

Stop overusing regular expressions!

Franklin Chen

<http://franklinchen.com/>

Pittsburgh Tech Fest 2013

June 1, 2013

Introduction

Jamie Zawinski famously wrote:

*Some people, when confronted with a problem, think,
“I know, I’ll use regular expressions.”
Now they have two problems.*

- ▶ Why the controversy over using regular expressions (regexes)?

Introduction

Jamie Zawinski famously wrote:

*Some people, when confronted with a problem, think,
“I know, I’ll use regular expressions.”
Now they have two problems.*

- ▶ Why the controversy over using regular expressions (regexes)?
- ▶ What are better regex practices?

Introduction

Jamie Zawinski famously wrote:

*Some people, when confronted with a problem, think,
“I know, I’ll use regular expressions.”
Now they have two problems.*

- ▶ Why the controversy over using regular expressions (regexes)?
- ▶ What are better regex practices?
- ▶ What alternatives?

An infamous regex for email

A big Perl regex for email address based on RFC 822 grammar:

/(?: (?: \r\n)? [\t]) * (?: (?: (?: [^()<>@,;: \\". \[\] \000-\031] -) + | \Z | (?: [\"()<>@,;: \\". \[\]])) | " (?: [^\" \r\n \\\\ | \\. | (?: (?: \r\n) \r\n)? [\t]) *) (?: \\. (?: (?: \r\n)? [\t]) * (?: [^()<>@,;: \\". \[\] \000-\031] + (?: (?: (?: \r\n)? [\t]) + | \Z | (?: [\"()<>@,;: \\". \[\]])) | \. (?: (?: \r\n)? [\t]) * (?: [^()<>@,;: \\". \[\] \000-\031] + (?: (?: (?: \r\n)? [\t]) + | \Z | (?: [\"()<>@,;: \\". \[\]])) | \. (?: (?: \r\n)? [\t]) *) * | (?: [^()<>@,;: \\". \[\] \000-\031] + (?: (?: (?: \r\n)? [\t]) + | \Z | (?: [\"()<>@,;: \\". \[\]])) | " (?: [^\" \r\n \\\\ | \\. | (?: (?: \r\n) \r\n)? [\t]) *) * \< (?: (?: \r\n)? [\t]) * (?: @ (?: [^()<>@,;: \\". \[\] \000-\031] + (?: (?: (?: \r\n)? [\t]) + | \Z | (?: [\"()<>@,;: \\". \[\]])) | \. (?: (?: \r\n)? [\t]) * (?: [^()<>@,;: \\". \[\] \000-\031] + (?: (?: (?: \r\n)? [\t]) + | \Z | (?: [\"()<>@,;: \\". \[\]])) | \. (?: (?: \r\n)? [\t]) *) *) * (?: , @ (?: [^()<>@,;: \\". \[\] \000-\031] + (?: (?: (?: \r\n)? [\t]) + | \Z | (?: [\"()<>@,;: \\". \[\]])) | \. (?: (?: \r\n)? [\t]) * (?: [^()<>@,;: \\". \[\] \000-\031] + (?: (?: (?: \r\n)? [\t]) + | \Z | (?: [\"()<>@,;: \\". \[\]])) | \. (?: (?: \r\n)? [\t]) *) *) (?: \\. (?: (?: \r\n)? [\t]) * (?: [^()<>@,;: \\". \[\] \000-\031] + (?: (?: (?: \r\n)? [\t]) + | \Z | (?: [\"()<>@,;: \\". \[\]])) | \. (?: (?: \r\n)? [\t]) *) *)

Outline

- ▶ Regular expressions
- ▶ Parsing more powerful grammars
- ▶ Practical considerations
- ▶ Concepts *not* language-specific

Example code in which languages?

Runnable, tested code projects: <https://github.com/franklinchen/talk-on-overusing-regular-expressions>

Time limitations: take two representatives.

- ▶ Dynamic typing: **Ruby**

Example code in which languages?

Runnable, tested code projects: <https://github.com/franklinchen/talk-on-overusing-regular-expressions>

Time limitations: take two representatives.

- ▶ Dynamic typing: **Ruby**
 - ▶ Perl, Python, JavaScript, Clojure, etc.

Example code in which languages?

Runnable, tested code projects: <https://github.com/franklinchen/talk-on-overusing-regular-expressions>

Time limitations: take two representatives.

- ▶ Dynamic typing: **Ruby**
 - ▶ Perl, Python, JavaScript, Clojure, etc.
- ▶ Static typing: **Scala**

Example code in which languages?

Runnable, tested code projects: <https://github.com/franklinchen/talk-on-overusing-regular-expressions>

Time limitations: take two representatives.

