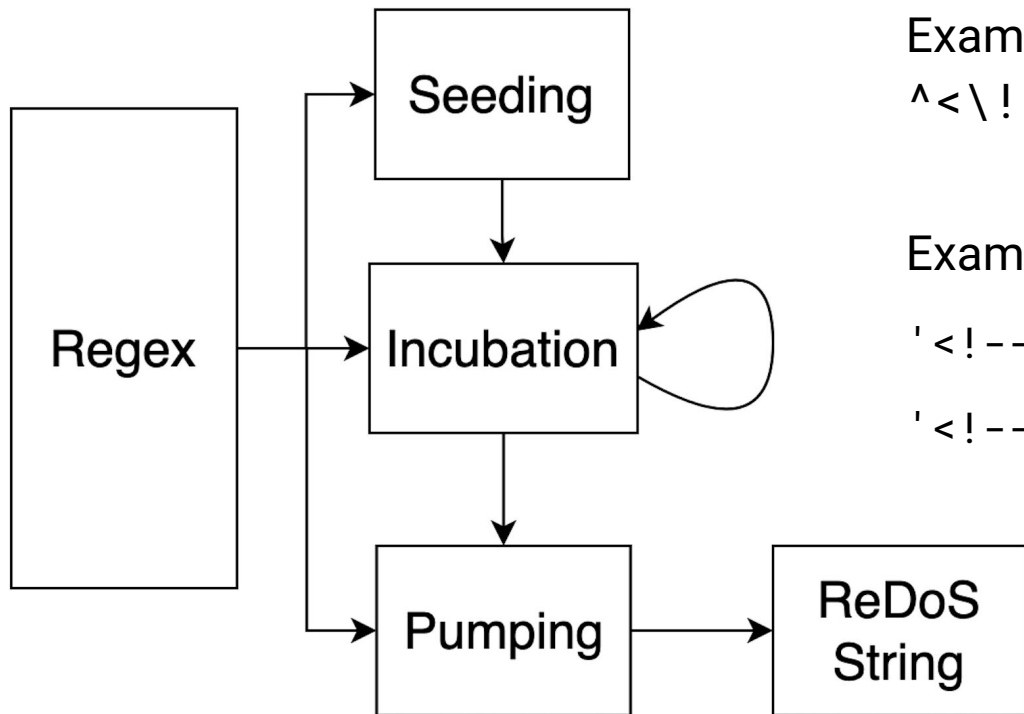
Pumping (Phase 3)



`XslowZ -> XslowslowslowZ`

Pumped output causes ReDoS
attack if used as input

ReScue Algorithm



Example Regex:

```
^<\!\-\- (.*)+(\!\/){0,1}\-\->$
```

Example Output:

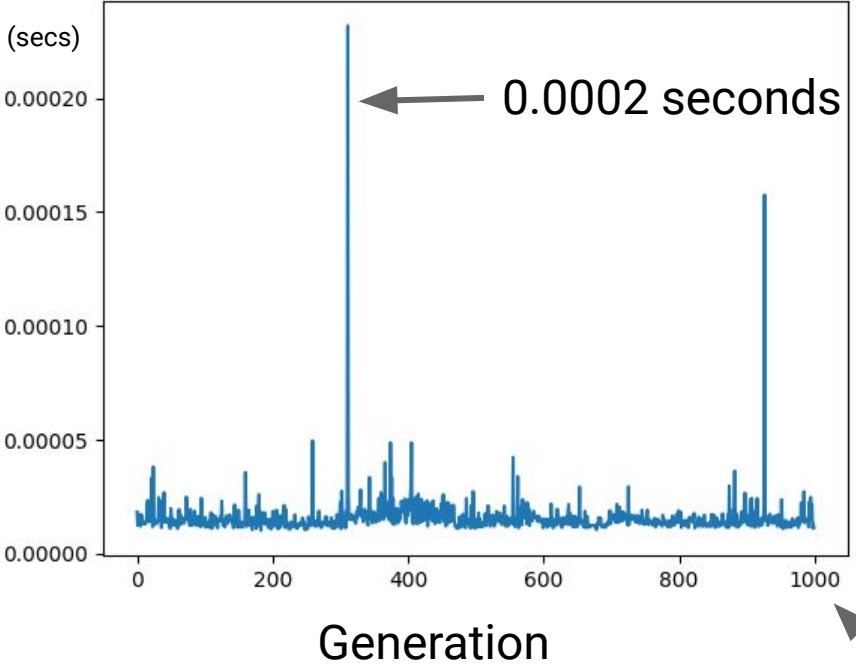
```
'<!-->`!--VvFcwh\hh^b!--^b!--' (6s)
```

```
'<!-->`!--VvFcwh\h^b!--h^b!--^b!--'
```

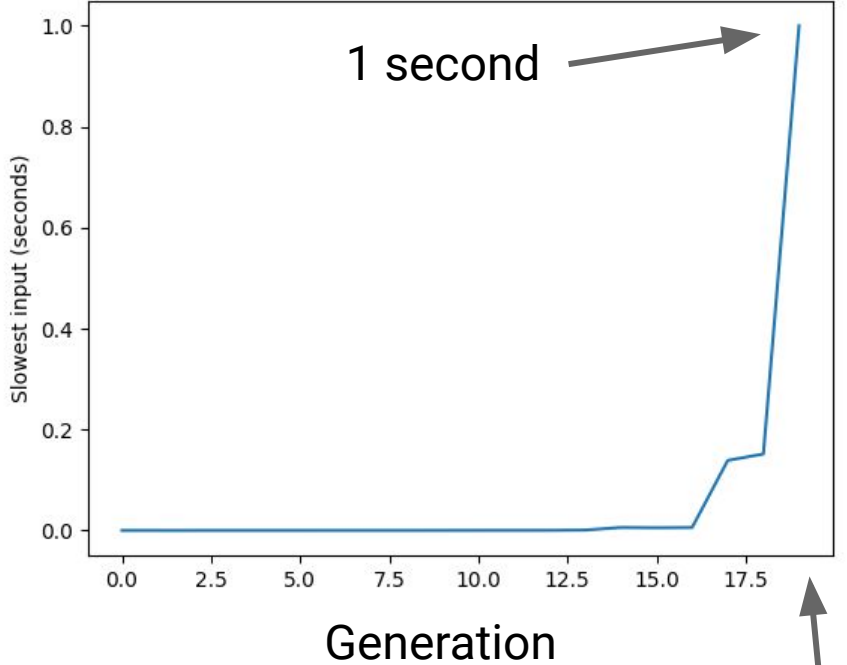
(4m)

`^(?=hello)[a-z]{5}`

`^(([a-z])+.)+[A-Z]([a-z])+$`



1000



18

Evolutionary Tracks

References

- Yuju Shen, Yanyan Jiang, Chang Xu, Ping Yu, Xiaoxing Ma, and Jian Lu. 2018. ReScue: crafting regular expression DoS attacks. Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering. Association for Computing Machinery, New York, NY, USA, 225–235. DOI: <https://doi.org/10.1145/3238147.3238159>
- Carl Chapman, Peipei Wang, and Kathryn T. Stolee. 2017. Exploring regular expression comprehension. In Proceedings of the 32nd IEEE/ACM International Conference on Automated Software Engineering (ASE 2017). IEEE Press, 405–416.

