php|works Thursday 11:30 - 13:00

Andrei's Regex Clinic



Thursday 11:30 - 13:00

Andrei's Regex Clinic

Download this presentation from http://www.gravitonic.com/talks/







Thursday 11:30 - 13:00

Andrei's Regex Clinic

Download this presentation from http://www.gravitonic.com/talks/







Thursday 11:30 - 13:00

Andrei's Regex Clinic

Download this presentation from http://www.gravitonic.com/talks/







-(-

Thursday 11:30 - 13:00

Andrei's Regex Clinic

Download this presentation from http://www.gravitonic.com/talks/







about me

- PHP core developer since 1999
- Infrastructure software engineer at Yahoo! Inc.
- Email: andrei@gravitonic.com

 $\uparrow \uparrow \uparrow \uparrow \uparrow \uparrow \uparrow \uparrow \uparrow$

what's the plan?





-



regular... expressions

-\$

-



 \oplus









Regular expressions can be traced back to early research on the human nervous system





15

Neurophysiologists Warren McCulloch and Walter Pitts developed a mathematical way of describing the neural networks



Later, mathematician Stephen Kleene published a paper that introduced the concept of regular expressions that were used to describe "the algebra of regular sets"



grep "/[A-Z]"

17

Subsequently, Ken Thompson, one of the fathers of Unix, found a practical application for them in the various tools of the early OS



what's it good for?

- Literal string searches are fast but inflexible
- **With regular expressions you can:**
 - Find out whether a certain pattern occurs in the text
 - Locate strings matching a pattern and remove them or replace them with something else
 - Extract the strings matching the pattern



and the statute for some or where, dow

and the second s

And the second s

A stretch given on a particular a movery particular to move the perturbance a part of state of the stretcher perturbance to base or allowed to a A compression of the to base of the state of the state of the stretcher to base of the state of the st

A side (if rW) and Solidarovy, exercit, manually, it and bath yes, a support of and bath yes, a support of a support of the support a support of a support

and herry of a New Jackson from the state of the

Apple on the set of standing other, apple Apple on the standing other standing apple of the standing of the standing apple of the standing of

Annual of level and a fee Permanent

and with standard & particular in the same

and the form of the formation of the for

convertige there are no to trady and the well and the second state of the second state

alpenase of band i see a second and i see a second and i second a second a

Regex

a pattern describing a set of characters

a b c d e f

apple

Subject String

text that the regex is applied to

January a Transmiss or aller, day |

and the second set of the seco

and the set of about a pro-

The second second second second

A strain of the second second

and a second solid with other states of a second solid solid line bases manufactures in the advantage. In the second solid line bases with associative and as to advant in with support in solid

all school (A rel) and Analy in the second of an and a second of a second of any in the second of a s

and here's of a West Indone from the send here's of a West Indone from a redde Line, in the constant of the bottom settle indone for the constant of the bottom date with it of the prior has been

Are a first of the band of the approximation of the second second

Married in fairs and adds. Permanning

a sub-sample by treater in many

A state and a design of the state of the sta

atoms in him ? why only the structure of a second state ? and a second state ? Without a second

of the photon of the photon of the state of

and the state state of the stat

the product of their second are like the bases, of pakes of the second are like the bases, and the second s

the loss

apple

Match

a portion of the string that is successfully described by the regex

Engine

A program or a library that obtains matches given a regex and a string

PCRE

regex flavors

23

Regular expressions are like ice cream

- 🔒 Common base
- Many flavors



regex flavors

- Three types of engines that affect how matching is done:
 - BFA
 - Traditional NFA
 - POSIX NFA
- For our purposes we discuss the regex flavor that PHPs' Perl-compatible regular expressions use

how an NFA engine works

The engine bumps along the string trying to match the regex

25

Sometimes it goes back and tries again



how an NFA engine works

Two basic things to understand about the engine

It will always return the earliest (leftmost) match it finds

The topic of the day is isotopes.