- ▶ Dynamic typing: **Ruby**
 - ▶ Perl, Python, JavaScript, Clojure, etc.
- ▶ Static typing: **Scala**
 - ▶ Java, C++, C#, ML (OCaml, F#), Haskell, etc.

My email address fiasco

To track and prevent spam: `FranklinChen+spam@cmu.edu`

- ▶ Some sites claimed invalid (wrongly)

My email address fiasco

To track and prevent spam: `FranklinChen+spam@cmu.edu`

- ▶ Some sites claimed invalid (wrongly)
- ▶ Some sites allowed registration

My email address fiasco

To track and prevent spam: `FranklinChen+spam@cmu.edu`

- ▶ Some sites claimed invalid (wrongly)
- ▶ Some sites allowed registration
 - ▶ I caught spam

My email address fiasco

To track and prevent spam: `FranklinChen+spam@cmu.edu`

- ▶ Some sites claimed invalid (wrongly)
- ▶ Some sites allowed registration
 - ▶ I caught spam
 - ▶ Unsubscribing failed

My email address fiasco

To track and prevent spam: `FranklinChen+spam@cmu.edu`

- ▶ Some sites claimed invalid (wrongly)
- ▶ Some sites allowed registration
 - ▶ I caught spam
 - ▶ Unsubscribing failed
 - ▶ Some claimed invalid (wrongly since I registered!)

My email address fiasco

To track and prevent spam: `FranklinChen+spam@cmu.edu`

- ▶ Some sites claimed invalid (wrongly)
- ▶ Some sites allowed registration
 - ▶ I caught spam
 - ▶ Unsubscribing failed
 - ▶ Some claimed invalid (wrongly since I registered!)
 - ▶ Some silently failed to unsubscribe

My email address fiasco

To track and prevent spam: `FranklinChen+spam@cmu.edu`

- ▶ Some sites claimed invalid (wrongly)
- ▶ Some sites allowed registration
 - ▶ I caught spam
 - ▶ Unsubscribing failed
 - ▶ Some claimed invalid (wrongly since I registered!)
 - ▶ Some silently failed to unsubscribe
 - ▶ Had to set up spam filter

Don't Do It Yourself: find libraries

Find a library of regexes!

Example: Perl has `Mail::RFC822::Address`, the infamous regex!

```
# Install the library
```

```
$ cpanm Mail::RFC822::Address
```

Other regexes ready for use:

<http://search.cpan.org/dist/Regexp-Common/Regexp::Common>

Build your own libraries

Excerpt from [source code](#) of Mail::RFC822::Address

```
# ...some excerpts only
my $atom = "[^$specials $controls]+(?:$lws+|\\Z|(?=[\\[\\\"'@]))";
my $word = "(?:$atom|$quoted_string)";
my $localpart = "$word(?:\\.$lws*$word)*";
my $sub_domain = "(?:$atom|$domain_literal)";
my $domain = "$sub_domain(?:\\.$lws*$sub_domain)*";
my $addr_spec = "$localpart\\@$lws*$domain";
my $phrase = "$word*";
my $route = "(?:\\@$domain(?:,\\@$lws*$domain)*:$lws*)";
my $route_addr = "\\<$lws*$route?$addr_spec\\>$lws*";
my $mailbox = "(?:$addr_spec|$phrase$route_addr)";
my $group = "$phrase:$lws*(?:$mailbox(?:,\\s*$mailbox)*)?";
my $address = "(?:$mailbox|$group)";
```

Regexes do not report errors usefully

- ▶ If success: we can extract matching/grouping information
- ▶ If failure: get no information about why!

We have a dilemma:

- ▶ Bug in the regex?
- ▶ The data really doesn't match?

Better error reporting: how?

Would prefer something at least minimally like:

```
[1.13] failure: '@' expected but '+' found
```

```
FranklinChen+spam@cmu.edu  
          ^
```

Features of *parser* libraries:

- ▶ Extraction of line, column information of error
- ▶ Automatically generated error explanation
- ▶ Hooks for customizing error explanation
- ▶ Hooks for error recovery

Modularized regex: Scala

```
val user = "[a-zA-Z]+"
val at = "@+"
val domainSegment = "[^@\\.]+"
val dot = "\\."
val email = raw"""(?x)
  \A
  $user
  $at
  (
    $domainSegment $dot
  )+
  $domainSegment
  \z"""

val emailRegex = email.r
```

Modularized regex: Ruby

```
LOCAL_PART = /[a-zA-Z]+/x
```

```
AT = /@/x
```

```
DOMAIN_CHAR = /^[^@\.]/x
```

```
SUB_DOMAIN = /#{DOMAIN_CHAR}+/x
```

```
DOT = /\./x
```

```
DOMAIN = /(#{SUB_DOMAIN}#{DOT})+ #{SUB_DOMAIN}/x
```

```
EMAIL = /
```

```
  \A
```

```
  #{LOCAL_PART}
```

```
  #{AT}
```

```
  #{DOMAIN}
```

```
  \z
```

```
/x
```

