



Fast, Functional Text Matching: Rosie Pattern Language

Dr. Jamie A. Jennings
Department of Computer Science
NC State University
27 April 2019

On the interwebs:
@jamietheriveter
<https://rosie-lang.org>
<https://gitlab.com/rosie-pattern-language>

Expression-based language

[a-z]

{ [a-z]+ "-" [0-9] }

ipv4 / ipv6

"INFO" (ipv4 / ipv6) hostname

{ (ipv4 / ipv6) port }

Expression-based language

Functions are values, can be composed

```
ci:"INFO" (ipv4 / ipv6) hostname
```

```
find:ci:"INFO" (ipv4 / ipv6) hostname
```


Expression-based language

Functions are values, can be composed

Lexical scope

```
grammar
  member = key ":" value
  object = "{" (member ("," member)*)? "}"
  array = "[" (value ("," value)*)? "]"
in
  value = ~ string / number / object / array / true / false / null
end
```

Expression-based language

Functions are values, can be composed

Lexical scope

Packages like in scheme48, Go, others

```
package json
```

```
import word, num
```

Expression-based language

Functions are values, can be composed

Lexical scope

Packages like in scheme48, Go, others

Compile-time meta-programming

```
$ rosie expand 'find:"INFO"'
Expression      find:"INFO"
Parses as      find:{"INFO"}
At top level    find:{"INFO"}
Expands to {
  grammar
    alias <search> = {!"INFO" .}*
    <anonymous> = {"INFO"}

  in
    alias find = {<search> <anonymous>}

  end
}
$
```


Expression-based language

Functions are values, can be composed

Lexical scope

Packages like in schem

Compile-time meta-pro

Prelude like in Haskell

```
$ rosie list
Rosie 1.1.0
```

Name	Cap?	Type	Color	Source
\$		pattern	default;bold	builtin/prelude
.		pattern	default;bold	builtin/prelude
^		pattern	default;bold	builtin/prelude
ci		macro		builtin/prelude
error		function		builtin/prelude
find		macro		builtin/prelude
findall		macro		builtin/prelude
kepto		macro		builtin/prelude
message		function		builtin/prelude
~		pattern	default;bold	builtin/prelude

```
10/10 names shown
$
```

Expression-based language

Functions are values, can be composed

Lexical scope

Packages like in scheme

Compile-time meta-pro

Prelude like in Haskell

Environment model
like a “Lisp-1”

```
$ rosie list
Rosie 1.1.0

Name                Cap?  Type      Color      Source
-----
$                   pattern  default;bold  builtin/prelude
.                   pattern  default;bold  builtin/prelude
^                   pattern  default;bold  builtin/prelude
ci                  macro    builtin/prelude
error              function builtin/prelude
find               macro    builtin/prelude
findall           macro    builtin/prelude
keepsto           macro    builtin/prelude
message           function builtin/prelude
~                  pattern  default;bold  builtin/prelude

10/10 names shown
$
```


Rosie Pattern Language

Expression-based language

Functions are values, can be composed

Lexical scope

Packages like in scheme48, Go, others

Compile-time meta-programming

Prelude like in Haskell

Environment model
like a “Lisp-1”



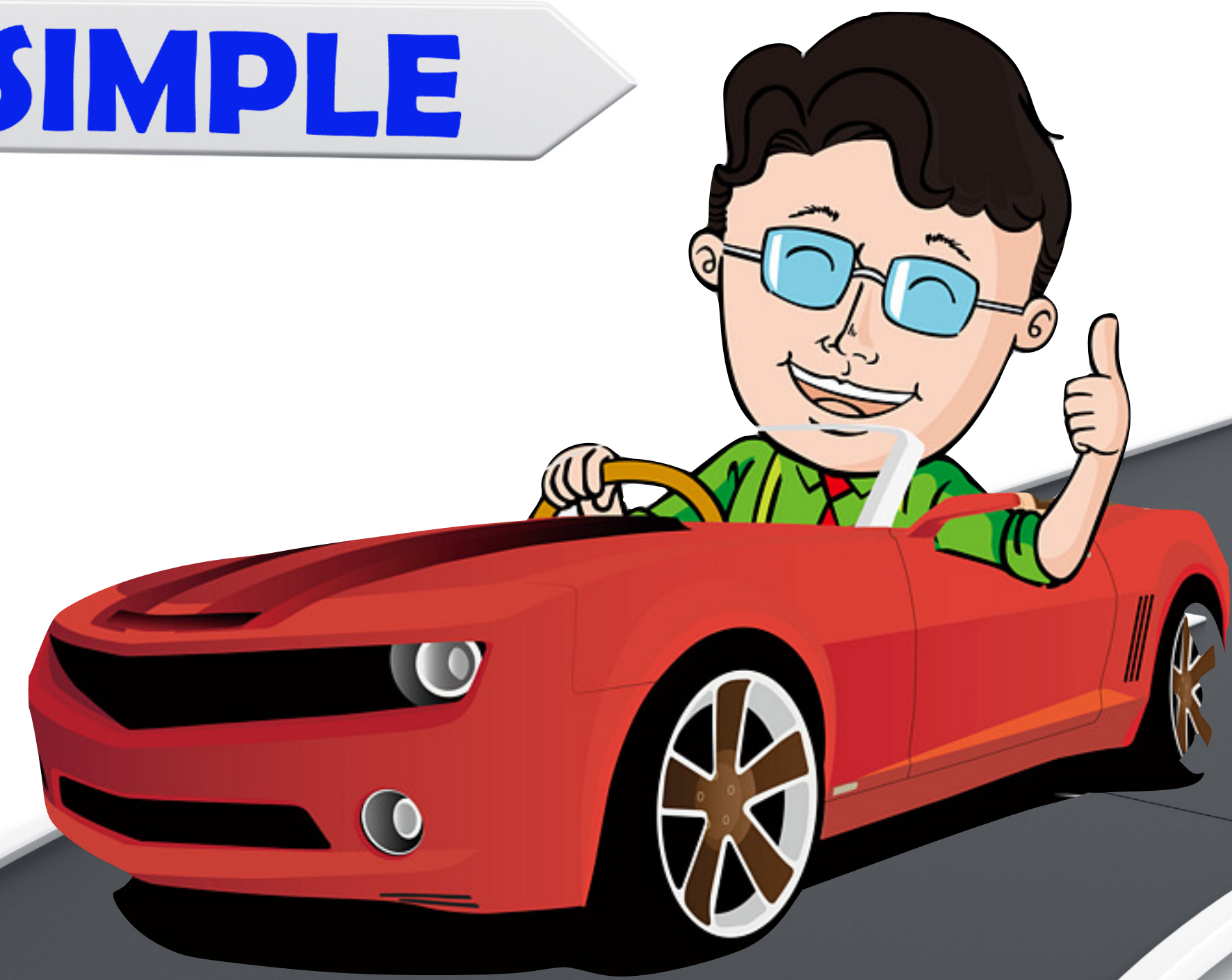
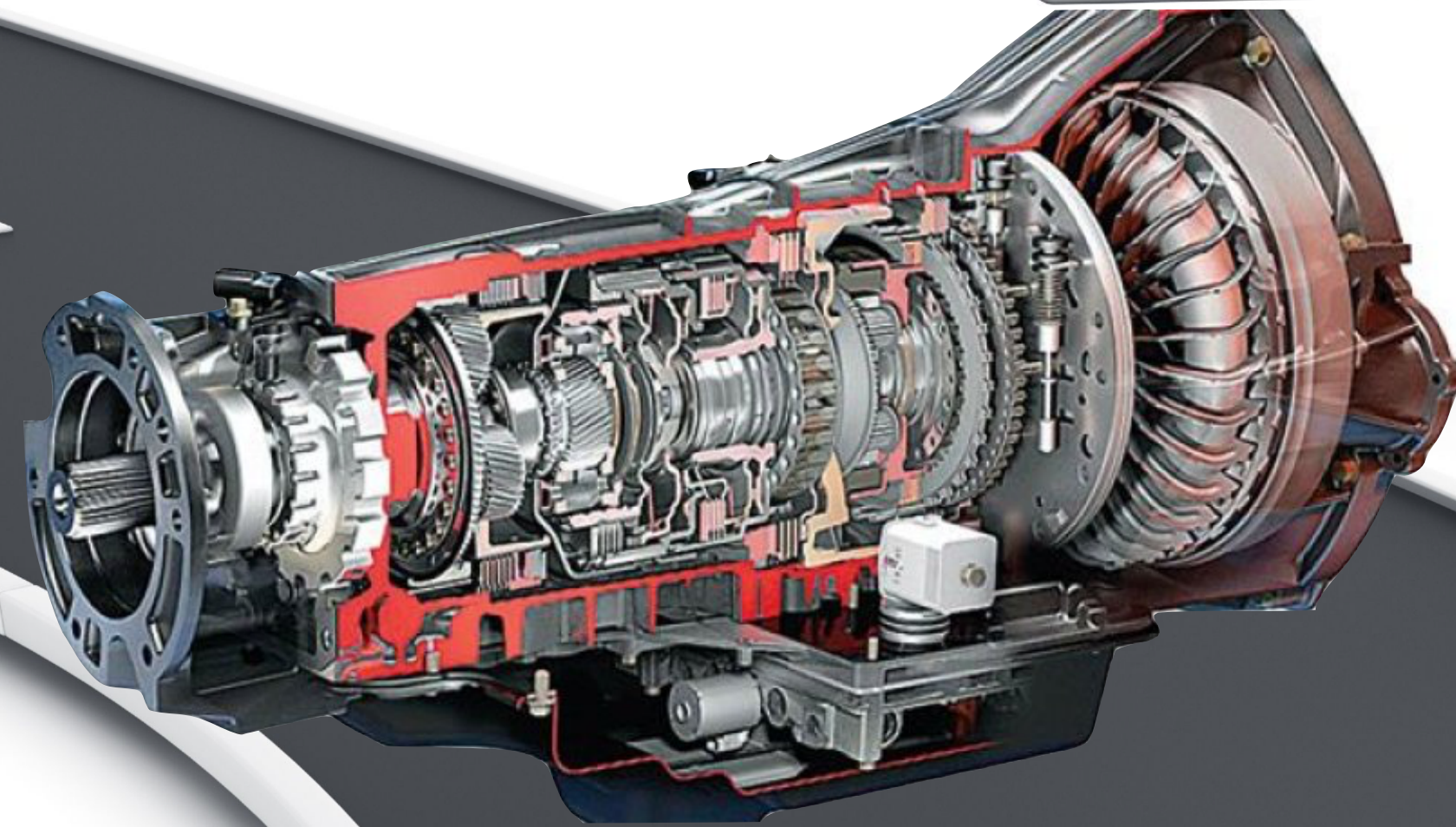
rosie-lang.org

But...

why???

