

Usage and Refactoring Studies Of Python Regular Expressions

Carl Chapman

Iowa State University

carlallenchapman@gmail.com

13 April, 2016

Overview

1 Regex Usage Studies

- Feature Analysis
- Behavioral Clustering
- Developer Survey

2 Regex Refactoring Studies

- Equivalence Model
- Community Support
- Understandability

3 Conclusion

- Refactoring Recommendations
- Future Work
- References

What Features Are Used Most Frequently?

Regex Feature Usage References Are Missing

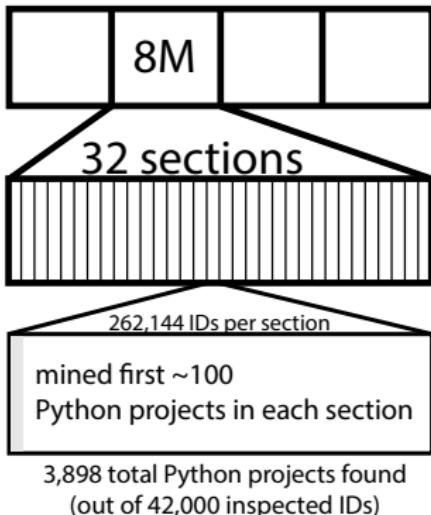
- feature usage statistics
- feature set summaries for variants and tools

Why Python?

Python has most commonly shared features, few rarely shared features

Code	Example	Python	Perl	.Net	Ruby	Java RE2	JavaScript	POSIX ERE
ADD	<code>z+</code>	•	•	•	•	•	•	•
CG	<code>(caught)</code>	•	•	•	•	•	•	•
KLE	<code>.*</code>	•	•	•	•	•	•	•
CCC	<code>[aeiou]</code>	•	•	•	•	•	•	•
ANY	<code>.</code>	•	•	•	•	•	•	•
RNG	<code>[a-z]</code>	•	•	•	•	•	•	•
STR	<code>-</code>	•	•	•	•	•	•	•
END	<code>\$</code>	•	•	•	•	•	•	•
NCCC	<code>['qwxzf']</code>	•	•	•	•	•	•	•
WSP	<code>\s</code>	•	•	•	•	•	•	•
OR	<code>a b</code>	•	•	•	•	•	•	•
DEC	<code>\d</code>	•	•	•	•	•	•	•
WRD	<code>\w</code>	•	•	•	•	•	•	•
QST	<code>z?</code>	•	•	•	•	•	•	•
LZY	<code>z+?</code>	•	•	•	•	•	•	•
NCG	<code>a(?:b)c</code>	•	•	•	•	•	•	•
PNG	<code>(?P<name>x)</code>	•	•	•	•	•	•	•
SNG	<code>x{8}</code>	•	•	•	•	•	•	•
NWSP	<code>\S</code>	•	•	•	•	•	•	•
DBB	<code>z{3,8}</code>	•	•	•	•	•	•	•
NLKA	<code>a(?iyz)</code>	•	•	•	•	•	•	•
WNW	<code>\b</code>	•	•	•	•	•	•	•
NWRD	<code>\W</code>	•	•	•	•	•	•	•
LWB	<code>z{15,}</code>	•	•	•	•	•	•	•
LKA	<code>a(?!bc)</code>	•	•	•	•	•	•	•
OPT	<code>(?i)CanE</code>	•	•	•	•	•	•	•
NLKB	<code>(?<:x)yz</code>	•	•	•	•	•	•	•
LKB	<code>(?<=a)bc</code>	•	•	•	•	•	•	•
ENDZ	<code>\Z</code>	•	•	•	•	•	•	•
BKR	<code>\1</code>	•	•	•	•	•	•	•
NDEC	<code>\d</code>	•	•	•	•	•	•	•
BKRN	<code>(?P<name></code>	•	•	•	•	•	•	•
VWSP	<code>\v</code>	•	•	•	•	•	•	•
NWNW	<code>\B</code>	•	•	•	•	•	•	•

Project Selection



Out of 3,898 pseudo-randomly selected Python projects, 1,645 contained one or more utilization.