Given a choice it always favors match over a nonmatch



color legend







match

Syntax



building blocks

- Regexes are like LEGOs
- Small pieces combined into larger ones using connectors
- Arbitrarily complex

-

characters

30

-0

8

characters

31

- Special set is a well-defined subset of ASCII
- Ordinary set consist of all characters not designated special
- Special characters are also called <u>metacharacters</u>

matching literals

The most basic regex consists of a single ordinary character

matching literals

The most basic regex consists of a single ordinary character

It matches the first occurrence of that character in the string

matching literals

- The most basic regex consists of a single ordinary character
- It matches the first occurrence of that character in the string
- Characters can be added together to form longer regexes

extended characters

- To match an extended character, use \xhh notation where hh are hexadecimal digits
- To match Unicode characters (in UTF-8 mode) mode use \x{hhh..} notation

extended characters

For example, the following regex matches my name in Cyrillic:

Андрей
metacharacters

To use one of these literally, escape it, that is prepend it with a backslash





metacharacters

To escape a sequence of characters, put them between \Q and \E

Price is \Q\$12.36\E

will match

Price is \$12.36



metacharacters

So will the backslashed version

Price is \\$12.36

will match

Price is \$12.36



- Consist of a set of characters placed inside square brackets
- Matches one and only one of the characters specified inside the class



matches an English vowel (lowercase)

[aeiou]

matches an English vowel (lowercase)

[aeiou]

[st]urf

matches an English vowel (lowercase)

matches surf or turf

[aeiou]

[st]urf

negated classes

- Placing a caret as the first character after the opening bracket negates the class
- Will match any character not in the class, including newlines
- [^<>] would match a character that is not left or right bracket

character ranges

- Placing a dash (-) between two characters creates a range from the first one to the second one
- Useful for abbreviating a list of characters

[**a**-z]



character ranges



Ranges can be reversed





47

character ranges

- Ranges can be reversed
- A class can have more than one range and combine ranges with normal lists

[a-z0-9:]



character ranges

[0-9/]

48

[**z**-w]

matches a digit or a slash

matches z, y, x, or w

[a-z0-9]

matches digits and lowercase letters

X01-X1f matches co

matches control characters

shortcuts for ranges

Some ranges are so frequently used that it would be nice to have...





\w	word character	[A-Za-z0-9_]
\ d	decimal digit	[0–9]
\S	whitespace	[\n\r\t\f]
\W	not a word character	[^A-Za-z0-9_]
\ D	not a decimal digit	[^0-9]
\S	not whitespace	[^ \n\r\t\f]

shortcuts for ranges



Inside a character class, most metacharacters lose their meaning

Inside a character class, most metacharacters lose their meaning

Exceptions are:



Exceptions are:

closing bracket



Inside a character class, most metacharacters lose their meaning

Exceptions are:

closing bracket

🔒 backslash



Inside a character class, most metacharacters lose their meaning

Exceptions are:

closing bracket

🔒 backslash

🔒 caret



Inside a character class, most metacharacters lose their meaning

Exceptions are:

closing bracket

🔒 backslash

🔒 caret

🔒 dash

[ab\]]

[ab^]

 $\begin{bmatrix} a - z - \end{bmatrix}$

To use them literally, either escape them with a backslash or put them where they do not have special meaning

 \oplus

dot metacharacter

By default matches any single character

- By default matches any single character
- Except a newline!

- By default matches any single character
- Except a newline!



Is equivalent to



- Use dot carefully it might match something you did not intend
- 12.45 will match literal 12.45
- But it will also match these:

12345 12945 12a45 12-45

78812 45839

quantifiers

Or, Hit Me Baby One More Time

i i i i i i i i i i i i i i i i i i i	
i i i i i i i i i i i i i i i i i i i	
1	
1	
i i i i i i i i i i i i i i i i i i i	
1	
1	
i i i i i i i i i i i i i i i i i i i	i i
1	
1	
1	
i i i i i i i i i i i i i i i i i i i	
	1
1	
1	
 i	

quantifiers

64



Confucius said,

"Real knowledge is to know the extent of one's ignorance."



quantifiers



We are almost never sure about the contents of the text.