SIMPLE

COMPLEX



Type inference

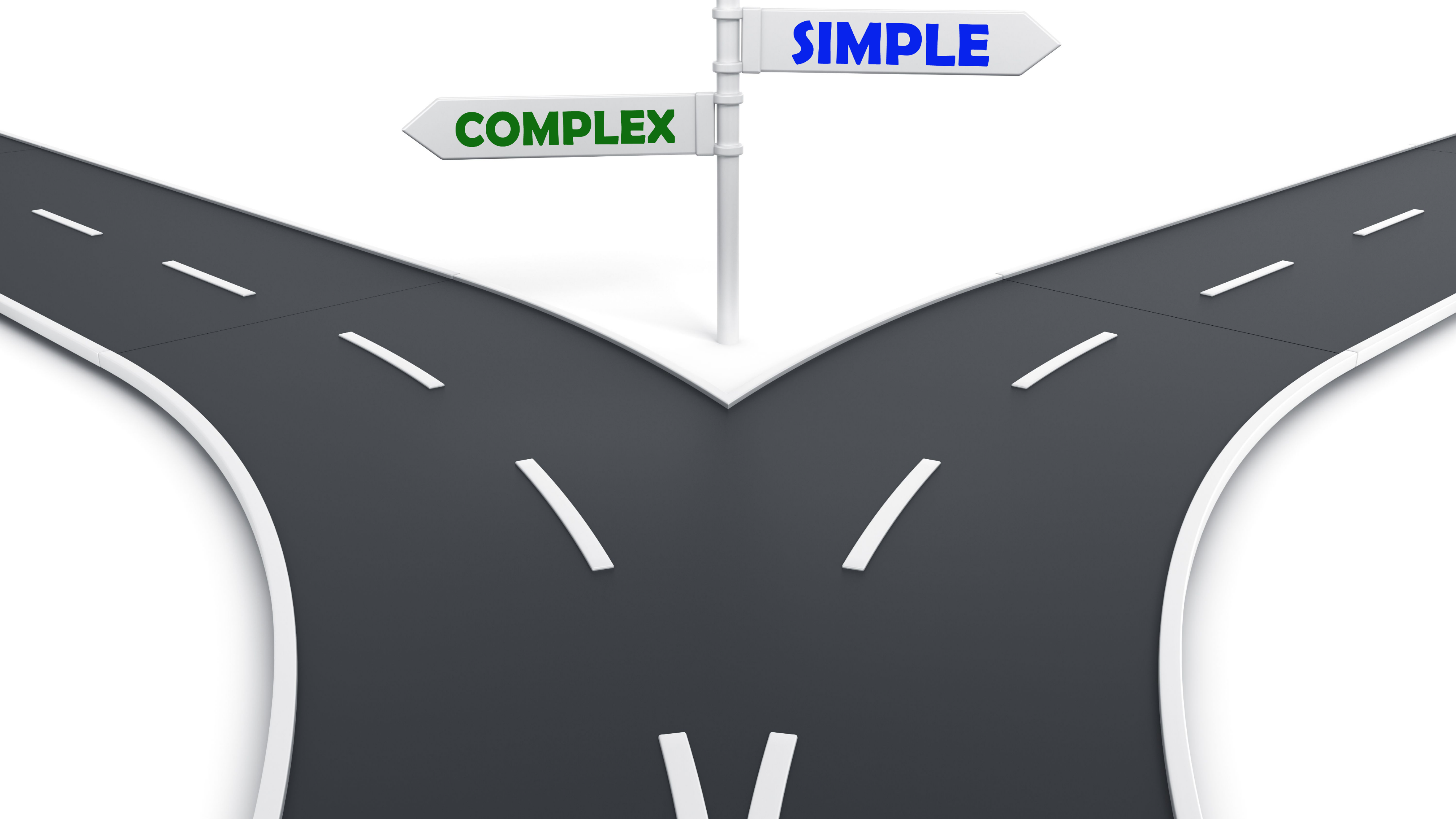
$\frac{x : \sigma \in \Gamma}{\Gamma \vdash_D x : \sigma}$	[Var]
$\frac{\Gamma \vdash_D e_0 : \tau \rightarrow \tau' \quad \Gamma \vdash_D e_1 : \tau}{\Gamma \vdash_D e_0 e_1 : \tau'}$	[App]
$\frac{\Gamma, x : \tau \vdash_D e : \tau'}{\Gamma \vdash_D \lambda x. e : \tau \rightarrow \tau'}$	[Abs]
$\frac{\Gamma \vdash_D e_0 : \sigma \quad \Gamma, x : \sigma \vdash_D e_1 : \tau}{\Gamma \vdash_D \text{let } x = e_0 \text{ in } e_1 : \tau}$	[Let]
$\frac{\Gamma \vdash_D e : \sigma' \quad \sigma' \sqsubseteq \sigma}{\Gamma \vdash_D e : \sigma}$	[Inst]
$\frac{\Gamma \vdash_D e : \sigma \quad \alpha \notin \text{free}(\Gamma)}{\Gamma \vdash_D e : \forall \alpha. \sigma}$	[Gen]

Using type inference

```
# let f x = x + 1;;  
val f : int -> int = <fun>  
#
```


SIMPLE

COMPLEX

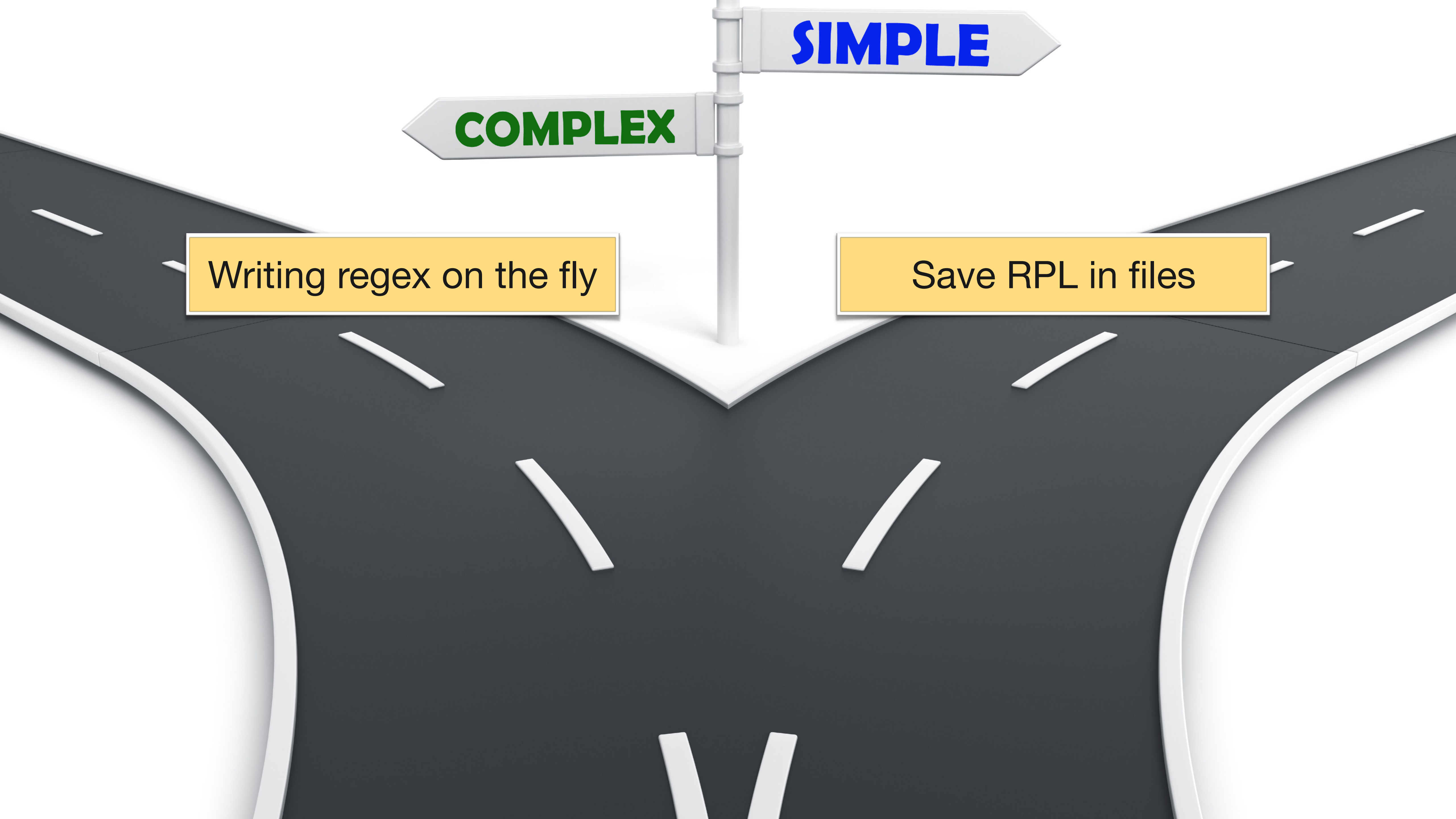


SIMPLE

COMPLEX

Writing regex on the fly

Save RPL in files



SIMPLE

COMPLEX

Writing regex on the fly

Save RPL in files

Reading cryptic syntax

PL-like syntax

SIMPLE

COMPLEX

Writing regex on the fly

Save RPL in files

Reading cryptic syntax

PL-like syntax

Exceptions to rules

Few special cases

SIMPLE

COMPLEX

Writing regex on the fly

Save RPL in files

Reading cryptic syntax

PL-like syntax

Exceptions to rules

Few special cases

Using an ad hoc
collection of tools

Tooling included (and
extensible)

Raison d'être



To Do List

1. Mine data from tools
2. Make predictions that help developers

Raison d'être



To Do List

1. Mine data from tools
2. Make predictions that help developers



Raison d'être



To Do List

1. Mine data from tools
2. Make predictions that help developers



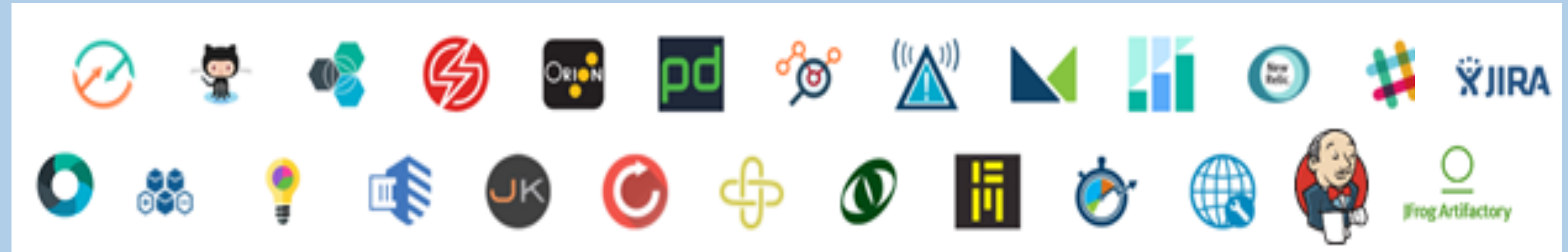
➔ My team had to write **lots of regex**

Raison d'être



To Do List

1. Mine data from tools
2. Make predictions that help developers



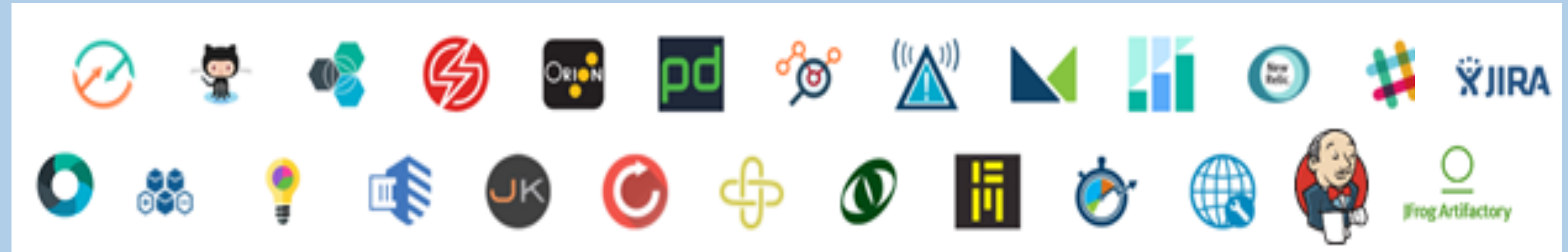
- ➔ My team had to write **lots of regex**
- ➔ We found that regex technology **does not scale**
 1. # of people, over time
 2. # of patterns
 3. data volume and velocity

Raison d'être



To Do List

1. Mine data from tools
2. Make predictions that help developers



- ➔ My team had to write **lots of regex**
- ➔ We found that regex technology **does not scale**
 1. # of people, over time
 2. # of patterns
 3. data volume and velocity
- ➔ So I designed **Rosie Pattern Language**

Raison d'être

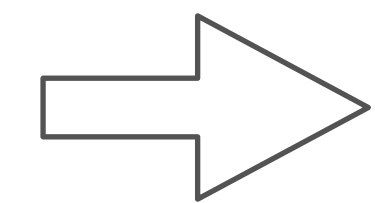


To Do List

1. Mine data from tools
2. Make predictions that help developers



- ➔ My team had to write **lots of regex**
- ➔ We found that regex technology **does not scale**
 1. # of people, over time
 2. # of patterns
 3. data volume and velocity
- ➔ So I designed **Rosie Pattern Language**

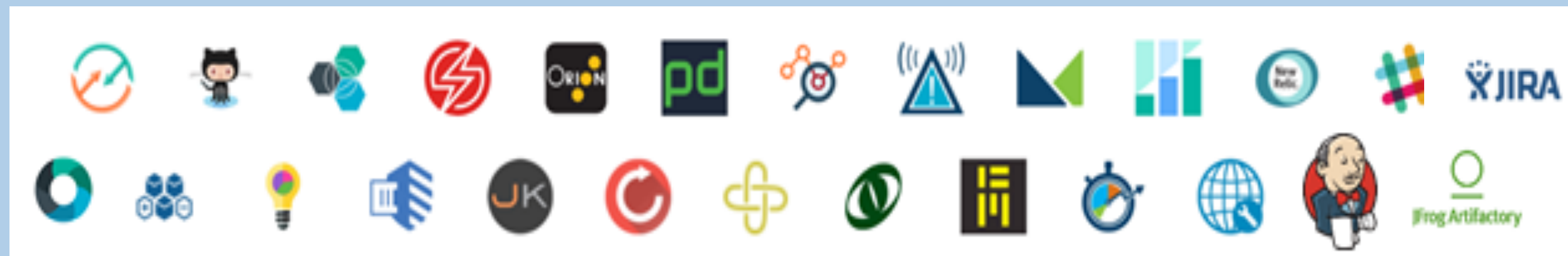


Raison d'être

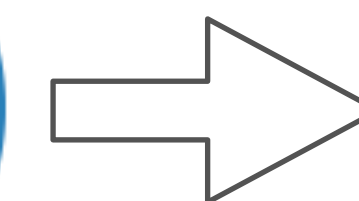
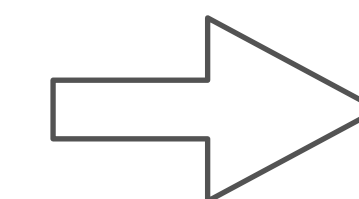


To Do List

1. Mine data from tools
2. Make predictions that help developers



- ➔ My team had to write **lots of regex**
- ➔ We found that regex technology **does not scale**
 1. # of people, over time
 2. # of patterns
 3. data volume and velocity
- ➔ So I designed **Rosie Pattern Language**



Current approach: regex

“If the only tool you have is a hammer...”