Utilizations of the re module

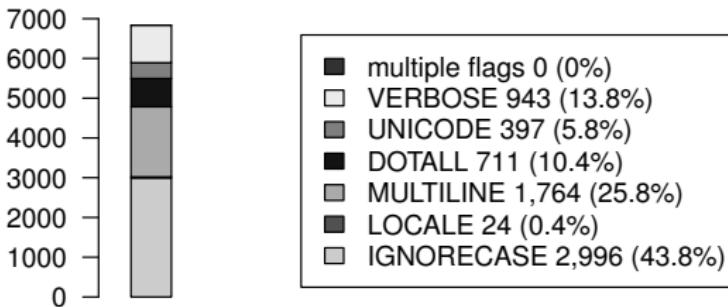
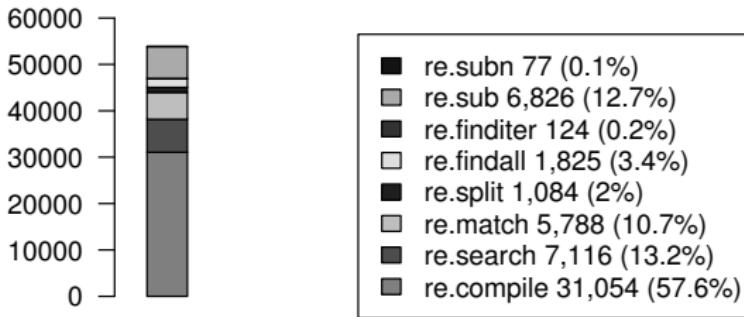
```
function          pattern          flags
r1 = re.compile("(0|-[1-9][0-9]*)$", re.MULTILINE)
```

function which function of the re module is called?

pattern string used to specify regex behavior

`flags` modifies the regex engine

re module insights



Filtering Utilizations And Patterns

53,894 unique utilizations observed.

12.7% use behavioral flags

6.5% were non-static patterns

43,525 utilizations remain

13,711 distinct normalized patterns

73 had unsupported Unicode characters

17 had non-Python features

22 had various errors

2 had ECOM feature - too rare to include

13,597 usable patterns remain for analysis

PCRE Parsing Patterns

$^m + (f(z)^*) +$	→	<table border="1"><tr><td>0</td><td>1</td><td>2</td><td>2</td><td>1</td><td>0</td></tr></table>	0	1	2	2	1	0
0	1	2	2	1	0			
$(ab^*c \mid yz^*) \$$	→	<table border="1"><tr><td>1</td><td>2</td><td>0</td><td>1</td><td>0</td><td>1</td></tr></table> <small>OR KLE ADD CG STR END</small>	1	2	0	1	0	1
1	2	0	1	0	1			

All Python features are recognizable by PCRE

Feature Statistics

Rank	Code	Example	% Projects	NProjects	NFiles	NPatterns	MaxTokens
1	ADD	z+	73.2	1,204	9,165	6,003	30
2	CG	(caught)	72.6	1,194	9,559	7,130	17
3	KLE	.*	66.8	1,099	8,163	6,017	50
4	CCC	[aeiou]	62.4	1,026	7,648	4,468	42
5	ANY	.	61.1	1,005	6,277	4,657	60
6	RNG	[a-z]	51.6	848	5,092	2,631	50
7	STR	^	51.4	846	5,458	3,563	12
8	END	\$	50.3	827	5,393	3,169	12
9	NCCC	[`qxwf]	47.2	776	3,947	1,935	15
10	WSP	\s	46.3	762	4,704	2,846	32
11	OR	a b	43	708	3,926	2,102	15
12	DEC	\d	42.1	692	4,198	2,297	24
13	WRD	\w	39.5	650	2,952	1,430	13
14	QST	z?	39.2	645	3,707	1,871	35
15	LZY	z+?	36.8	605	2,221	1,300	12
16	NCG	a(?:b)c	24.6	404	1,709	791	28
17	PNG	(?P<name>x)	21.5	354	1,475	915	16