65

×

{ }

-\$

quantifiers



Quantifiers help us deal with this uncertainty

×

{ }

quantifiers

67



They specify how many times a regex component must repeat in order for the match to be successful

repeatable components



\w \d \s \W \D \S

range shortcuts

subpattern

backreference



zero-or-one

Indicates that the preceding component is optional

- Regex welcome!? will match either welcome or welcome!
- Regex Super\s?strong means that super and strong may have an optional whitespace character between them
- Regex hello[!?]? Will match hello, hello!, or hello?

Indicates that the preceding component has to appear once or more

- Regex a+h will match ah, aah, aaah, etc
- Regex –\d+ will match negative integers, such as -33
- Regex [^"]+ means to match a sequence (more than one) of characters until the next quote

one-or-more



zero-or-more

- Indicates that the preceding component can match zero or more times
- Regex \d+\.\d* will match 2., 3.1, 0.001
- Regex <[a-z][a-z0-9]*> will match an opening HTML tag with no attributes, such as or <h2>, but not <> or </i>

general repetition

72

- Specifies the minimum and the maximum number of times a component has to match
- Regex ha{1,3} matches ha, haa, haaa
- Regex \d{8} matches exactly 8 digits
general repetition

- If second number is omitted, no upper range is set
- Regex go{2,}al matches goal, gooal, gooal, etc

general repetition



75



"One of the weaknesses of our age is our apparent inability to distinguish our needs from our greeds." — Don Robinson

greediness n., matching as much as possible, up to a limit

-

greediness

PHP 5?

PHP 5 is still buggy

\d{2,4}



-

- Quantifiers try to grab as much as possible by default
- Applying <.+> to <i>greediness</i>matches the whole string rather than just <i>



If the entire match fails because they grabbed too much, then they are forced to give up as much as needed to make the rest of regex succeed



- To find words ending in ness, you will probably use \w+ness
- On the first run \w+ takes the whole word
- But since ness still has to match, it gives up the last 4 characters and the match succeeds



- The simplest solution is to make the repetition operators non-greedy, or lazy
- Lazy quantifiers grab as little as possible
- If the overall match fails, they grab a little more and the match is tried again



- To make a greedy quantifier lazy, append ?
- Note that this use of the question mark is different from its use as a regular quantifier



***?**

+?

{,}?

??





- Another option is to use negated character classes
- More efficient and clearer than lazy repetition



.+?> can be turned into <[^>]+>

- Note that the second version will match tags spanning multiple lines
- Single-line version: <[^>\r\n]+>

assertions and anchors

- An assertion is a regex operator that
 - expresses a statement about the current matching point
 - consumes no characters

assertions and anchors

- The most common type of an assertion is an anchor
- Anchor matches a certain position in the subject string

Λ

caret

- Caret, or circumflex, is an anchor that matches at the beginning of the subject string
- **^F** basically means that the subject string has to start with an F

^F

Fandango

\$

dollar sign

- Dollar sign is an anchor that matches at the end of the subject string or right before the string-ending newline
- \d\$ means that the subject string has to end with a digit
- The string may be top 10 or top 10\n, but either one will match

d\$



top 10

multiline matching

- Often subject strings consist of multiple lines
- If the multiline option is set:
 - Caret (^) also matches immediately after any newlines
 - Dollar sign (\$) also matches immediately before any newlines

Cone two three

absolute start/end

- Sometimes you really want to match the absolute start or end of the subject string when in the multiline mode
- These assertions are always valid:
 - A matches only at the very beginning
 - Z matches only at the very end
 - Z matches like \$ used in single-line mode



b

word boundaries

\bto\b



right |to| vote

- A word boundary is a position in the string with a word character (\w) on one side and a non-word character (or string boundary) on the other
- b matches when the current position is a word boundary
- B matches when the current position is not a word boundary

word boundaries

\bto\b



come together

- A word boundary is a position in the string with a word character (\w) on one side and a non-word character (or string boundary) on the other
- b matches when the current position is a word boundary
- B matches when the current position is not a word boundary

word boundaries

B2B



doc2html

- A word boundary is a position in the string with a word character (\w) on one side and a non-word character (or string boundary) on the other
- b matches when the current position is a word boundary
- B matches when the current position is not a word boundary

subpatterns

- Parentheses can be used group a part of the regex together, creating a subpattern
- You can apply regex operators to a subpattern as a whole



grouping

95

- Regex is(land)? matches both is and island
- Regex (\d\d,)*\d\d will match a comma-separated list of double-digit numbers



capturing subpatterns

- All subpatterns by default are capturing
- A capturing subpattern stores the corresponding matched portion of the subject string in memory for later use

capturing subpatterns

- Subpatterns are numbered by counting their opening parentheses from left to right
- Regex (\d\d-(\w+)-\d{4}) has two subpatterns