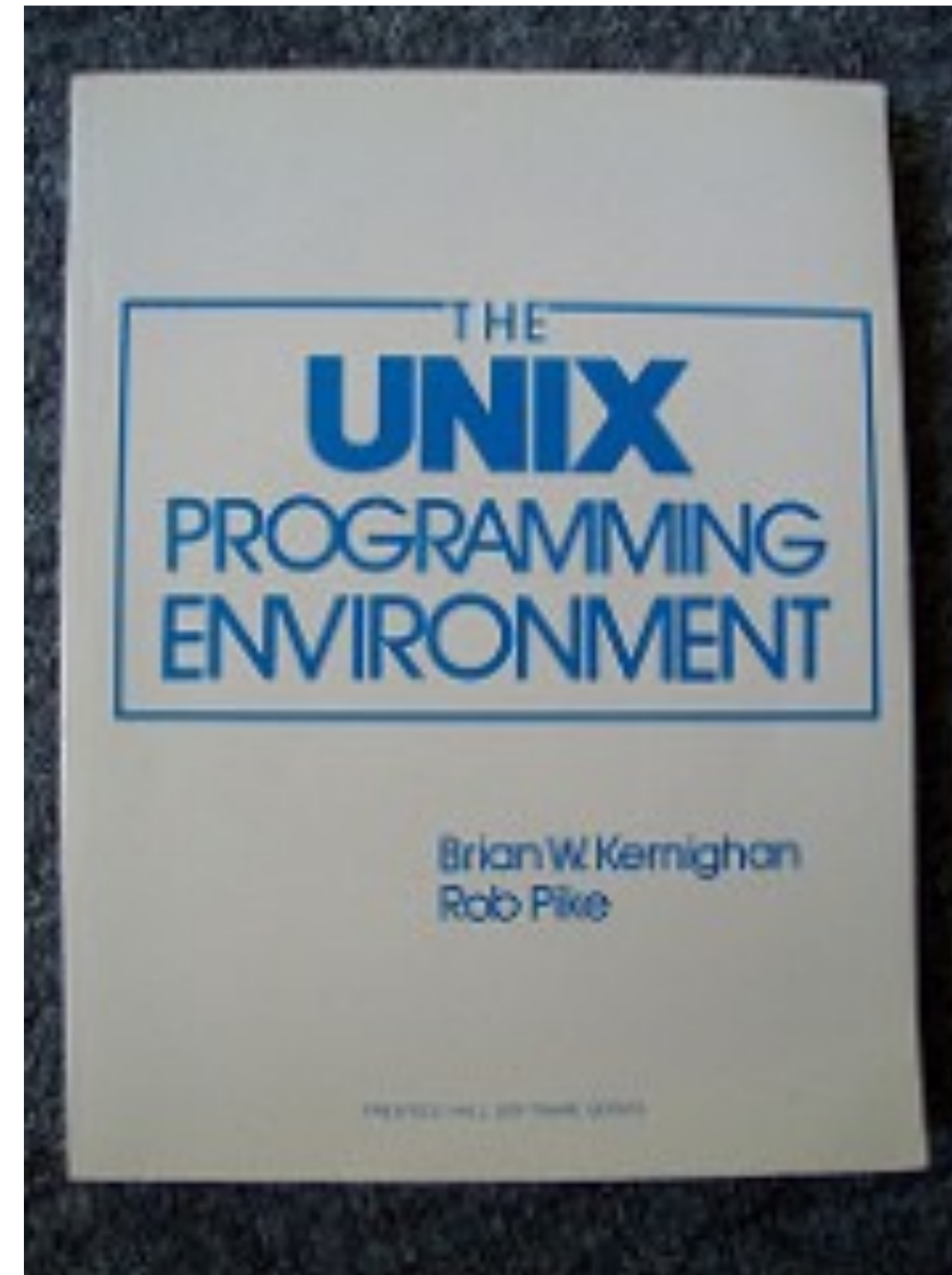
Abraham Maslow





Regular expressions as tools:

70s



On the command line:

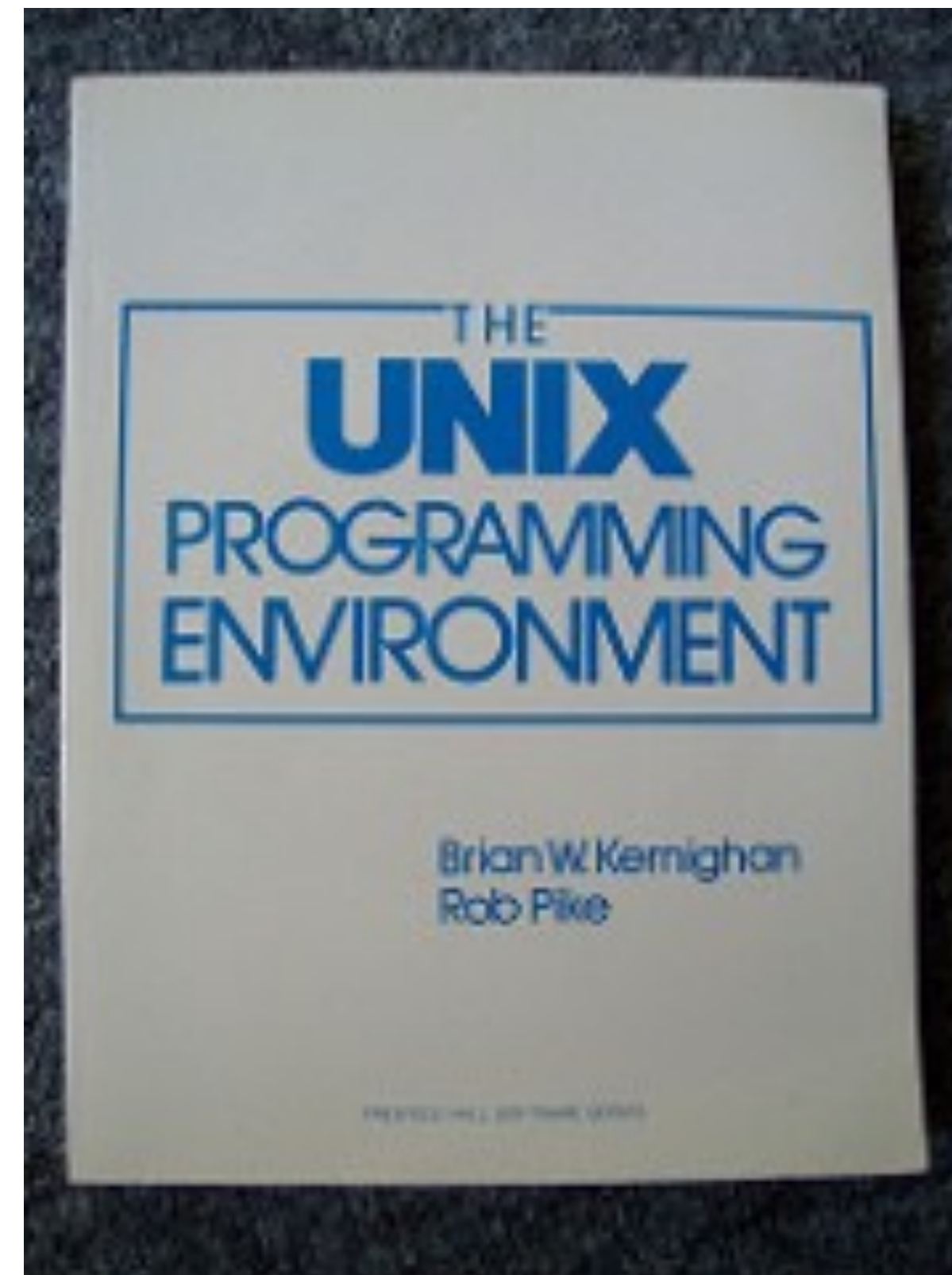
```
grep -v “^#\|^’\|^\/\|”
```

```
egrep -o '((\d{1,3})([.]\d{1,3}){2}|\w+(\.[.]\w+)+)'
```

```
sed -e ':a' -e 'N' -e '$!ba' -e 's/\n/ /g'
```


Regular expressions as tools:

70s



2017

Languages & Libraries

- [Boost](#)
- [Delphi](#)
- [GNU \(Linux\)](#)
- [Groovy](#)
- [Java](#)
- [JavaScript](#)
- [.NET](#)
- [PCRE \(C/C++\)](#)
- [PCRE2 \(C/C++\)](#)
- [Perl](#)
- [PHP](#)
- [POSIX](#)
- [PowerShell](#)
- [Python](#)
- [R](#)
- [Ruby](#)
- [std::regex](#)
- [Tcl](#)
- [VBScript](#)
- [Visual Basic 6](#)
- [wxWidgets](#)
- [XML Schema](#)
- [Xojo](#)
- [XQuery & XPath](#)
- [XRegExp](#)

<http://www.regular-expressions.info/tools.html>

On the command line:

```
grep -v “^#\|^’\|^\/\|”
```

```
egrep -o '((\d{1,3})([.]\d{1,3}){2}|\w+([.]\w+)+)'
```

```
sed -e ':a' -e 'N' -e '$!ba' -e 's/\n/ /g'
```

Regex are notoriously hard to read & maintain

- Dense, **cryptic** syntax
- Semantics **vary** across implementations
- Flags that **affect** the semantics are *not part of the pattern*
- Regex **do not easily compose**



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

Jamie Zawinski

<http://regex.info/blog/2006-09-15/247>





Regular expressions

Match a date with slashes, like 1/1/1970:

```
^\d{1,2}\/\d{1,2}\/\d{4}$
```

Match an email address (obviously!):

```
^(?>[a-zA-Z\d!#$%&'*+\/=?^_`{|}~]+\x20*|"(?  
=[\x01-\x7f])["\]|\\[\x01-\x7f])*"\x20*)*(?  
<angle><))?(?!\.)(?>\.?[a-zA-Z\d!#$%&'*+\/=?  
^_`{|}~]+)+|"(?=[\x01-\x7f])["\]|\\[\x01-  
\x7f])*")@((?!-)[a-zA-Z\d\_-]+(?<!--)\.)+[a-zA-Z]  
{2,}|\[((?<!--\[\]\.)(25[0-5]|2[0-4]\d|[01]?\d?  
\d))\]{4}|[a-zA-Z\d\_-]*[a-zA-Z\d]:((=[\x01-\x7f])  
[^\[\]\]|\\[\x01-\x7f])+)\](?(angle)>)$
```

Regular expressions

Rosie Pattern Language

Match a date with slashes, like 1/1/1970:

```
^\d{1,2}\/\d{1,2}\/\d{4}$
```

```
date.slashed
```