Email parser: Scala

Scala comes with **standard parser combinator library**.

```
object EmailParsers extends RegexParsers {  
  override def skipWhitespace = false  
  
  def localPart: Parser[String] = "[a-zA-Z]+"  
  def at = "@"  
  def domainChar = "[^@\\.]"  
  def subDomain = rep1(domainChar)  
  def dot = "."  
  def domain = rep1(subDomain ~ dot) ~ subDomain  
  def email = localPart ~ at ~ domain  
}
```

Inheriting from `RegexParsers` allows the implicit conversions from regexes into parsers.

Email parser: Ruby

Ruby does not have a standard parser combinator library. One that is popular is **Parslet**.

```
require 'parslet'

# A simplified email address parser
class EmailValidator::Parser < Parslet::Parser
  rule(:local_part) { match['a-zA-Z'].repeat(1) }
  rule(:at)         { str('@') }
  rule(:domain_char) { match['~@\\\\.'] }
  rule(:sub_domain) { domain_char.repeat(1) }
  rule(:dot)        { str('.') }
  rule(:domain)     { (sub_domain >> dot).repeat(1) >>
                      sub_domain }
  rule(:email)      { local_part >> at >> domain }
end
```

Error reporting: Scala

We have achieved the goal of decent error reporting:

```
parseAll(email, address) match {  
  case Success(_, _) =>  
    println(s"Successfully parsed $address")  
  case failure: NoSuccess =>  
    println("Parse failed, and here's why:")  
    println(failure)  
}
```

Parse failed, and here's why:

[1.13] failure: '@' expected but '+' found

FranklinChen+spam@cmu.edu

^

Error reporting: Ruby

We have achieved the goal of decent error reporting:

```
begin
  parser.email.parse(address)
  puts "Successfully parsed #{address}"
rescue Parslet::ParseFailed => error
  puts "Parse failed, and here's why:"
  puts error.cause.ascii_tree
end
```

\$ bundle exec bin/validate_email 'FranklinChen+spam@cmu.edu'

Parse failed, and here's why:

Failed to match sequence (LOCAL_PART AT DOMAIN) at line 1 of <stdin>

– Expected "@", but got "+" at line 1 char 13.

Infamous email regex: revisited

The Perl regex in `Mail::RFC822::Address` was actually *manually* back-converted from a parser generated by the Perl library

`Parse::RecDescent`.

Reason: speed.

There are always tradeoffs.

Combinator library vs. generator

We have seen the elegance of parser combinator libraries. What about the other approach?

- ▶ Combinator library: internal DSL
 - ▶ Just code: easiest for getting started
 - ▶ Easier to debug
 - ▶ May be slow: interpreted

Combinator library vs. generator

We have seen the elegance of parser combinator libraries. What about the other approach?

- ▶ Combinator library: internal DSL
 - ▶ Just code: easiest for getting started
 - ▶ Easier to debug
 - ▶ May be slow: interpreted
- ▶ Parser generator: external DSL
 - ▶ Requires running a separate program to generate source (compile)
 - ▶ yacc: `grammar.y` to `grammar.c`
 - ▶ **Racc**: `grammar.y` to `grammar.rb`
 - ▶ **ANTLR**: `grammar.g` to `grammar.java` (I used for a project years ago)

Why validate email address anyway?

Look at the bigger picture.

This guy advocates **simply sending a user an activation email**.

Engineering tradeoff: the email sending and receiving programs need to handle the email address anyway.

Email example wrapup

- ▶ It is possible to write a regex
- ▶ But you may want to use someone else's regex
- ▶ If you write a regex, modularize it
- ▶ For error reporting, use a parser: convert from modularized regex
- ▶ Do you even need to solve the problem?

Example: toy JSON parsing

```
{  
  "distance" : 5.6,  
  "address" : {  
    "street" : "0 Nowhere Road",  
    "neighbors" : ["X", "Y"],  
    "garage" : null  
  }  
}
```

- ▶ Would you use a regex to parse?
- ▶ Could you use a regex to parse?

Toy JSON parser: Scala

(To save time: no more Ruby code in this presentation.)

```
object ToyJSONParsers extends JavaTokenParsers {  
  def value: Parser[Any] = obj |  
    arr |  
    stringLiteral |  
    floatingPointNumber |  
    "null" |  
    "true" | "false"  
  def obj = "{" ~ repsep(member, ",") ~ "}"  
  def arr = "[" ~ repsep(value, ",") ~ "]"  
  def member = stringLiteral ~ ":" ~ value  
}
```

Inheriting from JavaTokenParsers allows the reuse of stringLiteral and floatingPointNumber parsers.