(\d\d-(\w+)-\d{4})



12-May-2004

capturing subpatterns

- Subpatterns are numbered by counting their opening parentheses from left to right
- Regex (\d\d-(\w+)-\d{4}) has two subpatterns
- When run against 12-May-2004 the second subpattern will capture May

 $(d^{\psi+})-d^{4})$

12-May-2004

non-capturing subpatterns

- The capturing aspect of subpatterns is not always necessary
- It requires more memory and more processing time

non-capturing subpatterns

- Using ?: after the opening parenthesis makes a purely grouping subpattern
- Regex box(?:ers)? will match boxers but will not capture anything
- The (?:) subpatterns are not included in the subpattern numbering

named subpatterns

- It can be hard to keep track of subpattern numbers in a complicated regex
- Using ?P<name> after the opening parenthesis creates a named subpattern
- Named subpatterns are still assigned numbers
- Pattern (?P<number>\d+) will match and capture 99 into subpattern named number when run against 99 bottles

Alternation operator allows testing several sub-expressions at a given point

- The branches are tried in order, from left to right, until one succeeds
- Empty alternatives are permitted
- Regex sailing cruising will match either sailing or cruising

alternation

- Since alternation has the lowest precedence, grouping is often necessary
- sixth|seventh sense will match the word sixth or the phrase seventh sense
- (sixth|seventh) sense will match sixth sense or seventh sense

alternation

- Remember that the regex engine is eager
- It will return a match as soon as it finds one
- Camel came camera will only match came when run against camera
- Put more likely regex as the first alternative

alternation



Also known as "if at first you don't succeed, try, try again"

When faced with several options it could try to achieve a match, the engine picks one and remembers the others





If the picked option does not lead to an overall successful match, the engine backtracks to the decision point and tries another option



- This continues until an overall match succeeds or all the options are exhausted
 - The decision points include quantifiers and alternation



Two important rules to remember

- With greedy quantifiers the engine always attempts the match, and with lazy ones it delays the match
- If there were several decision points, the engine always goes back to the most recent one


backtracking example

d+00

12300

start

backtracking example





backtracking example



11



backtracking example





backtracking example





backtracking example







d+00



string exhausted still need to match 00

-0

backtracking example

d+00



give up 0

 \oplus

backtracking example

d+00

12300

give up 0

backtracking example

d+00



-

backtracking example

d+00



success

backtracking example

\d+ff

123dd

start

 $\uparrow \uparrow \uparrow \uparrow \uparrow \uparrow \uparrow \uparrow \uparrow$

backtracking example





backtracking example





backtracking example





\d+ff



cannot match f here

\d+ff



give up 3 still cannot match f

\d+ff



give up 2 still cannot match f

\d+ff



cannot give up more because of +

backtracking example

\d+ff

123dd

failure

-

backtracking example

ab??c





backtracking example

<u>ab??c</u>



add a

 $\uparrow \uparrow \uparrow \uparrow \uparrow \uparrow \uparrow \uparrow \uparrow$

backtracking example





skip matching **b** at first

 $\uparrow \uparrow \uparrow \uparrow \uparrow \uparrow \uparrow \uparrow \uparrow$

backtracking example

ab??<u>c</u>



c cannot match here





go back and try matching b now

ab??<u>c</u>



c can be matched

ab??c



SUCCESS

atomic grouping

- Disabling backtracking can be useful
- The main goal is to speed up failed matches, especially with nested quantifiers



atomic grouping

- (?>regex) will treat regex as a single atomic token, no backtracking will occur inside it
- All the saved states are forgotten