Rank	Code	Example	% Projects	NProjects	NFiles	NPatterns	MaxTokens
18	SNG	z{8}	20.7	340	1,267	581	17
19	NWSP	\S	16.4	270	776	484	10
20	DBB	z{3,8}	14.5	238	647	367	11
21	NLKA	a(?!yz)	11.1	183	489	131	3
22	WNW	\b	10.1	166	438	248	36
23	NWRD	\W	10	165	305	94	6
24	LWB	z{15,}	9.6	158	281	91	3
25	LKA	a(?:bc)	9.6	158	358	112	4
26	OPT	(?<1>CasE	9.4	154	377	231	2
27	NLKB	(?<1>x)yz	8.3	137	296	94	4
28	LKB	(?<=>a)bc	7.3	120	255	80	4
29	ENDZ	\Z	5.5	90	149	89	1
30	BKR	\1	5.1	84	129	60	4
31	NDEC	\D	3.5	58	92	36	6
32	BKRN	(P?:=name)	1.7	28	44	17	2
33	VWSP	\v	0.9	15	16	13	2
34	NWNW	\B	0.7	11	11	4	2

Feature Statistics - Top 8

Rank	Code	Example	% Projects	NProjects	NFiles	NPatterns	MaxTokens
1	ADD	z+	73.2	1,204	9,165	6,003	30
2	CG	(caught)	72.6	1,194	9,559	7,130	17
3	KLE	.*	66.8	1,099	8,163	6,017	50
4	CCC	[aeiou]	62.4	1,026	7,648	4,468	42
5	ANY	.	61.1	1,005	6,277	4,657	60
6	RNG	[a-z]	51.6	848	5,092	2,631	50
7	STR	^	51.4	846	5,458	3,563	12
8	END	\$	50.3	827	5,393	3,169	12

Ranked features: Languages

Rank	Code	Example	Python	Perl	.Net	Ruby	Java	RE2	JavaScript	POSIX	ERE
1	ADD	z+	●	●	●	●	●	●	●	●	●
2	CG	(caught)	●	●	●	●	●	●	●	●	●
3	KLE	.*	●	●	●	●	●	●	●	●	●
4	CCC	[aeiou]	●	●	●	●	●	●	●	●	●
5	ANY	.	●	●	●	●	●	●	●	●	●
6	RNG	[a-z]	●	●	●	●	●	●	●	●	●
7	STR	^	●	●	●	●	●	●	●	●	●
8	END	\$	●	●	●	●	●	●	●	●	●
9	NCCC	[~qwxzf]	●	●	●	●	●	●	●	●	●
10	WSP	\s	●	●	●	●	●	●	●	●	●
11	OR	a b	●	●	●	●	●	●	●	●	●
12	DEC	\d	●	●	●	●	●	●	●	●	●
13	WRD	\w	●	●	●	●	●	●	●	●	●
14	QST	z?	●	●	●	●	●	●	●	●	●
15	LZY	z+?	●	●	●	●	●	●	●	●	●
16	NCG	a(?:b)c	●	●	●	●	●	●	●	●	●
17	PNG	(?P<name>x)	●	●	●	●	●	●	●	●	●