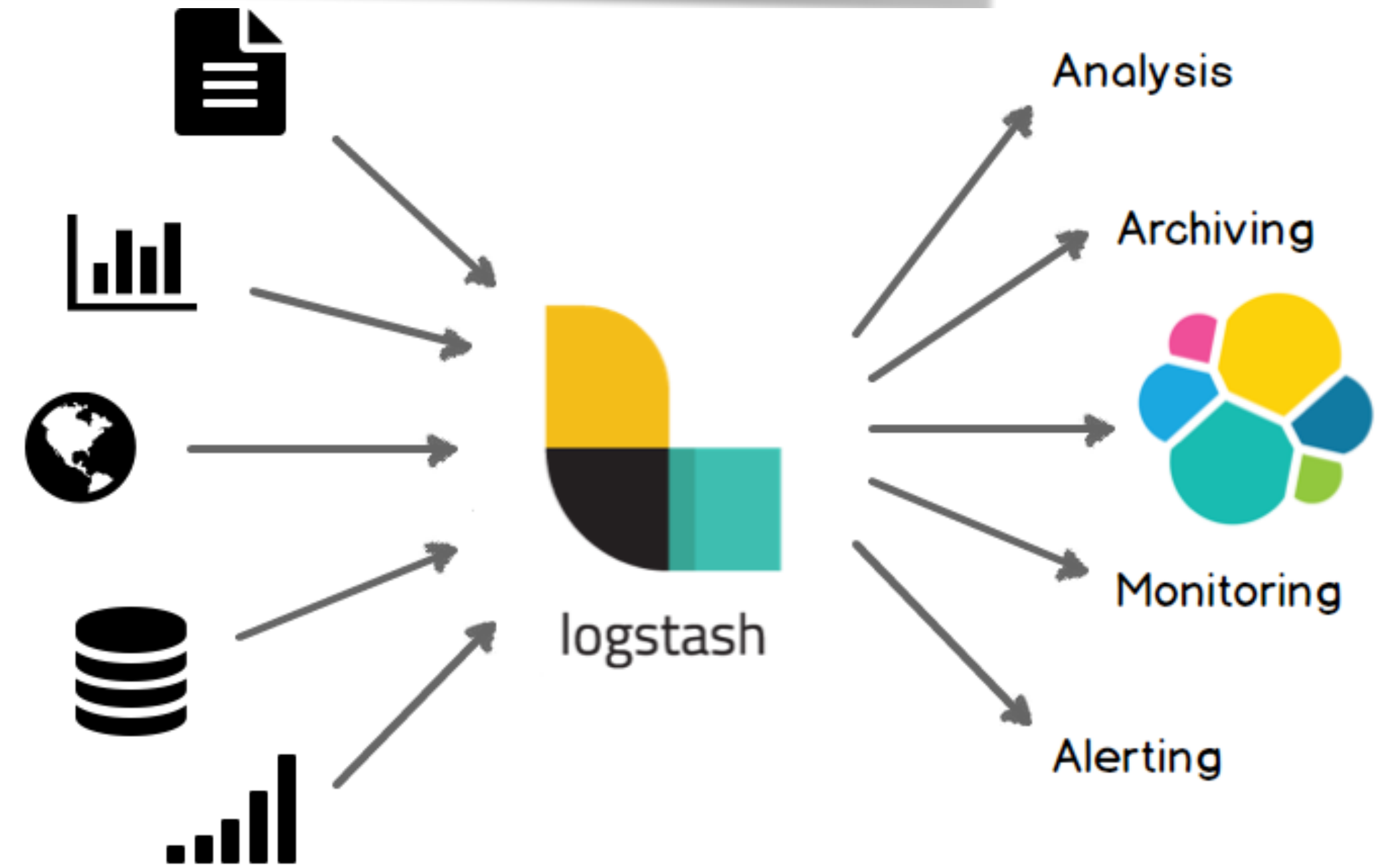
Match an email address (obviously!):

```
^(?>[a-zA-Z\d!#$%&'*+\/=?^_`{|}~]+\x20*|"(?  
=[\x01-\x7f])["\]|\\[\x01-\x7f])*"\x20*)*(?  
<angle><))?(?!\.)(?>\.?[a-zA-Z\d!#$%&'*+\/=?  
^_`{|}~]+)+|"((?=[\x01-\x7f])["\]|\\[\x01-  
\x7f])*")@((?!-)[a-zA-Z\d\_-]+(?<!--)\.)+[a-zA-Z]  
{2,}|\[((?<!--\[\]\.)(25[0-5]|2[0-4]\d|[01]?\d?  
\d))\{4\}|[a-zA-Z\d\_-]*[a-zA-Z\d]:((?=[\x01-\x7f])  
["\]|\\[\x01-\x7f])+\)](?<angle>)$
```

```
net.email
```

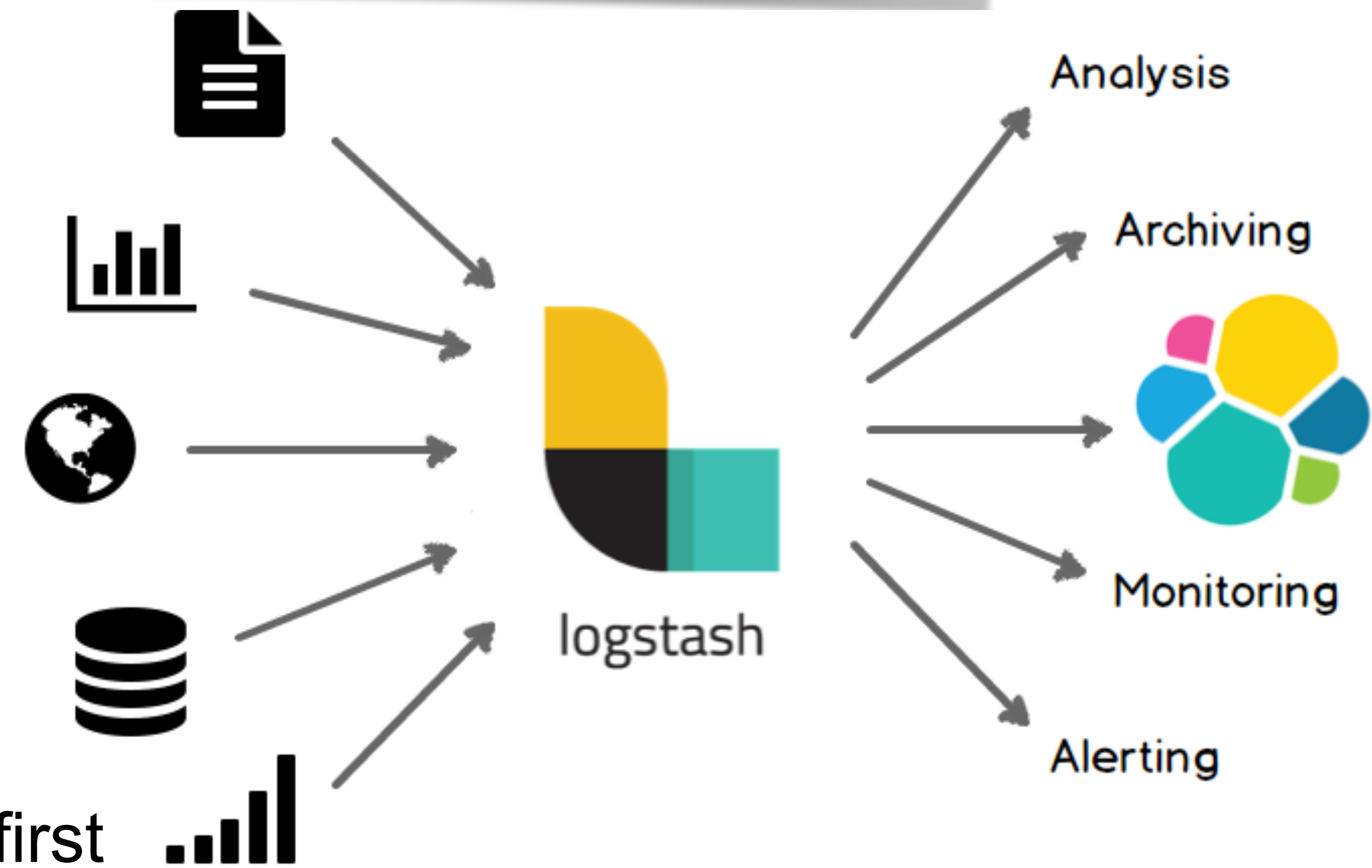

Other regex collections? Grok does this.

Grok sits on top of regular expressions, so any regular expressions are valid in grok as well. The regular expression library is Oniguruma, and you can see the full supported regexp syntax [on the Oniguruma site](#).



Other regex collections? Grok does this.

Grok sits on top of regular expressions, so any regular expressions are valid in grok as well. The regular expression library is Oniguruma, and you can see the full supported regexp syntax [on the Oniguruma site](#).



Caveats

- ✦ Name collisions? Some versions will use the first one, some the last
- ✦ No packages, hierarchy, or dependencies
- ✦ They are still **unreadable** and **unmaintainable!**

And they don't play well with dev tools

```
grok$ diff orig copy
18c18
< QUOTEDSTRING (?>(?!\\)(?>"(?>\\. | [^\\"]+)+|"|(?>'(?>\\. | [^\\']+)+')|'|(?>>`(?>\\. | [^\\`]+)+`)|``))
---
> QUOTEDSTRING (?>(?(?!\\)(?>"(?>\\. | [^\\"]+)+|"|(?>'(?>\\. | [^\\']+)+')|'|(?>>`(?>\\. | [^\\`]+)+`)|``))
26c26
< IPV6 ((([0-9A-Fa-f]{1,4}:){7}([0-9A-Fa-f]{1,4}|:))|(([0-9A-Fa-f]{1,4}:){6}(:[0-9A-Fa-f]{1,4}|((25[0-5]|2[0-4]\\d|1\\d\\d|[1-9]
?\\d)(\\. (25[0-5]|2[0-4]\\d|1\\d\\d|[1-9]?\\d)){3})|:))|(([0-9A-Fa-f]{1,4}:){5}((:[0-9A-Fa-f]{1,4}){1,2})|:(25[0-5]|2[0-4]\\d|1\\d\\
d|[1-9]?\\d)(\\. (25[0-5]|2[0-4]\\d|1\\d\\d|[1-9]?\\d)){3})|:))|(([0-9A-Fa-f]{1,4}:){4}((:[0-9A-Fa-f]{1,4}){1,3})|((:[0-9A-Fa-f]{1,
4})?:((25[0-5]|2[0-4]\\d|1\\d\\d|[1-9]?\\d)(\\. (25[0-5]|2[0-4]\\d|1\\d\\d|[1-9]?\\d)){3}))|:))|(([0-9A-Fa-f]{1,4}:){3}((:[0-9A-Fa-f]{
1,4}){1,4})|((:[0-9A-Fa-f]{1,4}){0,2}:((25[0-5]|2[0-4]\\d|1\\d\\d|[1-9]?\\d)(\\. (25[0-5]|2[0-4]\\d|1\\d\\d|[1-9]?\\d)){3}))|:))|(([0-9
A-Fa-f]{1,4}:){2}((:[0-9A-Fa-f]{1,4}){1,5})|((:[0-9A-Fa-f]{1,4}){0,3}:((25[0-5]|2[0-4]\\d|1\\d\\d|[1-9]?\\d)(\\. (25[0-5]|2[0-4]\\d
|1\\d\\d|[1-9]?\\d)){3}))|:))|(([0-9A-Fa-f]{1,4}:){1}((:[0-9A-Fa-f]{1,4}){1,6})|((:[0-9A-Fa-f]{1,4}){0,4}:((25[0-5]|2[0-4]\\d|1\\
d\\d|[1-9]?\\d)(\\. (25[0-5]|2[0-4]\\d|1\\d\\d|[1-9]?\\d)){3}))|:))|(:((:[0-9A-Fa-f]{1,4}){1,7})|((:[0-9A-Fa-f]{1,4}){0,5}:((25[0-5]
|2[0-4]\\d|1\\d\\d|[1-9]?\\d)(\\. (25[0-5]|2[0-4]\\d|1\\d\\d|[1-9]?\\d)){3}))|:)))((%.+)?
---
> IPV6 ((([0-9A-Fa-f]{1,4}:){7}([0-9A-Fa-f]{1,4}|:))|(([0-9A-Fa-f]{1,4}:){6}(:[0-9A-Fa-f]{1,4}|((25[0-5]|2[0-4]\\d|1\\d\\d|[1-9]
?\\d)(\\. (25[0-5]|2[0-4]\\d|1\\d\\d|[1-9]?\\d)){3})|:))|(([0-9A-Fa-f]{1,4}:){5}((:[0-9A-Fa-f]{1,4}){1,2})|:(25[0-5]|2[0-4]\\d|1\\d\\
d|[1-9]?\\d)(\\. (25[0-5]|2[0-4]\\d|1\\d\\d|[1-9]?\\d)){3})|:))|(([0-9A-Fa-f]{1,4}:){4}((:[0-9A-Fa-f]{1,4}){1,3})|((:[0-9A-Fa-f]{1,
4})?:((25[0-5]|2[0-4]\\d|1\\d\\d|[1-9]?\\d)(\\. (25[0-5]|2[0-4]\\d|1\\d\\d|[1-9]?\\d)){3}))|:))|(([0-9A-Fa-f]{1,4}:){3}((:[0-9A-Fa-f]{
1,4}){1,4})|((:[0-9A-Fa-f]{1,4}){0,3}:((25[0-5]|2[0-4]\\d|1\\d\\d|[1-9]?\\d)(\\. (25[0-5]|2[0-4]\\d|1\\d\\d|[1-9]?\\d)){3}))|:))|(([0-9
A-Fa-f]{1,4}:){2}((:[0-9A-Fa-f]{1,4}){1,5})|((:[0-9A-Fa-f]{1,4}){0,3}:((25[0-5]|2[0-4]\\d|1\\d\\d|[1-9]?\\d)(\\. (25[0-5]|2[0-4]\\d
|1\\d\\d|[1-9]?\\d)){3}))|:))|(([0-9A-Fa-f]{1,4}:){1}((:[0-9A-Fa-f]{1,4}){1,6})|((:[0-9A-Fa-f]{1,4}){0,3}:((25[0-5]|2[0-4]\\d|1\\
d\\d|[1-9]?\\d)(\\. (25[0-5]|2[0-4]\\d|1\\d\\d|[1-9]?\\d)){3}))|:))|(:((:[0-9A-Fa-f]{1,4}){1,7})|((:[0-9A-Fa-f]{1,4}){0,5}:((25[0-5]
|2[0-4]\\d|1\\d\\d|[1-9]?\\d)(\\. (25[0-5]|2[0-4]\\d|1\\d\\d|[1-9]?\\d)){3}))|:)))((%.+)?
grok$
```