atomic grouping

- (?>\d+)ff will lock up all available digits and fail right away if the next two characters are not ff
- Atomic groups are not capturing



possessive quantifiers



- Atomic groups can be arbitrarily complex and nested
- Possessive quantifiers are simpler and apply to a single repeated item

possessive quantifiers



- To make a quantifier possessive append a single +
- \d++ff is equivalent to (?>\d+)ff

possessive quantifiers



Other ones are *+, ?+, and {m,n}+

Possessive quantifiers are always greedy

do not over-optimize

- Keep in mind that atomic grouping and possessive quantifiers can change the outcome of the match
- When run against string abcdef
 - \w+d will match abcd
 - \w++d will not match at all
 - \w+ will match the whole string

h

backreferences

- A backreference is an alias to a capturing subpattern
- It matches whatever the referent capturing subpattern has matched



backreferences

- (re|le)\w+\1 matches words that start with re or le and end with the same thing
- For example, <u>retire</u> and <u>legible</u>, but not <u>revocable</u> or <u>lecture</u>
- Reference to a named subpattern can be made with (?P=name)


lookaround

- Assertions that test whether the characters before or after the current point match the given regex
- Consume no characters
- Do not capture anything
- Includes lookahead and lookbehind

positive lookahead



- Tests whether the characters after the current point match the given regex
- (\w+)(?=:)(.*) matches surfing: a sport but semicolon ends up in the second subpattern

negative lookahead

- Tests whether the characters after the current point do not match the given regex
- fish(?!ing) matches fish not followed by ing
- Will match fisherman and fished



(?!)

negative lookahead

Difficult to do with character classes

- fish[^i][^n][^g] might work but will consume more than needed and fail on subjects shorter than 7 letters
- Don't even go there if you want something like fish(?!hook|ing)



positive lookbehind

- Tests whether the characters immediately preceding the current point match the given regex
- The regex must be of fixed size, but branches are allowed
- (?<=foo)bar matches bar only if preceded by foo, e.g. my foobar

(?<!)

negative lookbehind

- Tests whether the characters immediately preceding the current point do not match the given regex
- Once again, regex must be of fixed size
- (?<!foo)bar matches bar only if not preceded by foo, e.g. in the bar but not my foobar



inline options

The matching can be modified by options you put in the regular expression

(?i)

(?m)

(?s)

(?x)

(?U)

enables case-insensitive mode

enables multiline matching for ^ and \$

makes dot metacharacter match newline also

ignores literal whitespace

makes quantifiers ungreedy (lazy) by default

inline options

(?i)

(?m)

(?s)

(?x)

(?U)

Options can be combined and unset (?im-sx)

At top level, apply to the whole pattern

Localized inside subpatterns

(a(?i)b)c

comments

Here's a regex I wrote when working on Smarty templating engine

153



comments

Let me blow that up for you

^\\$\w+(?>(\[(\d+|\\$\w+|\w+(\.\w+)?)\])| ((\.|->)\\$?\w+))*(?>\|@?\w+(:(?>"[^"\\\]* (?:\\\\.[^"\\\]*)*"|\'[^\'\\\]* (?\\\.[^\'\\\]*)*\'|[^|]+))*)*\$

Would you like some comments with that?

2#

comments

- Most regexes could definitely use some comments
- (?#...) specifies a comment

\d+(?# match some digits)



?#

comments

- If (?x) option is set, anything after # outside a character class and up to the next newline is considered a comment
- To match literal whitespace, escape it

(?x) \w+ # start with word characters [?!] # and end with ? or !



Regex API



Regex API

- Perl-compatible regex API (PCRE) was introduced in PHP 3.0.9
- **Starting with PHP 4.2.0 the API is enabled by default**
- Uses consistent pattern syntax
- All functions start with preg_ prefix

'/[abc]+/'

The regex must be enclosed in delimiters and passed as a single- or double-quoted string

"/[abc]+/"

<u>z[abc]+z</u>

NO!

