

Writing Nice Regular Expressions

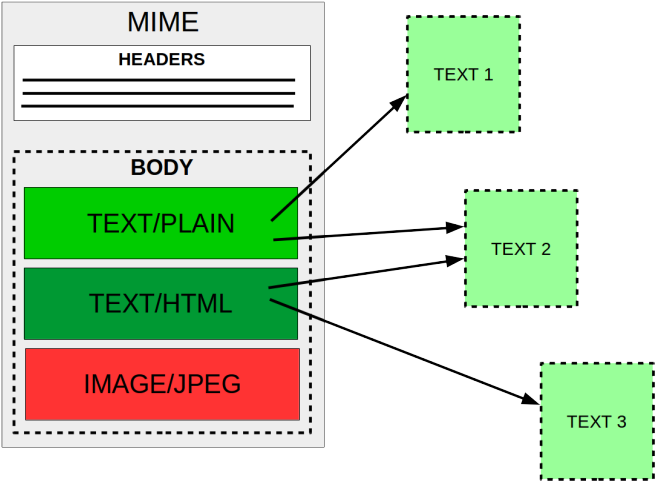
Dmitry Shiltsov

April 14, 2015

Me and regular expressions

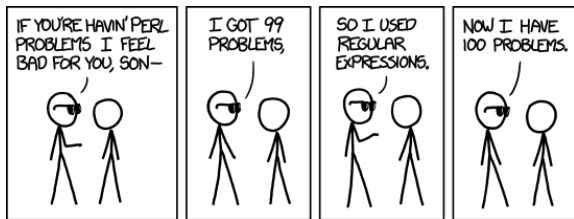


Kaspersky Anti-Spam (2006 – 2014)



Readability vs performance vs coding time

- ▶ Simple regexps
- ▶ Optimized regexps
- ▶ Not using regexps



Regular expressions in Perl

```
1  m/PATTERN/;  
2  s/PATTERN/REPLACEMENT/;  
3  split /PATTERN/, EXPR;  
4  qr/STRING/;
```

Matching phone numbers

Wellington region land line:

'(04) 529-9999',

'(04) 239-9999',

'(04) 569-9999',

'(04) 479-9999',

'(04) 389-9999',

Matching phone numbers: regexp 1

Wellington region land line:

'(04) 529-9999',

'(04) 239-9999',

'(04) 569-9999',

'(04) 479-9999',

'(04) 389-9999',

```
1 qr/\(04\) \d{3}-\d{4}\b/;
```

Matching phone numbers: regexp 1

Wellington region land line:

'(04) 529-9999',

'(04) 239-9999',

'(04) 569-9999',

'(04) 479-9999',

'(04) 389-9999',

'+64-4-123 4567',

'64 4 1234 567',

'(04) 1234 401',

```
1 qr/\(04\) \d{3}-\d{4}\b/;
```

Matching phone numbers: regexp 2

Wellington region land line:

'(04) 529-9999',
'(04) 239-9999',
'(04) 569-9999',
'(04) 479-9999',
'(04) 389-9999',
'+64-4-123 4567',
'64 4 1234 567',
'(04) 1234 401',

```
1 qr/(?:(04\ ) |\+64-4-)\d{3}[ -]\d{4}\b/;
```


Matching phone numbers: regexp 3

Wellington region land line:

'(04) 529-9999',
'(04) 239-9999',
'(04) 569-9999',
'(04) 479-9999',
'(04) 389-9999',
'+64-4-123 4567',
'64 4 1234 567',
'(04) 1234 401',

```
1 qr/(?:(04\ ) |\+64-4-)(?:\d{3}[ -]\d{4}|\d{4}[ -]\d{3})\b/;
```

Matching phone numbers: regexp 4

Wellington region land line:

'(04) 529-9999',
'(04) 239-9999',
'(04) 569-9999',
'(04) 479-9999',
'(04) 389-9999',
'+64-4-123 4567',
'64 4 1234 567',
'(04) 1234 401',

```
1 qr/(?:(?:\d{3}\d{4}|\d{4}\d{3})|(\d{4})\d{3})\b/;
```