And they don't play well with dev tools

```
grok$ diff orig copy
```

```
18c18
```

```
< QUOTEDSTRING (?>(
```

```
---
```

```
> QUOTEDSTRING (?>(
```

```
26c26
```

```
< IPV6 ((([0-9A-Fa-
```

```
?\d)(\.(25[0-5]|2[0-
```

```
d|[1-9]?\d)(\.(25[0-
```

```
4})?:((25[0-5]|2[0-
```

```
1,4}){1,4})|((:[0-9
```

```
A-Fa-f){1,4}:){2}((
```

```
|1\d\d|[1-9]?\d)){3
```

```
d\d|[1-9]?\d)(\.(25
```

```
|2[0-4]\d|1\d\d|[1-
```

```
---
```

```
> IPV6 ((([0-9A-Fa-
```

```
?\d)(\.(25[0-5]|2[0-
```

```
d|[1-9]?\d)(\.(25[0-
```

```
4})?:((25[0-5]|2[0-
```

```
1,4}){1,4})|((:[0-9
```

```
A-Fa-f){1,4}:){2}((
```

```
|1\d\d|[1-9]?\d)){3
```

```
d\d|[1-9]?\d)(\.(25
```

```
|2[0-4]\d|1\d\d|[1-
```

```
grok$
```

RPL has syntax like a programming language

- It reads like code
- It diffs like code
- It debugs like code

```
0-4]\d|1\d\d|[1-9]
```

```
0-5]|2[0-4]\d|1\d\
```

```
|((:[0-9A-Fa-f){1,
```

```
3}(((:[0-9A-Fa-f){
```

```
d)){3}))|:))|(([0-9
```

```
\.(25[0-5]|2[0-4]\d
```

```
25[0-5]|2[0-4]\d|1\
```

```
4}){0,5}:((25[0-5]
```

```
0-4]\d|1\d\d|[1-9]
```

```
0-5]|2[0-4]\d|1\d\
```

```
|((:[0-9A-Fa-f){1,
```

```
3}(((:[0-9A-Fa-f){
```

```
))){3}))|:))|(([0-9
```

```
\.(25[0-5]|2[0-4]\d
```

```
25[0-5]|2[0-4]\d|1\
```

```
4}){0,5}:((25[0-5]
```

```
0-4]\d|1\d\d|[1-9]
```

```
0-5]|2[0-4]\d|1\d\
```

```
|((:[0-9A-Fa-f){1,
```

```
3}(((:[0-9A-Fa-f){
```

```
))){3}))|:))|(([0-9
```

```
\.(25[0-5]|2[0-4]\d
```

```
25[0-5]|2[0-4]\d|1\
```






Regex performance is surprisingly variable



Regular expression matching can be very efficient: linear time in the size of the input.



“The worst-case exponential-time backtracking strategy [is] used almost everywhere [but grep and RE2], including ed, sed, Perl, PCRE, and Python.”

(Russ Cox <https://swtch.com/~rsc/regexp/regexp2.html>)

Regex performance is surprisingly variable



M

In RPL, expressions are greedy and possessive.

- Backtracking is explicit
- To get exponential backtracking, you write it that way
- Today (v1.1.x) such RPL patterns have exponential size

A

RPL makes it difficult to be accidentally inefficient.

Perl*

?a?a?a?

Perl*

4,25,26,

pre = (. : ,) 25, 26, ,

(*) Perl 5.16.3 darwin-thread-multi-2level

Rosie Pattern Language

“All progress depends on the unreasonable [woman]”

George Bernard Shaw, paraphrased

RPL is designed like a programming language

```
----- -*- Mode: rpl; -*-
-----
----- json.rpl    rpl patterns for processing json input
-----
----- © Copyright IBM Corporation 2016, 2017.
----- LICENSE: MIT License (https://opensource.org/licenses/mit-license.html)
----- AUTHOR: Jamie A. Jennings

package json

import word, num

local key = word.dq
local string = word.dq
local number = num.signed_number

local true = "true"
local false = "false"
local null = "null"

grammar
  value = ~ string / number / object / array / true / false / null
  member = key ":" value
  object = "{" (member ("," member)*)? "}"
  array = "[" (value ("," value)*)? "]"
end

-- test value accepts "true", "false", "null"
-- test value rejects "ture", "f", "NULL"
-- test value accepts "0", "123", "-1", "1.1001", "1.2e10", "1.2e-10", "+3.3"
-- test value accepts "\"hello\"", "\"this string has \\\"embedded\\\" double quotes\""
-- test value rejects "hello", "\"this string has no \\\"final quote\\\" "
-- test value rejects "--2", "9.1.", "9.1.2", "++2", "2E02."

-- test value accepts "[]", "[1, 2, 3.14, \"V\", 6.02e23, true]", "[1, 2, [7], [[8]]]"
-- test value rejects "[]", "[", "[[]", "{1, 2}"

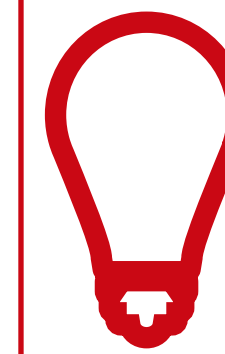
-- test value accepts "{\"one\":1}", "{ \"one\" :1}", "{ \"one\" : 1  }"
-- test value accepts "{\"one\":1, \"two\": 2}", "{\"one\":1, \"two\": 2, \"array\":[1,2]}"
-- test value accepts "[{\"v\":1}, {\"v\":2}, {\"v\":3}]"
```

RPL is designed like a programming language

Comments
Modules
Identifiers
Whitespace
Quoted literals
Unit tests
Macros
(not shown)

```
----- -*- Mode: rpl; -*-  
-----  
----- json.rpl    rpl patterns for processing json input  
-----  
----- © Copyright IBM Corporation 2016, 2017.  
----- LICENSE: MIT License (https://opensource.org/licenses/mit-license.html)  
----- AUTHOR: Jamie A. Jennings  
  
package json  
  
import word, num  
  
local key = word.dq  
local string = word.dq  
local number = num.signed_number  
  
local true = "true"  
local false = "false"  
local null = "null"  
  
grammar  
  value = ~ string / number / object / array / true / false / null  
  member = key ":" value  
  object = "{" (member ("," member)*)? "  
  array = "[" (value ("," value)*)? "]"  
end  
  
-- test value accepts "true", "false", "null"  
-- test value rejects "ture", "f", "NULL"  
-- test value accepts "0", "123", "-1", "1.1001", "1.2e10", "1.2e-10", "+3.3"  
-- test value accepts "\"hello\"", "\"this string has \\\"embedded\\\" double quotes\""  
-- test value rejects "hello", "\"this string has no \\\"final quote\\\" "  
-- test value rejects "--2", "9.1.", "9.1.2", "++2", "2E02."  
  
-- test value accepts "[]", "[1, 2, 3.14, \"V\", 6.02e23, true]", "[1, 2, [7], [[8]]]"  
-- test value rejects "[[]]", "[", "[[]", "{1, 2}"  
  
-- test value accepts "{\"one\":1}", "{ \"one\" :1}", "{ \"one\" : 1 }"  
-- test value accepts "{\"one\":1, \"two\": 2}", "{\"one\":1, \"two\": 2, \"array\":[1,2]}"  
-- test value accepts "[{\"v\":1}, {\"v\":2}, {\"v\":3}]"
```


Can your 'grep' do this?



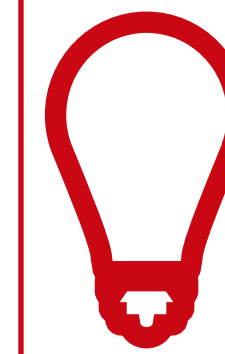
NAMED PATTERNS

```
$ curl -s www.google.com | rosie grep -o subs net.url  
http://schema.org/WebPage  
http://www.google.com/imghp?hl=en&tab=wi  
http://maps.google.com/maps?hl=en&tab=w1  
https://play.google.com/?hl=en&tab=w8  
http://www.youtube.com/?gl=US&tab=w1  
http://news.google.com/nwshp?hl=en&tab=wn  
https://mail.google.com/mail/?tab=wm  
https://drive.google.com/?tab=wo  
https://www.google.com/intl/en/options/  
http://www.google.com/history/optout?hl=en  
https://accounts.google.com/ServiceLogin?hl=en&passive=true&continue=http://www.google.com/  
https://plus.google.com/116899029375914044550  
$
```

-o Output format
subs ==> sub-matches

pattern is net.url
==> namespace net, pattern url

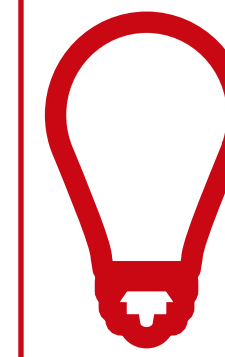
Can your 'grep' do this?



CUSTOMIZABLE
OUTPUT
HIGHLIGHTING

```
$ rosie match 'word.any (net.any)+' resolv.conf
domain abc.aus.example.com
search ibm.com mylocaldomain.myisp.net example.com
nameserver 192.9.201.1
nameserver 192.9.201.2
nameserver fde9:4789:96dd:03bd::1
$
```


Can your 'grep' do this?

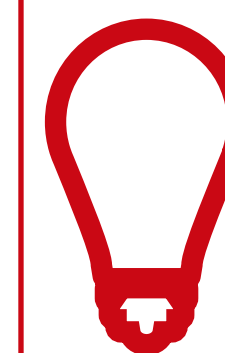


CUSTOMIZABLE
OUTPUT
HIGHLIGHTING

```
$ rosie match 'word.any (net.any)+' resolv.conf
domain abc.aus.example.com
search ibm.com mylocaldomain.myisp.net example.com
nameserver 192.9.201.1
nameserver 192.9.201.2
nameserver fde9:4789:96dd:03bd::1
$
```

```
$ rosie --colors='net.ipv4=blue:bold' match 'word.any (net.any)+' resolv.conf
domain abc.aus.example.com
search ibm.com mylocaldomain.myisp.net example.com
nameserver 192.9.201.1
nameserver 192.9.201.2
nameserver fde9:4789:96dd:03bd::1
$
```


Can your 'grep' do this?