Fancier toy JSON parser: use the data

Want to actually shape and use the data parsed.

```
{  
  "distance" : 5.6,  
  "address" : {  
    "street" : "0 Nowhere Road",  
    "neighbors" : ["X", "Y"],  
    "garage" : null  
  }  
}
```

Example: traverse the JSON to address and then to the second neighbor, Y.

Domain modeling: Scala

```
sealed trait ToyJSON

case class JObject(map: Map[String, ToyJSON]) extends ToyJSON
case class JArray(list: List[ToyJSON]) extends ToyJSON
case class JString(string: String) extends ToyJSON
case class JFloat(float: Float) extends ToyJSON
case object JNull extends ToyJSON
case class JBoolean(boolean: Boolean) extends ToyJSON
```

Fancier toy JSON parser: Scala

```
def value: Parser[ToyJSON] = obj | arr |  
  stringLiteralStripped ^^ { JString(_) } |  
  floatingPointNumber ^^ { s => JFloat(s.toFloat) } |  
  "null" ^^^ JNull | "true" ^^^ JBoolean(true) |  
  "false" ^^^ JBoolean(false)  
  
def stringLiteralStripped: Parser[String] =  
  stringLiteral ^^ { s => s.substring(1, s.length-1) }  
  
def obj: Parser[JObject] = "{" ~> (repsep(member, ",") ^^  
  { aList => JObject(aList.toMap) }) <~ "}"  
  
def arr: Parser[JArray] = "[" ~> (repsep(value, ",") ^^  
  { aList => JArray(aList) }) <~ "]"  
  
def member: Parser[(String, ToyJSON)] =  
  ((stringLiteralStripped <~ ":") ~ value) ^^  
  { case s ~ v => s -> v }
```

Using the parsed JSON

A test using **Specs2** Scala testing framework:

```
parsers.value must succeedOn("""{
  "distance" : 5.6,
  "address" : {
    "street" : "0 Nowhere Road",
    "neighbors" : ["X", "Y"],
    "garage" : null
  }
}""").withResult({ result: ToyJSON =>
  result.asInstanceOf[JObject].
    map("address").asInstanceOf[JObject].
    map("neighbors").asInstanceOf[JArray].
    list(1).asInstanceOf[JString].
    string
}) ^^ be_==("Y"))
```

JSON: reminder not to reinvent

- ▶ Use a standard JSON parsing library!

JSON: reminder not to reinvent

- ▶ Use a standard JSON parsing library!
 - ▶ **Scala has one.**

JSON: reminder not to reinvent

- ▶ Use a standard JSON parsing library!
 - ▶ **Scala has one.**
 - ▶ All languages have a standard JSON parsing library.

JSON: reminder not to reinvent

- ▶ Use a standard JSON parsing library!
 - ▶ **Scala has one.**
 - ▶ All languages have a standard JSON parsing library.
 - ▶ Shop around: alternate libraries have different tradeoffs.

JSON: reminder not to reinvent

- ▶ Use a standard JSON parsing library!
 - ▶ **Scala has one.**
 - ▶ All languages have a standard JSON parsing library.
 - ▶ Shop around: alternate libraries have different tradeoffs.
- ▶ Other standard formats: HTML, XML, CSV, etc.

JSON wrapup

- ▶ Just parsing is simple

JSON wrapup

- ▶ Just parsing is simple
- ▶ Domain modeling is trickier

JSON wrapup

- ▶ Just parsing is simple
- ▶ Domain modeling is trickier
- ▶ *Use an existing JSON parsing library already!*

Final example of real regex code

Recently in [Octopress Ruby code](#):

```
EXPRESSION = /(.*?)\s+(unless|if)\s+(.+)/i
```

```
TERNARY = /(.*?)\s*(.*?)\s+\?\s+(.*?)\s+:\s+(.*?)\s*(.*)\s*(.?)
```

```
def strip_expression(markup, context = false)
  if markup =~ TERNARY
    result = evaluate_ternary($2, $3, $4, context)
    markup = "#{$1} #{result} #{ $5}"
  end
  markup =~ EXPRESSION ? $1 : markup
end
```

Depending on the nature of corner cases and nesting, real parser might be preferable.

Conclusion

- ▶ Regular expressions
 - ▶ No error reporting
 - ▶ Flat data
- ▶ More general grammars
 - ▶ Composable structure (using combinator parsers)
 - ▶ Hierarchical, nested data
- ▶ Avoid reinventing

All materials for this talk available at <https://github.com/franklinchen/talk-on-overusing-regular-expressions>.

The hyperlinks on the slide PDFs are clickable.