Matching phone numbers: regexp 4

Wellington *and* Auckland region land line:

'(04) 529-9999',

'(04) 239-9999',

'(04) 569-9999',

'(04) 479-9999',

'(04) 389-9999',

'+64-4-123 4567',

'64 4 1234 567',

'(04) 1234 401',

'09 123 4567',

'(+64 9) 123 4567',

```
qr/(?:(04\\)|\\+?\\b64[ -]4)[ -](?:\\d{3}[ -]\\d{4}|\\d{4}[ -]\\d{3})\\b/;
```

Matching phone numbers: regexp 5

Wellington and Auckland region land line:

'(04) 529-9999',
'(04) 239-9999',
'(04) 569-9999',
'(04) 479-9999',
'(04) 389-9999',
'+64-4-123 4567',
'64 4 1234 567',
'(04) 1234 401',
'09 123 4567',
'(+64 9) 123 4567',

```
qr/(?:(?(\b0[49])?|\+?\b64[ -][49])[-](?:\d{3}[-]\d{4}|\d{4}[-]\d{3}))\b/;
```

Matching phone numbers: regexp 6

Wellington and Auckland region land line:

```
'(04) 529-9999',  
'(04) 239-9999',  
'(04) 569-9999',  
'(04) 479-9999',  
'(04) 389-9999',  
'+64-4-123 4567',  
'64 4 1234 567',  
'(04) 1234 401',  
'09 123 4567',  
'(+64 9) 123 4567',
```

```
1 qr/\((?:0[49]|\+?\b64[ -][49])\)?[ -](?:\d{3}[ -]\d{4}|\d{4}[ -]\d{3})\b/;
```

Matching phone numbers: regexp 7

Wellington and Auckland region land line:

'(04) 529-9999',
'(04) 239-9999',
'(04) 569-9999',
'(04) 479-9999',
'(04) 389-9999',
'+64-4-123 4567',
'64 4 1234 567',
'(04) 1234 401',
'09 123 4567',
'(+64 9) 123 4567',

```
qr/\((?:\+?\b64[ -]?\|\b0)[49]\)?[ -]?(?:\d{3}[ -]?\d{4}|\d{4}[ -]?\d{3})\b/;
```

Matching phone numbers: free formatting and comments

```
1 qr/\(?(?:\+?\b64[ -]?|\b0)[49]\)?[ -]?(?:\d{3}[ -]?\d{4}|\d{4}[ -]?\d{3})\b/;
```

```
1 qr/  
2     \(?  
3         (?:\+?\b64[ -]?|\b0) # prefix +64 or 0  
4         [49]                # 4 - Well, 9 - Auck  
5     \)?[ -]?  
6     (?:  
7         \d{3}[ -]?\d{4} # numbers like xxx-xxxx  
8         |  
9         \d{4}[ -]?\d{3} # numbers like xxxx-xxx  
10    )  
11    \b/x;
```

```
1 qr/(Abc(?:# Another way to put comments )Xyz)/;
```

Matching phone numbers: with variables

```
1 qr/\((?:\+?\b64[ -]?!\\b0)[49]\\)?[ -]?(?:\d{3}[ -]?\\d{4}|\\d{4}[ -]?\\d{3})\\b/;
```

```
1 my $sep = qr/[ -]+/;
2 my $land_prefix = qr/[49]/; # 4 - Well, 9 - Auck
3 my $prefix = qr/
4     \+?\b64 $sep $land_prefix # international prefix
5     |
6     \b0$land_prefix/x; # local prefix
7 my $main_number = qr/
8     \d{3} $sep \d{4} # numbers like xxx-xxxx
9     |
10    \d{4} $sep \d{3}/x; # numbers like xxxx-xxx
11
12 my $regexp= qr/\((? $prefix \)? $sep $main_number \b/x;
```