CUSTOMIZABLE
OUTPUT
HIGHLIGHTING

```
$ sed -n 46,49p /var/log/system.log
```

```
Jul 30 10:18:42 Jamies-Compabler com.apple.xpc.launchd[1] (com.apple.CoreSimulator.CoreSimulatorService [669]): Service exited due to signal: Killed: 9 sent by com.apple.CoreSimulator.CoreSimu[669]
```

```
Jul 30 10:18:42 Jamies-Compabler systemstats[71]: assertion failed: 17G65: systemstats + 914800 [D1E75C38-62CE-3D77-9ED3-5F6D38EF0676]: 0x40
```

```
Jul 30 10:18:43 Jamies-Compabler ContainerMetadataExtractor[92065]: objc[92065]: Class BRMangledID is implemented in both /System/Library/PrivateFrameworks/CloudDocs.framework/Versions/A/CloudDocs (0x7fff8b848c88) and /System/Library/PrivateFrameworks/CloudDocsDaemon.framework/XPCServices/ContainerMetadataExtractor.xpc/Contents/MacOS/ContainerMetadataExtractor (0x10a8e0528). One of the two will be used. Which one is undefined.
```

```
Jul 30 10:18:50 Jamies-Compabler systemstats[71]: assertion failed: 17G65: systemstats + 914800 [D1E75C38-62CE-3D77-9ED3-5F6D38EF0676]: 0x40
```

```
$
```

```
$ sed -n 46,49p /var/log/system.log | rosie match all.things
```

```
Jul 30 10:18:42 Jamies-Compabler com.apple.xpc.launchd[1] (com.apple.CoreSimulator.CoreSimulatorService [669]): Service exited due to signal: Killed: 9 sent by com.apple.CoreSimulator.CoreSimu[669]
```

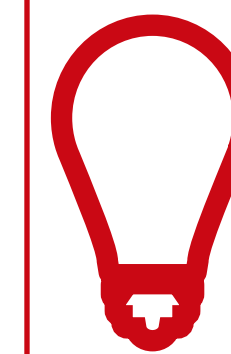
```
Jul 30 10:18:42 Jamies-Compabler systemstats[71]: assertion failed: 17G65: systemstats + 914800 [D1E75C38-62CE-3D77-9ED3-5F6D38EF0676]: 0x40
```

```
Jul 30 10:18:43 Jamies-Compabler ContainerMetadataExtractor[92065]: objc[92065]: Class BRMangledID is implemented in both /System/Library/PrivateFrameworks/CloudDocs.framework/Versions/A/CloudDocs (0x7fff8b848c88) and /System/Library/PrivateFrameworks/CloudDocsDaemon.framework/XPCServices/ContainerMetadataExtractor.xpc/Contents/MacOS/ContainerMetadataExtractor (0x10a8e0528). One of the two will be used. Which one is undefined.
```

```
Jul 30 10:18:50 Jamies-Compabler systemstats[71]: assertion failed: 17G65: systemstats + 914800 [D1E75C38-62CE-3D77-9ED3-5F6D38EF0676]: 0x40
```

```
$
```


Can your 'grep' do this?



**STRUCTURED
OUTPUT OPTION**

```
$ head -n 1 /var/log/system.log | rosie grep -o jsonpp num.denoted_hex
```

```
{ "s": 1,  
  "e": 80,  
  "data": "Jul 29 16:17:13 Jamies-Compabler timed[90268]: settimeofday({0x5b5e20c9,0x75bd3",  
  "subs":  
    [ { "s": 62,  
        "e": 72,  
        "data": "0x5b5e20c9",  
        "subs":  
          [ { "s": 64,  
              "e": 72,  
              "data": "5b5e20c9",  
              "type": "num.hex" } ],  
        "type": "num.denoted_hex" },  
      { "s": 73,  
        "e": 80,  
        "data": "0x75bd3",  
        "subs":  
          [ { "s": 75,  
              "e": 80,  
              "data": "75bd3",  
              "type": "num.hex" } ],  
        "type": "num.denoted_hex" } ],  
  "type": "*" }
```

Matching line

num.denoted_hex

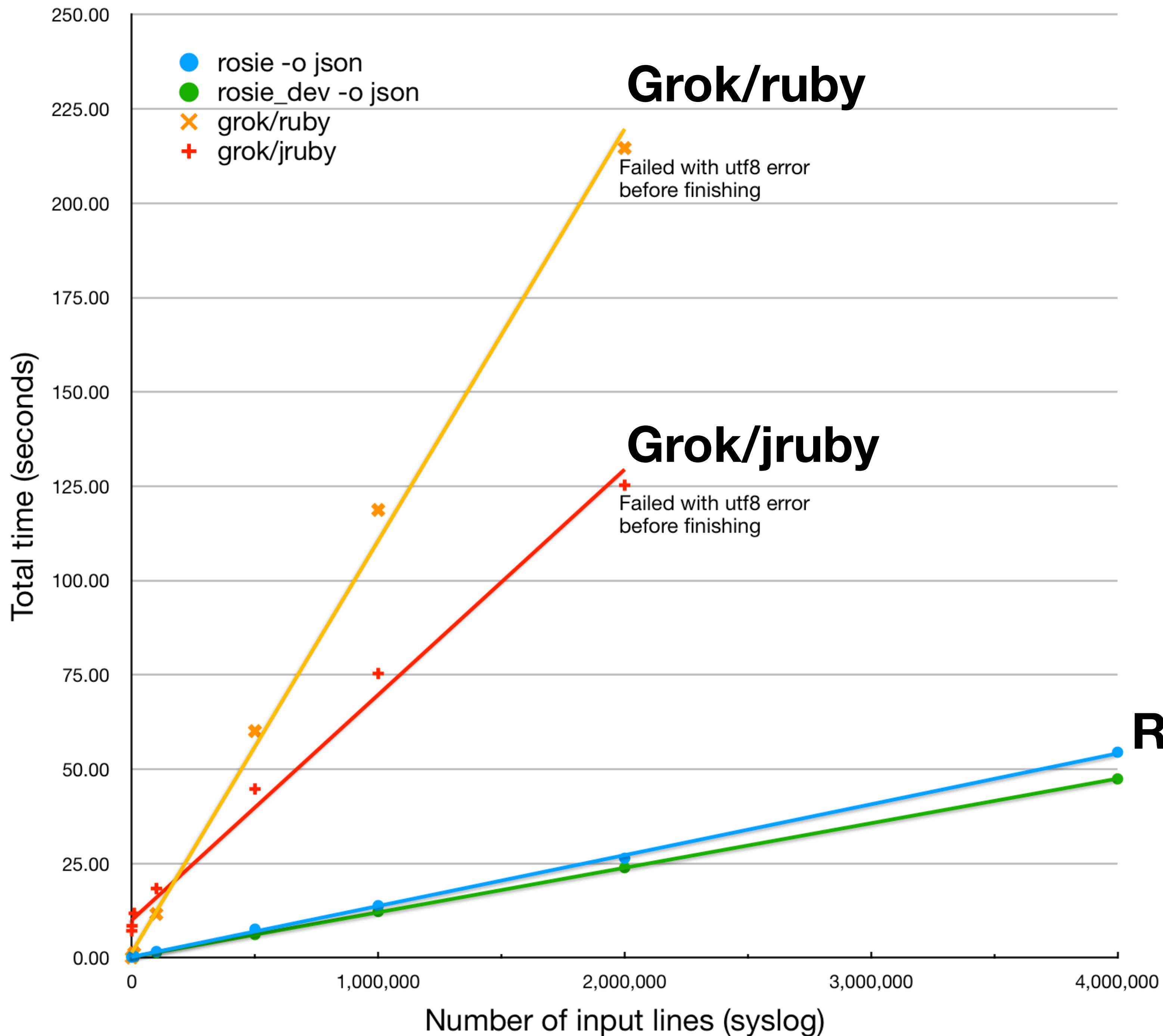
num.hex, a sub-match

\$

Performance

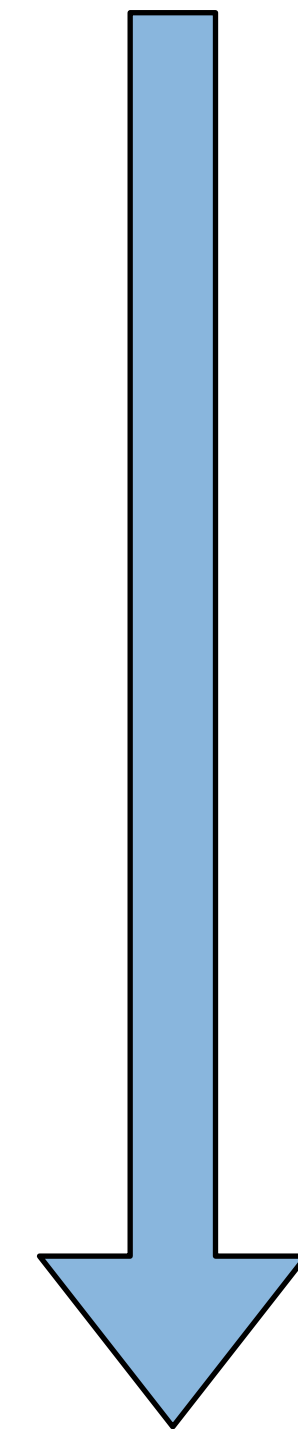
“I want to believe”

Fox Mulder, FBI



Performance

Worse



Better

Notes:

1. Log entry parsing is one narrow use case.
2. Hard to design fair comparisons.
3. Rosie output is nested JSON; Grok output is flat lists.

Debugging

“To err is human, but to really foul things up you need a computer.”

Paul R. Ehrlich

Trace a (mis-)match

```
$ date | rosie match date.us_dashed
```

```
$
```




```
$ date | rosie match date.us_dashed
```

```
$
```

```
$ date | rosie trace date.us_dashed
```

```
Expression: {month "-" day "-" short_long_year}
```

```
Looking at: 《Mon Jul 30 12:43:09 EDT 2018》 (input pos = 1)
```

```
No match
```

```
└─ Expression: month
```

```
Looking at: 《Mon Jul 30 12:43:09 EDT 2018》 (input pos = 1)
```

```
No match
```

```
└─└─ Expression: {"1" [0-2]} / {"0"}? [1-9]}
```

```
Looking at: 《Mon Jul 30 12:43:09 EDT 2018》 (input pos = 1)
```

```
No match
```

```
└─└─└─ Expression: "1" [0-2]}
```

```
Looking at: 《Mon Jul 30 12:43:09 EDT 2018》 (input pos = 1)
```

```
No match
```

```
└─└─└─└─ Expression: "1"
```

```
Looking at: 《Mon Jul 30 12:43:09 EDT 2018》 (input pos = 1)
```

```
No match
```

```
└─└─└─└─└─ Expression: [0-2]
```

```
Not attempted
```

```
└─└─└─└─ Expression: {"0"}? [1-9]}
```

```
Looking at: 《Mon Jul 30 12:43:09 EDT 2018》 (input pos = 1)
```

```
No match
```

```
└─└─ Expression: "-"
```

```
Not attempted
```

```
└─└─ Expression: day
```

```
Not attempted
```

```
└─└─ Expression: "-"
```

```
Not attempted
```

```
└─└─ Expression: short_long_year
```

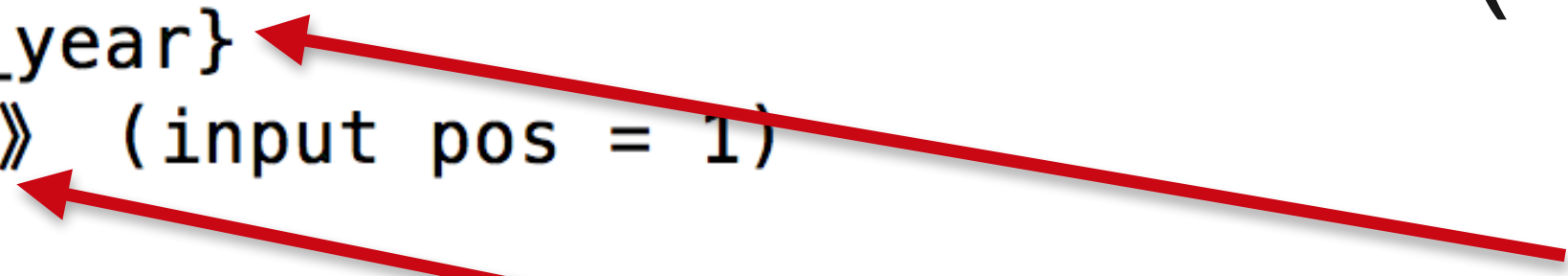
```
Not attempted
```

Trace a (mis-)match

Pattern definition

Input text

Matching steps



Read-eval-print loop

```
$ rosie repl
Rosie 1.0.0-sepcomp3
Rosie> import destructure as des
Rosie> .list des.*
```

Name	Cap?	Type	Color	Source
[snip]				
numalpha	Yes	pattern	default;bold	destructure
parentheses	Yes	pattern	default;bold	destructure
rest	Yes	pattern	default;bold	destructure
semicolons	Yes	pattern	default;bold	destructure
sep		pattern	default;bold	destructure
slashes	Yes	pattern	default;bold	destructure
term	Yes	pattern	default;bold	destructure
tryall		pattern	default;bold	destructure
~		pattern	default;bold	builtin/prelude