Rank	Code	Example	Python	Perl	.Net	Ruby	Java	RE2	JavaScript	POSIX	ERE
18	SNG	z{8}	●	●	●	●	●	●	●	●	●
19	NWSP	\S	●	●	●	●	●	●	●	●	●
20	DBB	z{3,8}	●	●	●	●	●	●	●	●	●
21	NLKA	a(?!yz)	●	●	●	●	●	●	●	●	●
22	WNW	\b	●	●	●	●	●	●	●	●	●
23	NWRD	\w	●	●	●	●	●	●	●	●	●
24	LWB	z{15,}	●	●	●	●	●	●	●	●	●
25	LKA	a(?:bc)	●	●	●	●	●	●	●	●	●
26	OPT	(?i)CasE	●	●	●	●	●	●	●	●	●
27	NLKB	(?<!x)yz	●	●	●	●	●	●	●	●	●
28	LKB	(?<=a)bc	●	●	●	●	●	●	●	●	●
29	ENDZ	\Z	●	●	●	●	●	●	●	●	●
30	BKR	\1	●	●	●	●	●	●	●	●	●
31	NDEC	\D	●	●	●	●	●	●	●	●	●
32	BKRN	(P?=name)	●	●	●	●	●	●	●	●	●
33	VWSP	\v	●	●	●	●	●	●	●	●	●
34	NWNW	\B	●	●	●	●	●	●	●	●	●

Ranked features: Languages - Notable Missing Features

Rank	Code	Example	Python	Perl	.Net	Ruby	Java	RE2	JavaScript	POSIX	ERE
21	NLKA	a(?!yz)	●	●	●	●	●	●	●	●	●
22	WNW	\b	●	●	●	●	●	●	●	●	●
23	NWRD	\W	●	●	●	●	●	●	●	●	●
24	LWB	z{15,}	●	●	●	●	●	●	●	●	●
25	LKA	a(?=bc)	●	●	●	●	●	●	●	●	●
26	OPT	(?i)CasE	●	●	●	●	●	●	●	●	●
27	NLKB	(?<!x)yz	●	●	●	●	●	●	●	●	●
28	LKB	(?<=a)bc	●	●	●	●	●	●	●	●	●
29	ENDZ	\Z	●	●	●	●	●	●	●	●	●
30	BKR	\1	●	●	●	●	●	●	●	●	●
31	NDEC	\D	●	●	●	●	●	●	●	●	●

Ranked features: Analysis Tools

Rank	Code	Example	Brics	Hampi	Rex	Automata.Z3
1	ADD	z+	●	●	●	●
2	CG	(caught)	●	●	●	●
3	KLE	.*	●	●	●	●
4	CCC	[aeiou]	●	●	●	●
5	ANY	.	●	●	●	●
6	RNG	[a-z]	●	●	●	●
7	STR	^	●	●	●	●
8	END	\$	●	●	●	●
9	NCCC	[^qwxrf]	●	●	●	●
10	WSP	\s	●	●	●	●
11	OR	a b	●	●	●	●
12	DEC	\d	●	●	●	●
13	WRD	\w	●	●	●	●
14	QST	z?	●	●	●	●
15	LZY	z+?	●	●	●	●
16	NCG	a(?:b)c	●	●	●	●
17	PNG	(?P<name>x)	●	●	●	●

Rank	Code	Example	Brics	Hampi	Rex	Automata.Z3
18	SNG	z{8}	●	●	●	●
19	NWSP	\S	●	●	●	●
20	DBB	z{3,8}	●	●	●	●
21	NLKA	a(?!yz)	●	●	●	●
22	WNW	\b	●	●	●	●
23	NWRD	\W	●	●	●	●
24	LWB	z{15,}	●	●	●	●
25	LKA	a(?=bc)	●	●	●	●
26	OPT	(?i)CasE	●	●	●	●
27	NLKB	(?<!x)yz	●	●	●	●
28	LKB	(?<=a)bc	●	●	●	●
29	ENDZ	\Z	●	●	●	●
30	BKR	\1	●	●	●	●
31	NDEC	\D	●	●	●	●
32	BKRN	\g<name>	●	●	●	●
33	VWSP	\v	●	●	●	●
34	NWNW	\B	●	●	●	●

What Are Regexes Used For?