The regex must be enclosed in delimiters and passed as a single- or double-quoted string

Delimiter character cannot be alphanumeric or backslash

/<\/i>/

- The regex must be enclosed in delimiters and passed as a single- or double-quoted string
- Delimiter character cannot be alphanumeric or backslash
- If the delimiter character has to be used in the regex, escape it with a backslash

/<a.+?>/is

- The regex must be enclosed in delimiters and passed as a single- or double-quoted string
- Delimiter character cannot be alphanumeric or backslash
- If the delimiter character has to be used in the regex, escape it with a backslash
- The ending delimiter may optionally be followed by pattern modifiers

pattern modifiers

The first five should be familiar

enables case-insensitive mode

enables multiline matching for ^ and \$

makes dot metacharacter match newline also

ignores literal whitespace and allows **#** comments

makes quantifiers ungreedy (lazy) by default



pattern modifiers

But there are some more

/A

/S

/u

/e

anchors the pattern at the beginning of string (similar to \A assertion)

performs additional analysis on the pattern

enables UTF-8 mode

explained in preg_replace() section

pattern examples

Valid:

- $\langle d{4}-d/d/d/d/$
- .*?<\/\1>/iU
- Interpretation of the second secon



pattern examples

Invalid:

166

- I/.+\$ missing end delimiter
- Ab(c|d)/J unknown modifier J
- \$\s?*/ compilation failure, misapplied quantifier *



- PHP can interpret regex metacharacters as its own
- **To avoid confusion:**
 - Backslash the common metacharacters
 - Use single quotes to make life easier

- Even with single quotes, the "leaning toothpick" syndrome may occur
- To match a single backslash, one has to use '/\\\/'



- Even with single quotes, the "leaning toothpick" syndrome may occur
- To match a single backslash, one has to use '/////'
- **First, PHP interprets it as '///**

'///'

- Even with single quotes, the "leaning toothpick" syndrome may occur
- To match a single backslash, one has to use '/\\\/'
 - ♣ First, PHP interprets it as '/\/'
 - Then, regex engine sees it as an escaped backslash metacharacter

locales

- Caseless matching and character class determination are affected by the current locale
- The locale can be changed via PHP's setlocale() function
- For example, set_locale('fr_FR') will set the French locale which will be taken into account by the engine

to save time...



Since all PCRE functions are described in the manual in exquisite detail, we'll just have a brief look at them...

173

preg_match(string regex, string subject, array matches, int flags, int offset)

- Tries to find the first occurrence of a pattern described by regex in the string
- Returns 0 or 1 (FALSE on error)
- If matches is provided, it is filled with the match results
- Stops after the first successful match
- Best used for validation

preg_match_all(string regex, string subject, array matches, *int flags, int offset*)

- **Tries to find all patterns described by regex in the string**
- Matching continues from the end of the last match
- Return number of successful matches or FALSE on error

175

preg_replace(mixed regex, mixed replacement, mixed subject, *int limit*)

- Applies regex to subject and replaces matches with replacement
- limit specifies how many matches to replace, -1 means no limit (the default)
- Returns modified subject if matches are found
- regex, subject, and replacement can be one-dimensional arrays
- Allows for multiple searches and replacements on multiple strings at once

preg_replace(mixed regex, mixed replacement, mixed subject, *int limit*)

- replacement may contain references of the form \\n or \$n (the preferred syntax)
- Such reference will be replaced by the text matched by the n'th capturing subpattern

preg_replace(mixed regex, mixed replacement, mixed subject, *int limit*)

- /e modifier on regex treats replacement as PHP code
- The references are resolved, the code is evaluated, and the result is used as the replacement

preg_replace_callback(mixed regex, mixed callback, mixed subject, *int limit*)

- Identical to preg_replace() except that the replacement is specified by a callback function
- For each match the callback is invoked with the match info and is supposed to return the replacement string

preg_split(string regex, string subject, int limit, int flags)

- Splits subject along boundaries matched by regex
- Returns an array of split pieces
- Iimit determines the maximum number of pieces, -1 means no limit (the default)
- The type of splitting can be controlled by flags

preg_grep(string regex, array input, int flags)