24/24 names shown

```
Rosie>
```



```
Rosie> .match des.tryall "(1.2; 3.77; 0)"
```

```
{ "data": "(1.2; 3.77; 0)",  
  "e": 15,  
  "s": 1,  
  "subs":  
    [ { "data": "(1.2; 3.77; 0)",  
        "e": 15,  
        "s": 1,  
        "subs":  
          [ { "data": "1.2; 3.77; 0",  
              "e": 14,  
              "s": 2,  
              "subs":  
                [ { "data": "1.2",  
                    "e": 5,  
                    "s": 2,  
                    "type": "des.find.<search>" },  
                  { "data": " 3.77",  
                    "e": 11,  
                    "s": 6,  
                    "type": "des.find.<search>" },  
                  { "data": " 0",  
                    "e": 14,  
                    "s": 12
```

snip

snip

Read-eval-print loop

- ◆ Define patterns
- ◆ Try them
- ◆ Debug (trace) them

```
Rosie> .match des.tryall "(1.2; 3.77; 0)"
```

```
{"data": "(1.2; 3.77; 0)",  
 "e": 15,  
 "s": 1,  
 "subs":  
   [{"data": "(1.2; 3.77; 0)",  
    "e": 15,  
    "s": 1,  
    "subs":  
      [{"data": "1.2; 3.77; 0",  
       "e": 14,  
       "s": 2,  
       "subs":  
         [{"data": "1.2",  
          "e": 5,  
          "s": 2,  
          "type": "des.find.<search>"}],  
         [{"data": "3.77",  
          "e": 11,  
          "s": 6,  
          "type": "des.find.<search>"}],  
         [{"data": "0",  
          "e": 14,  
          "s": 12
```

snip

snip

Read-eval-print loop

- ◆ Define patterns
- ◆ Try them
- ◆ Debug (trace) them

Executable unit tests

```
-----  
---- net.rpl      Rosie Pattern Language patterns for hostnames, ip addresses, and such  
-----  
  
package net  
import num  
  
[snip]  
  
ipv4 = ip_address_v4  
-- test ipv4 accepts "0.0.0.0", "1.2.234.123", "999.999.999.999"  
-- test ipv4 rejects "1234.1.2.3", "1.2.3", "111.222.333.", "111.222.333..444"  
  
ipv6 = ipv6_mixed / ip_address_v6  
-- test ipv6 includes ipv4 "::192.9.5.5", "::FFFF:129.144.52.38"  
-- test ipv6 excludes ipv4 "1080::8:800:200C:417A", "2010:836B:4179::836B:4179"
```

Executable unit tests

```
$ rosie test /usr/local/lib/rosie/rpl/*.rpl
/usr/local/lib/rosie/rpl/all.rpl
  all 4 tests passed
/usr/local/lib/rosie/rpl/csv.rpl
  no tests found
/usr/local/lib/rosie/rpl/date.rpl
  all 89 tests passed
/usr/local/lib/rosie/rpl/id.rpl
  all 51 tests passed
/usr/local/lib/rosie/rpl/json.rpl
  all 45 tests passed
/usr/local/lib/rosie/rpl/net.rpl
  all 125 tests passed
/usr/local/lib/rosie/rpl/num.rpl
  all 80 tests passed
/usr/local/lib/rosie/rpl/os.rpl
  no tests found
/usr/local/lib/rosie/rpl/time.rpl
  all 85 tests passed
/usr/local/lib/rosie/rpl/ts.rpl
  all 27 tests passed
/usr/local/lib/rosie/rpl/word.rpl
  all 20 tests passed
```

```
$
```

- Part of the documentation
- Regression when making changes
- Use them in app build/compile stage

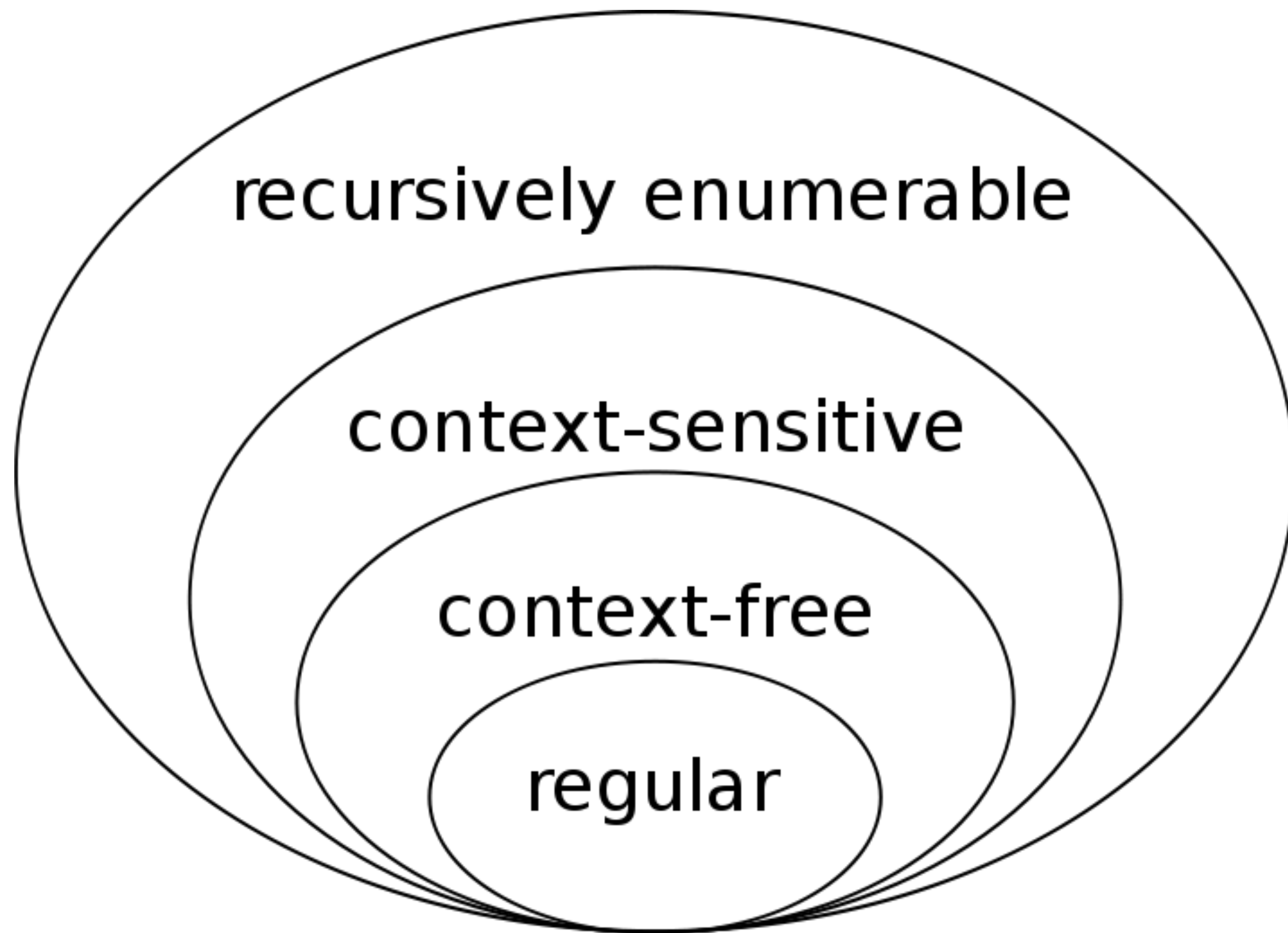
Formal basis

“Language is a process of free creation [though] its laws and principles are fixed”

Noam Chomsky

Formal basis

Chomsky hierarchy



Parsing Expression Grammars: A Recognition-Based Syntactic Foundation

Bryan Ford
Massachusetts Institute of Technology
Cambridge, MA
baford@mit.edu

Abstract

For decades we have been using Chomsky's generative system of grammars, particularly context-free grammars (CFGs) and regular expressions (REs), to express the syntax of programming languages and protocols. The power of generative grammars to express ambiguity is crucial to their original purpose of modelling natural languages, but this very power makes it unnecessarily difficult both to express and to parse machine-oriented languages using CFGs. Parsing Expression Grammars (PEGs) provide an alternative, recognition-based formal foundation for describing machine-oriented syntax, which solves the ambiguity problem by not introducing ambiguity in the first place. Where CFGs express nondeterministic choice between alternatives, PEGs instead use *prioritized choice*. PEGs address frequently felt expressiveness limitations of CFGs and REs, simplifying syntax definitions and making it unnecessary to separate their lexical and hierarchical components. A linear-time parser can be built for any PEG, avoiding both the complexity and fickleness of LR parsers and the inefficiency of generalized CFG parsing. While PEGs provide a rich set of operators for constructing grammars, they are reducible to two minimal recognition schemas developed around 1970, TS/TDPL and gTS/GTDPL, which are here proven equivalent in effective recognition power.

1 Introduction

Most language syntax theory and practice is based on *generative* systems, such as regular expressions and context-free grammars, in which a language is defined formally by a set of rules applied recursively to generate strings of the language. A *recognition-based* system, in contrast, defines a language in terms of rules or predicates that decide whether or not a given string is in the language. Simple languages can be expressed easily in either paradigm. For example, $\{s \in a^* \mid s = (aa)^n\}$ is a generative definition of a trivial language over a unary character set, whose strings are "constructed" by concatenating pairs of a's. In contrast, $\{s \in a^* \mid (|s| \bmod 2 = 0)\}$ is a recognition-based definition of the same language, in which a string of a's is "accepted" if its length is even.

While most language theory adopts the generative paradigm, most practical language applications in computer science involve the recognition and structural decomposition, or *parsing*, of strings. Bridging the gap from generative definitions to practical recognizers is the purpose of our ever-expanding library of parsing algorithms with diverse capabilities and trade-offs [9].

Chomsky's generative system of grammars, from which the ubiqui-

A Text Pattern-Matching Tool based on Parsing Expression Grammars

Roberto Ierusalimschy¹

¹ PUC-Rio, Brazil

This is a preprint of an article accepted for publication in Software: Practice and Experience; Copyright 2008 by John Wiley and Sons.

SUMMARY

Current text pattern-matching tools are based on regular expressions. However, pure regular expressions have proven too weak a formalism for the task: many interesting patterns either are difficult to describe or cannot be described by regular expressions. Moreover, the inherent non-determinism of regular expressions does not fit the need to capture specific parts of a match.

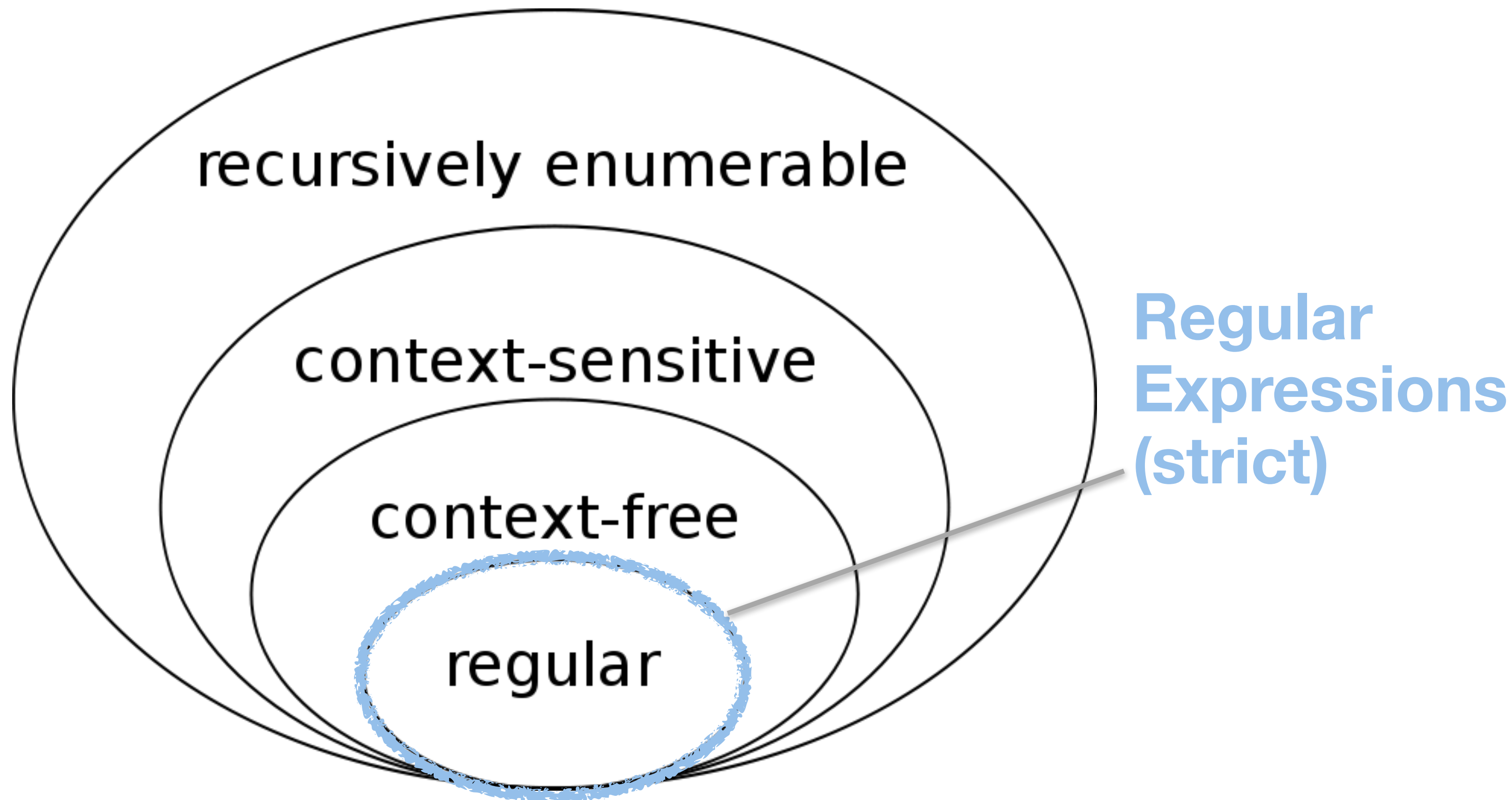
Motivated by these reasons, most scripting languages nowadays use pattern-matching tools that extend the original regular-expression formalism with a set of ad-hoc features, such as greedy repetitions, lazy repetitions, possessive repetitions, "longest match rule", lookahead, etc. These ad-hoc extensions bring their own set of problems, such as lack of a formal foundation and complex implementations.