Non-Anecdotal Knowledge About Usage Is Missing

- task categories
- behavioral categories



How to Categorize Regex Usages

- ① thorough inspection of 55K utilizations
- ② unguided manual categorization of 4,694 regexes (in 2 or more projects), without objective basis
- ③ cluster by syntactic similarity like Jaccard or longest substring
- ④ formal analytical subsumption, using hampi (94%?) cannot get it to work
- ⑤ formal analytical subsumption, using brics (30% or less)
- ⑥ Chosen technique: cluster by behavioral similarity using Rex (61%)

Measuring Behavioral Similarity

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

	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



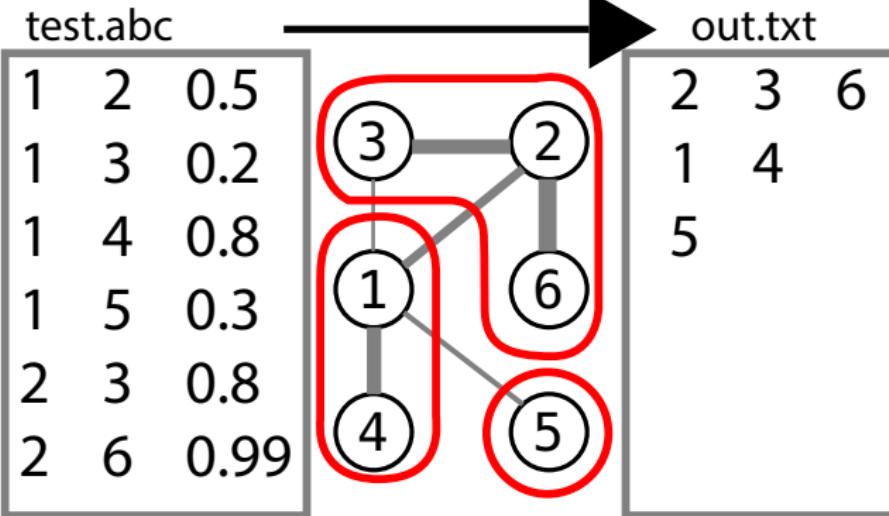
	A	B	C	D
A	1.0			
B	0.1	1.0		
C	0.75	0.8	1.0	
D	0.0	0.65	0.15	1.0

Rex generates 400 strings for each regex.

Convert scores to half-matrix to make *.abc file for mcl.

MCL example

```
>mcl test.abc --abc -o out.txt
```



Example Cluster

Index	Pattern	NProjects	Index	Pattern	NProjects
1	\s*([^\n:]*)\s*:.*	9	7	[:]	6
2	:+	8	8	([^\n:]+):.*	6
3	(:)	8	9	\s*:\s*	4
4	(:+)	8	10	\:	2
5	(:)(:*)	8	11	^([^\n:]*):[^\$]*\$	2
6	^([^\n:]*):.*	8	12	^[:]*:([^\$]*)\$	2

Six Categories Of Clusters

Category	Clusters	Patterns	Projects	% Projects
Multi Matches	21	237	295	40%
Specific Char	17	103	184	25%
Anchored Patterns	20	85	141	19%
Two or More Chars	16	40	120	16%
Content of Parens	10	46	111	15%
Code Search	15	27	92	13%