- Applies regex to each element of input array
- Return a new array consisting only of elements that matched
- If flags if PREG_GREP_INVERT, only the elements that did not match will be returned
Regex Toolkit



regex toolkit

- In your day-to-day development, you will frequently find yourself running into situations calling for regular expressions
- It is useful to have a toolkit from which you can quickly draw the solution
- It is also important to know how to avoid problems in the regexes themselves

matching vs. validation

- In matching (extraction) the regex must account for boundary conditions
- In validation your boundary conditions are known – the whole string



matching vs. validation

Matching an English word starting with a capital letter

\b[A-Z][a-zA-Z'-]*\b

Validating that a string fulfills the same condition

^[A-Z][a-zA-Z'-]*\$

Do not forget ^ and \$ anchors for validation!



using dot properly

- One of the most used
- One of the most <u>misused</u>
- Remember dot is a shortcut for [^\n]
- May match more than you really want
- Solution
- Be explicit about what you want
- <a><[a-z]> is better

using dot properly



- When dot is combined with quantifiers it becomes greedy
- <.+> will consume any characters between the first bracket in the line and the last one
- Including any other brackets!

using dot properly



It's better to use negated character class instead

<[^>]+> if bracketed expression spans lines

<[^>\r\n]+> otherwise

Lazy quantifier can be used, but they are not as efficient, due to backtracking

One of the most common problems is combining an inner repetition with an outer one (regex1|regex2|..)*

(regex*)+

(regex+)*

- One of the most common problems is combining an inner repetition with an outer one
- If the initial match fails, the number of ways to split the string between the quantifiers grows exponentially

(regex1|regex2|..)*

(regex*)+

(regex+)*

- One of the most common problems is combining an inner repetition with an outer one
- If the initial match fails, the number of ways to split the string between the quantifiers grows exponentially
- The problem gets worse when the inner regex contains a dot, because it can match anything!

(regex1|regex2|..)*

(regex*)+

(regex+)*

- PCRE has an optimization that helps in certain cases, and also a hardcoded limit for the backtracking
- The best way to solve this is to prevent unnecessary backtracking in the first place via atomic grouping or possessive quantifiers

(regex1|regex2|..)*

(regex*)+

(regex+)*

- Consider the expression
 - ["'](\w+|\s{1,2})*["']
- When applied to the string "aaaaaaaaaaa" (with final quote), it matches quickly

We can prevent backtracking from going back to the matched portion by adding a possessive quantifier:

["'](\w+|\s{1,2})*+["']

With nested unlimited repeats, you should lock up as much of the string as possible right away

- Naïve implementation:
 - Match ([a-z]+) \1
 - Replace with \$1
 - Has problems with This is island



- **Naïve implementation:**
 - Match ([a-z]+) \1
 - Replace with \$1
 - Has problems with This is island



Better approach:

- Match \b([a-z]+) \1\b
- Replace with \$1
- Handles This is island just fine





Even better, concentrate on delimiters



Even better, concentrate on delimiters

First match a non-delimiter sequence



- Even better, concentrate on delimiters
- First match a non-delimiter sequence, that is preceded by a delimiter or beginning of string



- Even better, concentrate on delimiters
- First match a non-delimiter sequence, that is preceded by a delimiter or beginning of string
- Then match atomically one or more duplicates of the first match, separated by delimiters



- Even better, concentrate on delimiters
- First match a non-delimiter sequence, that is preceded by a delimiter or beginning of string
- Then match atomically one or more duplicates of the first match, separated by delimiters
- And make sure it is followed by a delimiter or the end of the string



- Even better, concentrate on delimiters
- First match a non-delimiter sequence, that is preceded by a delimiter or beginning of string
- Then match atomically one or more duplicates of the first match, separated by delimiters
- And make sure it is followed by a delimiter or the end of the string
- Replace with \$1