Extracting substrings

```
1 my $date = 'Tue Mar 10 18:57:40 NZDT 2015';
2
3 my $date_pattern = qr /~
4     ([A-Z][a-z][a-z])\s # weekday
5     ([A-Z][a-z][a-z])\s # month
6     (\d{1,2})\s         # day
7     (\d\d:\d\d:\d\d)\s # time
8     (\S+)\s            # tz
9     (\d{4})/x;         # year
10
11 if ($date =~ $date_pattern) {
12     say "$3 $2 $6";
13 }
```

```
1 10 Mar 2015
```

Extracting substrings: named captures

```
1 my $date = 'Tue Mar 10 18:57:40 NZDT 2015';
2
3 my $date_pattern = qr /~
4     (?<weekday>[A-Z][a-z][a-z])\s
5     (?<month>[A-Z][a-z][a-z])\s
6     (?<day>\d{1,2})\s
7     (?<time>\d\d:\d\d:\d\d)\s
8     (?<tz>\S+)\s
9     (?<year>\d{4})/x;
10
11 if ($date =~ $date_pattern) {
12     say "${day} ${month} ${year}";
13 }
```

```
1 10 Mar 2015
```

Some regexp modifiers

```
1 /PATTERN/msixgo
```

m treat string as multiple lines.

s treat string as single line

i case-insensitive

x free formatting

o compile once

g iteration

```
1 /Abc (?i-sm)case insensitive here(?-i) Xyz /ms;  
2 /Abc (?i-sm:case insensitive here) Xyz /ms;
```

Email: From matches To

```
1 my $headers =
2 'Received: from somewhere ([123.123.123.123])
3     by someone-else with lots of details...
4 From: Canadian Pharmacy <doctor@superpharm.info>
5 Reply-To: Fake Name <email@server.com>
6 To: Another Name <doctor@superpharm.info>
7 Subject: Enlarge your pencil!
8 Content-Type: text/plain';
```

```
1 my $email = qr/<[\w.-]+@[ \w\.-]+\.[a-z]{2,4}>/;
2
3 my $from_in_to = qr/
4     \nFrom: .* ($email)\r?\n    # capturing from-address
5     (? : .* \r?\n)*          # any number of lines
6     To: .* \1                # To matches From
7 /x;
```

Email: From matches To (with /s and /m)

```
1 my $headers =
2 'Received: from somewhere ([123.123.123.123])
3     by someone-else with lots of details...
4 From: Canadian Pharmacy <doctor@superpharm.info>
5 Reply-To: Fake Name <email@server.com>
6 To: Another Name <doctor@superpharm.info>
7 Subject: Enlarge your pencil!
8 Content-Type: text/plain';
```

```
1 my $email = qr/<[\w.-]+@[ \w\.-]+\.[a-z]{2,4}>/;
2
3 my $from_in_to = qr/
4     ^From: \N* ($email)$ # capturing from-address
5     .*
6     ^To: \N* \1$/xms; # To matches From
```

Dealing with slashes

```
1 /PATTERN/  
2 m/PATTERN/  
3 m!PATTERN!
```

```
1 s/PATTERN/REPLACEMENT/  
2 s{PATTERN}[REPLACEMENT]
```

Examples:

```
1 m/^\\home\\/dmitry\\/Desktop\\/.*pdf$/  
2 m\\/home/dmitry/Desktop/. *pdf$\\
```

Beware of a special case!

```
1 m?PATTERN?
```

Literal search (quotemeta)

```
1 my $code = '... \d{1,3}\.\d{1,3}\.\d{1,3}\.\d{1,3}...'
;
2 $ip =~ m/\\d\\{1,3}\\\\.\\.\\d\\{1,3}\\\\.\\.\\d\\{1,3}\\\\.\\.\\d\\{1,3}\\\\.\\.\\d\\{1,3}\\/
3 $ip =~ m/\Q\d{1,3}\.\d{1,3}\.\d{1,3}\.\d{1,3}\E/
```

```
1 my $perl_line = 'say "${day} ${month} ${year}";';
2 my $pattern = qr/^\s*\Q$say_line\E$/m;
```