Code Search `.*rlen=([0-9]+)`

Two Or More Chars `@[a-z]+`

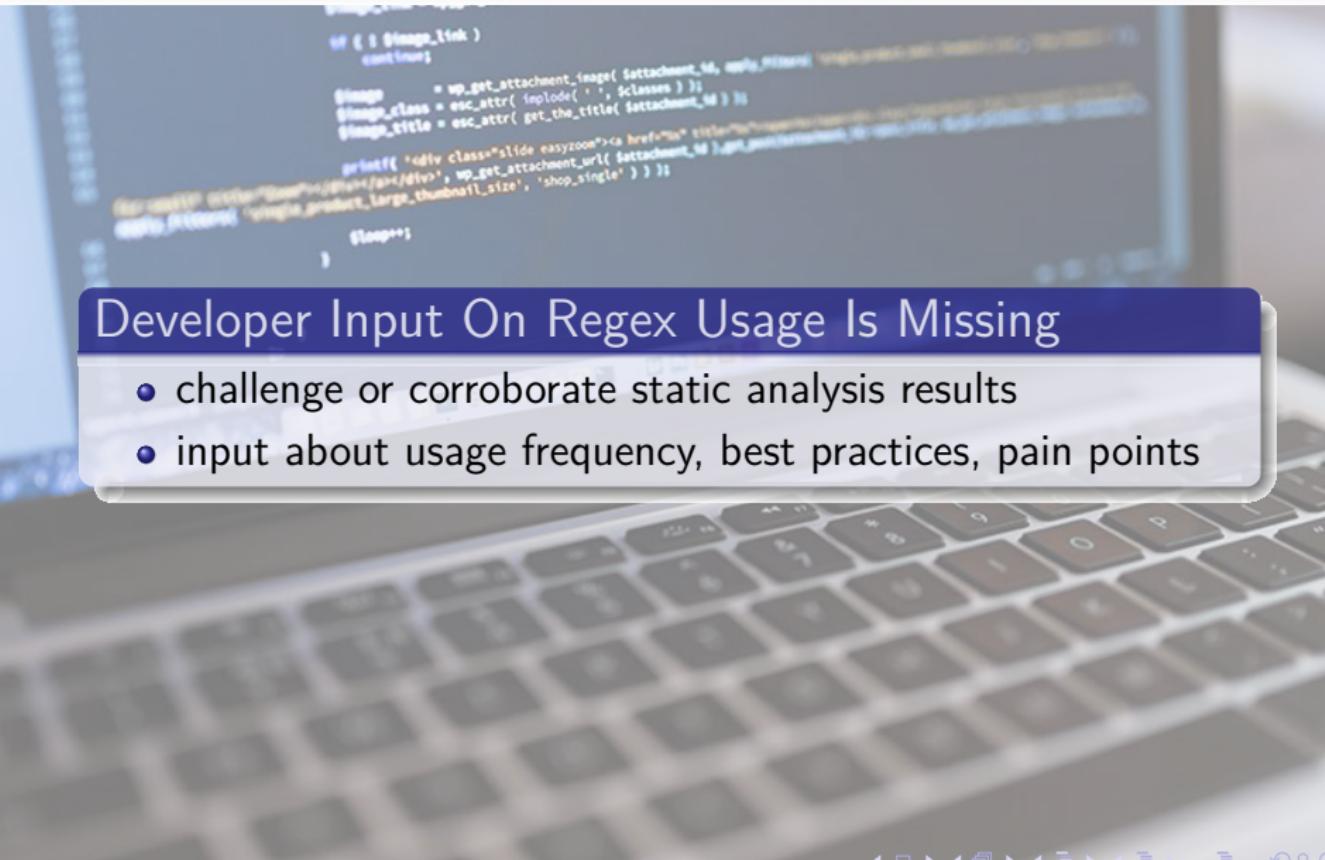
Bracket Parsing `<(.+)>`, `<[^>]*?>`

Specific Char `:+ , } , %`

Anchored `^[-_A-Za-z0-9]+$`

Multiple Alternatives `(\s) , , | ;`

How Do Developers Say They Use Regexes?



```
if (! $image_link )
    continue;
$image = wp_get_attachment_image( $attachment_id, apply_filters( 'single_product_gallery_size', 'shop_single' ) );
$image_class = esc_attr( implode( ' ', $classes ) );
$image_title = esc_attr( get_the_title( $attachment_id ) );
printf( 'div class="slide easyzoom">><a href="%s" title="%s" class="%s" data-zoom-image="%s"></a></div>', wp_get_attachment_url( $attachment_id ), $image_title, $image_class, $image->src, $image->alt, $image->src );
```