This is unwieldy. Let's use PHP interpolation.

```
$dlm = '[\s.,?!]';
$n_dlm = '[^\s.,?!]';
$s = preg_replace("/
        (?<=$dlm|^)
        ($n_dlm+)
        ($dlm++\\1)+
        (?=$dlm|$)
        /x", '$1', $s);
```



removing multiline comments

- Simple, if the comments are not allowed to nest
- I/*.*?*/!s replaced with an empty string will work for C-like comments
- General pattern: /start.*?end/s
- For nested comments, a recursive pattern is necessary

extracting markup

- Possible to use preg_match_all() for grabbing marked up portions
- But for tokenizing approach, preg_split() is better

```
$s = 'a <b><I>test</I></b> of <br /> markup';
$tokens = preg_split(
    '!( < /? [a-zA-Z][a-zA-Z0-9]* [^/>]* /? > ) !x', $s, -1,
    PREG_SPLIT_NO_EMPTY | PREG_SPLIT_DELIM_CAPTURE);
result is array('a','<b>','<I>','test','</I>',
```

```
'</b>','of','<br />','markup')
```

restricting markup

- Suppose you want to strip all markup except for some allowed subset, what are your possible approaches?
 - Use strip_tags() which is limited
 - Multiple invocations of str_replace() or preg_replace() to remove script blocks, etc
 - Custom tokenizer and processor, or..

}

 \oplus

restricting markup

\$s = preg_replace_callback(
 '! < (/?) ([a-zA-Z][a-zA-Z0-9]*) ([^/>]*) (/?) > !x',
 'my_strip', \$s);

```
function my_strip($match) {
   static $allowed_tags = array('b', 'i', 'p', 'br', 'a');
   $tag = $match[2];
   $attrs = $match[3];
   if (!in_array($tag, $allowed_tags)) return '';
   if (!empty($match[1])) return "</$tag>";
   /* strip evil attributes here */
   if ($tag == 'a') { $attrs = ''; }
   /* any other kind of processing here */
   return "<$tag$attrs$match[4]>";
```

matching numbers

- Integers are easy: \b\d+\b
- Floating point numbers:
 - integer.fractional
 - .fractional
- Can be covered by (\b\d+)?\.\d+\b



matching numbers

- To match both integers and floating point numbers, either combine them with alternation or use: ((\b\d+)?\.)?\b\d+\b
- [+-]? can be prepended to any of these, if sign matching is needed
- b can be substituted by more appropriate assertions based on the required delimiters



matching quoted strings

- A simple case is a string that does not contain escaped quotes inside it
- Matching a quoted string that spans lines:

"[**^**"]*"

Matching a quoted string that does not span lines:

"[**^**"\r\n]*"

matching quoted strings

opening quote

a component that is

a segment without any quotes

or

a quote preceded by a backslash

component repeated zero or more times without backtracking closing quote

Matching a string with escaped quotes inside

"([^"]+|(?<=\\\\)")*+"

| (?<=\\\\)")*+

"

[^"]+

"

matching e-mail addresses

- Do I look crazy to you?
- The complete regex is about a book page long in 10-point type
- Buy a copy of Jeffrey Friedl's book and steal it from there

matching phone numbers

Assuming we want to match US/Canada-style phone numbers

- 800-555-1212 1-800-555-1212
- 800.555.1212 1.800.555.1212
- (800) 555-1212 1 (800) 555-1212

How would you do it?

matching phone numbers

- The simplistic approach could be:
 (1[.-])? \(? \d{3} \)? [.-] \d{3} [.-] \d{4}
 But this would result in a lot of false positives:
 - 1. (800) 555 1212 800) . 555 1212
 - 1-800
 555-1212
 (800
 555-1212

matching phone numbers

^(?:	anchor to the start of the string
(?:1([]))?	may have 1. or 1- (remember the separator)
\d{3}	three digits
((?(1)	if we had a separator
\1	match the same (and remember), otherwise
[]))	match . or - as a separator (and remember)
\d{3}	another three digits
\2	same separator as before
\d{4}	final four digits
	or
1[]?\(\d{3}\)[]\d{3}-\d{4}	just match the third format
)\$	anchor to the end of the string

tips

- Don't do everything in regex a lot of tasks are best left to PHP
- Use string functions for simple tasks
- Make sure you know how backtracking works


tips



- Be aware of the context
- Capture only what you intend to use
- Don't use single-character classes



tips

Lazy vs. greedy, be specific

Put most likely alternative first in the alternation list

Think!



Thank You!

Questions?