Developer Input On Regex Usage Is Missing

- challenge or corroborate static analysis results
- input about usage frequency, best practices, pain points

Feature Usage Is Consistent With Analysis

	CG	STR or END	LZY	WNW	look-arounds
very frequently	2	1	0	1	0
frequently	4	9	2	3	1
occasionally	9	5	6	6	2
rarely	2	2	2	2	5
very rarely	1	1	4	6	7
never	0	0	4	0	3
avg	5.8	6.1	4	4.8	3.5

Ranked Order: CG (2), STR/END (7,8), LZY (15), WNW (22),
 look-arounds (21, 25, 27, 28)

Task Frequencies are Mostly Consistent With Behavioral Categories

	Capturing	Counting Lines	Counting All	Finding	Filtering	Single Char	Parse User Input	Parse Generated	Other
v. freq	1	1	1	3	0	0	2	2	0
freq.	9	2	3	7	1	0	5	1	1
occ.	3	5	4	3	8	1	5	4	0
rarely	5	3	3	4	2	3	3	3	0
v. rarely	0	3	4	1	5	5	3	5	1
never	0	4	3	0	2	9	0	3	16
avg	3.3	2.0	2.2	3.4	2.1	0.8	3	2.1	0.3

Developers said they did not frequently search for a single character.

Regex Testing

	Always	V. Freq	Freq.	Occ.	Rarely	V. Rarely	Never
test code	4	7	5	1	0	0	1
test regex	3	4	5	5	1	1	0

The screenshot shows the regex101.com web application interface. On the left, there's a sidebar with buttons for 'SAVE & SH...', 'FLAVOR' (set to PCRE), 'TEST STRING' (containing 'abbbbbbcde'), and 'TOOLS'. The main area has tabs for 'REGULAR EXPRESSION' (selected) and 'EXPLANATION'. The regular expression input field contains '/(ab*)(cd)*e?/'. The explanation panel details the match: it finds 1 match in 11 steps. The match is explained as follows:

- 1st Capturing group (ab*)
- a matches the character a literally (case sensitive)
- b* matches the character b literally

The 'MATCH INFORMATION' panel shows two matches found in the test string:

- MATCH 1**
 - 1. [0-7] `abb`
bbbb
 - 2. [7-9] `cd`

50% say they use testing tools like www.regex101.com

Usage Frequency - By Technical Environment

Heaviest regex use is in Command line tools and text editors.

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

Ephemeral vs Persistent Users

Task	Persistence Freq.	Ephemeral Freq.	Difference
Counting substrings that match a pattern	3	1.7	1.2
Parsing user input	3.6	2.7	0.9
Capturing parts of strings	3.8	3.1	0.7
Parsing generated text	2.4	1.9	0.5
Locating content within a file or files	3.6	3.2	0.4
Filtering collections (lists, tables, etc.)	2.2	1.9	0.3
Counting lines that match a pattern	1.8	2.1	-0.3

Code	Persistent Freq.	Ephemeral Freq.	Difference
LKA, NLKA, LKB, NLKB	3.2	2.2	1.0
LZY	3	2.8	0.2
STR, END	4.4	4.4	0
CG	4.2	4.2	0
WNW	3.4	3.5	-0.1

Pain Points

hard to compose (11)

...very difficult to write them since I've never read up on them.

...trickiness to getting the expression right

hard to read (7)

long ones can be hard to read

Readability. Edge cases.

It is terrible to read (especially later after initial development)

inconsistency across implementations (3)

Differences in implementation across languages

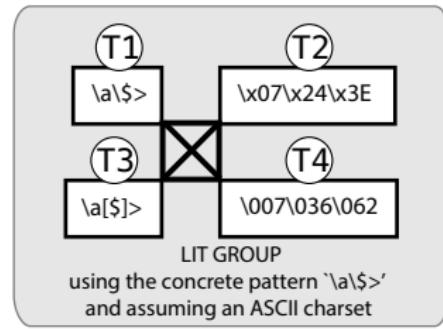
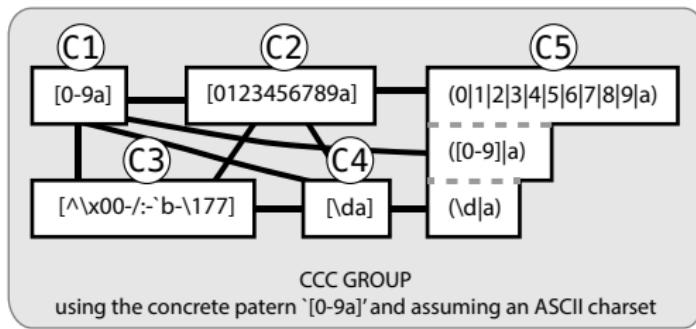
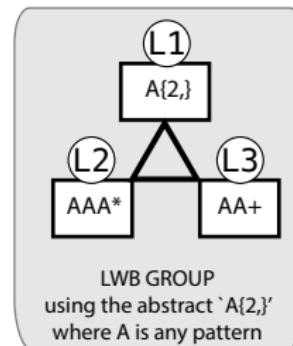
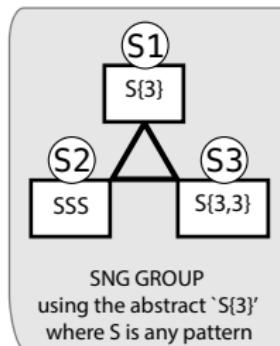
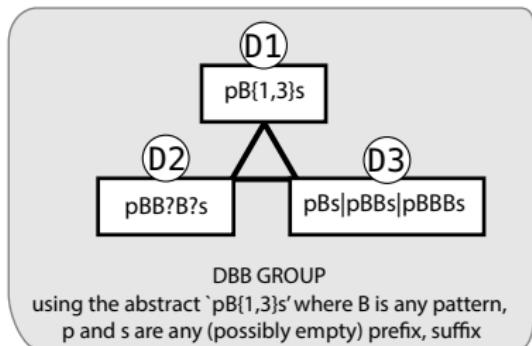
Some regexes work differently (or don't work) in some languages.

What Different Ways Can A Regex Be Represented?

Regex Equivalence Model Is Missing

- regex refactoring is new
- equivalence model required for refactoring

Regex Equivalence Classes



Example Equivalences

LIT : `[\072\073]` \equiv `[{}{)}`

DBB : `((q4f)?ab)` \equiv `(q4fab|ab)`

CCC : `tri[a-f]3` \equiv `tri(a|b|c|d|e|f)3`

LWB : `zaa*` \equiv `za+`

SNG : `[aeiou]{2}` \equiv `[aeiou] [aeiou]`

How to decide which representation is preferred?

How Does The Community Prefer To Represent Regexes?

Knowledge Of Community Standards Is Missing

- programmers probably choose the best representation
- conformance to community standards eases comprehension

Another Point

What Representations Are More Understandable?

Knowledge Of Regex Understandability Is Missing

- understandability is a major pain point
- comprehension tests can determine understandability

Another Point

Citation

An example of the \cite command to cite within the presentation:

This statement requires citation Smith (2012).

Verbatim

Example (Theorem Slide Code)

```
\begin{frame}
\frametitle{Theorem}
\begin{theorem}[Mass--energy equivalence]
$E = mc^2$
\end{theorem}
\end{frame}
```

regex formatting test

ab*c

References

John Smith (2012) Title of the publication Journal Name 12(3), 45 – 678.

Questions?