In this paper, we propose the use of Parsing Expression Grammars (PEGs) as a basis for pattern matching. Following this proposal, we present LPEG, a pattern-matching tool based on PEGs for the Lua scripting language. LPEG unifies the ease of use of pattern-matching tools with the full expressive power of PEGs. Because of this expressive power, it can avoid the myriad of ad-hoc constructions present in several current pattern-matching tools. We also present a Parsing Machine that allows a small and efficient implementation of PEGs for pattern matching.

KEY WORDS: pattern matching, Parsing Expression Grammars, scripting languages

Formal basis

Chomsky hierarchy



Parsing Expression Grammars: A Recognition-Based Syntactic Foundation

Bryan Ford
Massachusetts Institute of Technology
Cambridge, MA
baford@mit.edu

Abstract

For decades we have been using Chomsky's generative system of grammars, particularly context-free grammars (CFGs) and regular expressions (REs), to express the syntax of programming languages and protocols. The power of generative grammars to express ambiguity is crucial to their original purpose of modelling natural languages, but this very power makes it unnecessarily difficult both to express and to parse machine-oriented languages using CFGs. Parsing Expression Grammars (PEGs) provide an alternative, recognition-based formal foundation for describing machine-oriented syntax, which solves the ambiguity problem by not introducing ambiguity in the first place. Where CFGs express nondeterministic choice between alternatives, PEGs instead use *prioritized choice*. PEGs address frequently felt expressiveness limitations of CFGs and REs, simplifying syntax definitions and making it unnecessary to separate their lexical and hierarchical components. A linear-time parser can be built for any PEG, avoiding both the complexity and fickleness of LR parsers and the inefficiency of generalized CFG parsing. While PEGs provide a rich set of operators for constructing grammars, they are reducible to two minimal recognition schemas developed around 1970, TS/TDPL and gTS/GTDPL, which are here proven equivalent in effective recognition power.

1 Introduction

Most language syntax theory and practice is based on *generative* systems, such as regular expressions and context-free grammars, in which a language is defined formally by a set of rules applied recursively to generate strings of the language. A *recognition-based* system, in contrast, defines a language in terms of rules or predicates that decide whether or not a given string is in the language. Simple languages can be expressed easily in either paradigm. For example, $\{s \in a^* \mid s = (aa)^n\}$ is a generative definition of a trivial language over a unary character set, whose strings are "constructed" by concatenating pairs of a's. In contrast, $\{s \in a^* \mid (|s| \bmod 2 = 0)\}$ is a recognition-based definition of the same language, in which a string of a's is "accepted" if its length is even.

While most language theory adopts the generative paradigm, most practical language applications in computer science involve the recognition and structural decomposition, or *parsing*, of strings. Bridging the gap from generative definitions to practical recognizers is the purpose of our ever-expanding library of parsing algorithms with diverse capabilities and trade-offs [9].

Chomsky's generative system of grammars, from which the ubiqui-



A Text Pattern-Matching Tool based on Parsing Expression Grammars

Roberto Ierusalimsky¹

¹ PUC-Rio, Brazil

This is a preprint of an article accepted for publication in Software: Practice and Experience; Copyright 2008 by John Wiley and Sons.

SUMMARY

Current text pattern-matching tools are based on regular expressions. However, pure regular expressions have proven too weak a formalism for the task: many interesting patterns either are difficult to describe or cannot be described by regular expressions. Moreover, the inherent non-determinism of regular expressions does not fit the need to capture specific parts of a match.

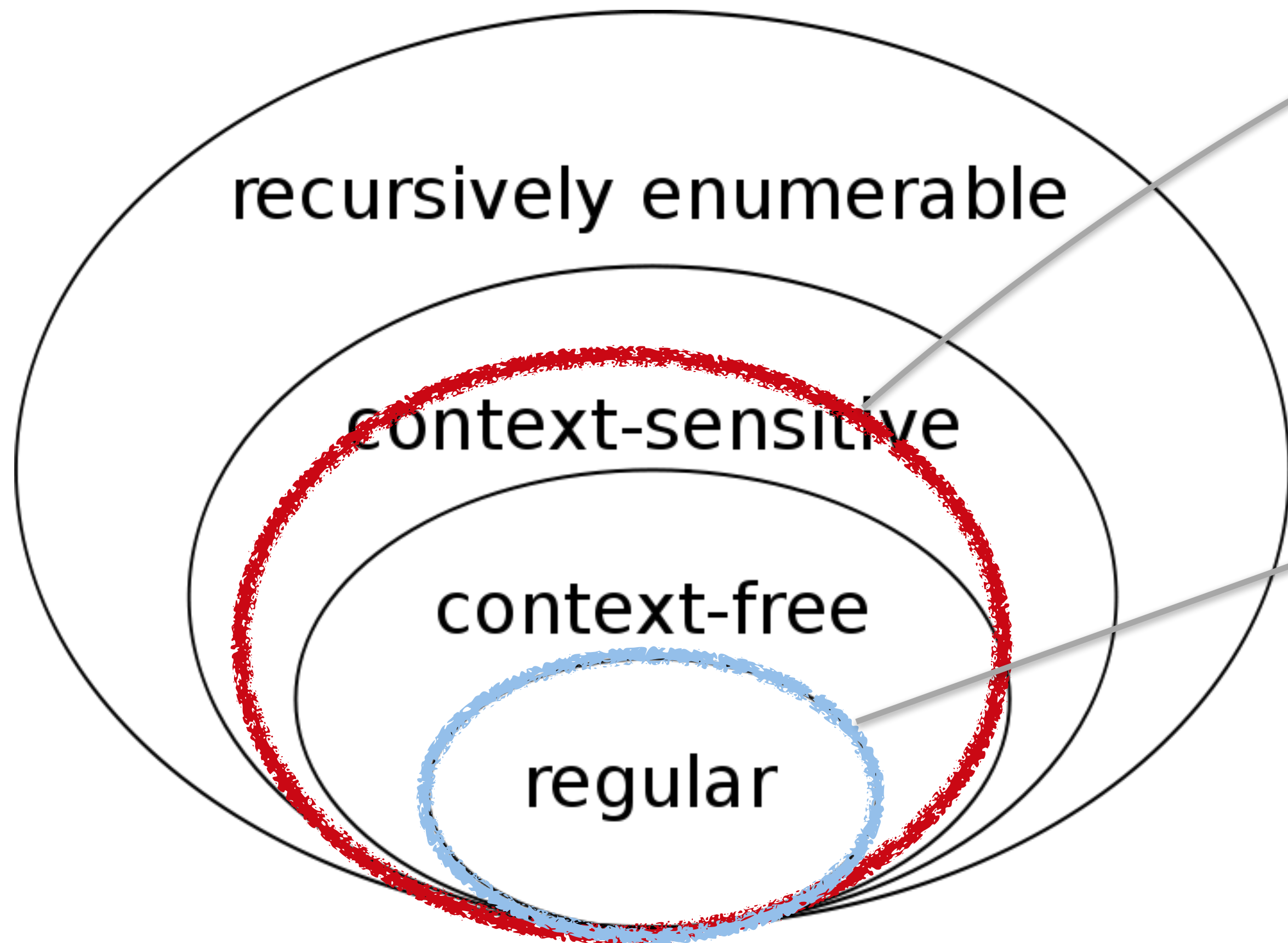
Motivated by these reasons, most scripting languages nowadays use pattern-matching tools that extend the original regular-expression formalism with a set of ad-hoc features, such as greedy repetitions, lazy repetitions, possessive repetitions, "longest match rule", lookahead, etc. These ad-hoc extensions bring their own set of problems, such as lack of a formal foundation and complex implementations.

In this paper, we propose the use of Parsing Expression Grammars (PEGs) as a basis for pattern matching. Following this proposal, we present LPEG, a pattern-matching tool based on PEGs for the Lua scripting language. LPEG unifies the ease of use of pattern-matching tools with the full expressive power of PEGs. Because of this expressive power, it can avoid the myriad of ad-hoc constructions present in several current pattern-matching tools. We also present a Parsing Machine that allows a small and efficient implementation of PEGs for pattern matching.

KEY WORDS: pattern matching, Parsing Expression Grammars, scripting languages

Formal basis

Chomsky hierarchy



**Rosie
Pattern
Language**
(and all PEG grammars)

**Regular
Expressions**
(strict)

Parsing Expression Grammars: A Recognition-Based Syntactic Foundation

Bryan Ford
Massachusetts Institute of Technology
Cambridge, MA
baford@mit.edu

Abstract

For decades we have been using Chomsky's generative system of grammars, particularly context-free grammars (CFGs) and regular expressions (REs), to express the syntax of programming languages and protocols. The power of generative grammars to express ambiguity is crucial to their original purpose of modelling natural languages, but this very power makes it unnecessarily difficult both to express and to parse machine-oriented languages using CFGs. Parsing Expression Grammars (PEGs) provide an alternative, recognition-based formal foundation for describing machine-oriented syntax, which solves the ambiguity problem by not introducing ambiguity in the first place. Where CFGs express nondeterministic choice between alternatives, PEGs instead use *prioritized choice*. PEGs address frequently felt expressiveness limitations of CFGs and REs, simplifying syntax definitions and making it unnecessary to separate their lexical and hierarchical components. A linear-time parser can be built for any PEG, avoiding both the complexity and fickleness of LR parsers and the inefficiency of generalized CFG parsing. While PEGs provide a rich set of operators for constructing grammars, they are reducible to two minimal recognition schemas developed around 1970, TS/TDPL and gTS/GTDPL, which are here proven equivalent in effective recognition power.

1 Introduction

Most language syntax theory and practice is based on *generative* systems, such as regular expressions and context-free grammars, in which a language is defined formally by a set of rules applied recursively to generate strings of the language. A *recognition-based* system, in contrast, defines a language in terms of rules or predicates that decide whether or not a given string is in the language. Simple languages can be expressed easily in either paradigm. For example, $\{s \in a^* \mid s = (aa)^n\}$ is a generative definition of a trivial language over a unary character set, whose strings are "constructed" by concatenating pairs of a's. In contrast, $\{s \in a^* \mid (|s| \bmod 2 = 0)\}$ is a recognition-based definition of the same language, in which a string of a's is "accepted" if its length is even.

While most language theory adopts the generative paradigm, most practical language applications in computer science involve the recognition and structural decomposition, or *parsing*, of strings. Bridging the gap from generative definitions to practical recognizers is the purpose of our ever-expanding library of parsing algorithms with diverse capabilities and trade-offs [9].

Chomsky's generative system of grammars, from which the ubiquitous



A Text Pattern-Matching Tool based on Parsing Expression Grammars

Roberto Ierusalimsky¹

¹ PUC-Rio, Brazil

This is a preprint of an article accepted for publication in Software: Practice and Experience; Copyright 2008 by John Wiley and Sons.

SUMMARY

Current text pattern-matching tools are based on regular expressions. However, pure regular expressions have proven too weak a formalism for the task: many interesting patterns either are difficult to describe or cannot be described by regular expressions. Moreover, the inherent non-determinism of regular expressions does not fit the need to capture specific parts of a match.

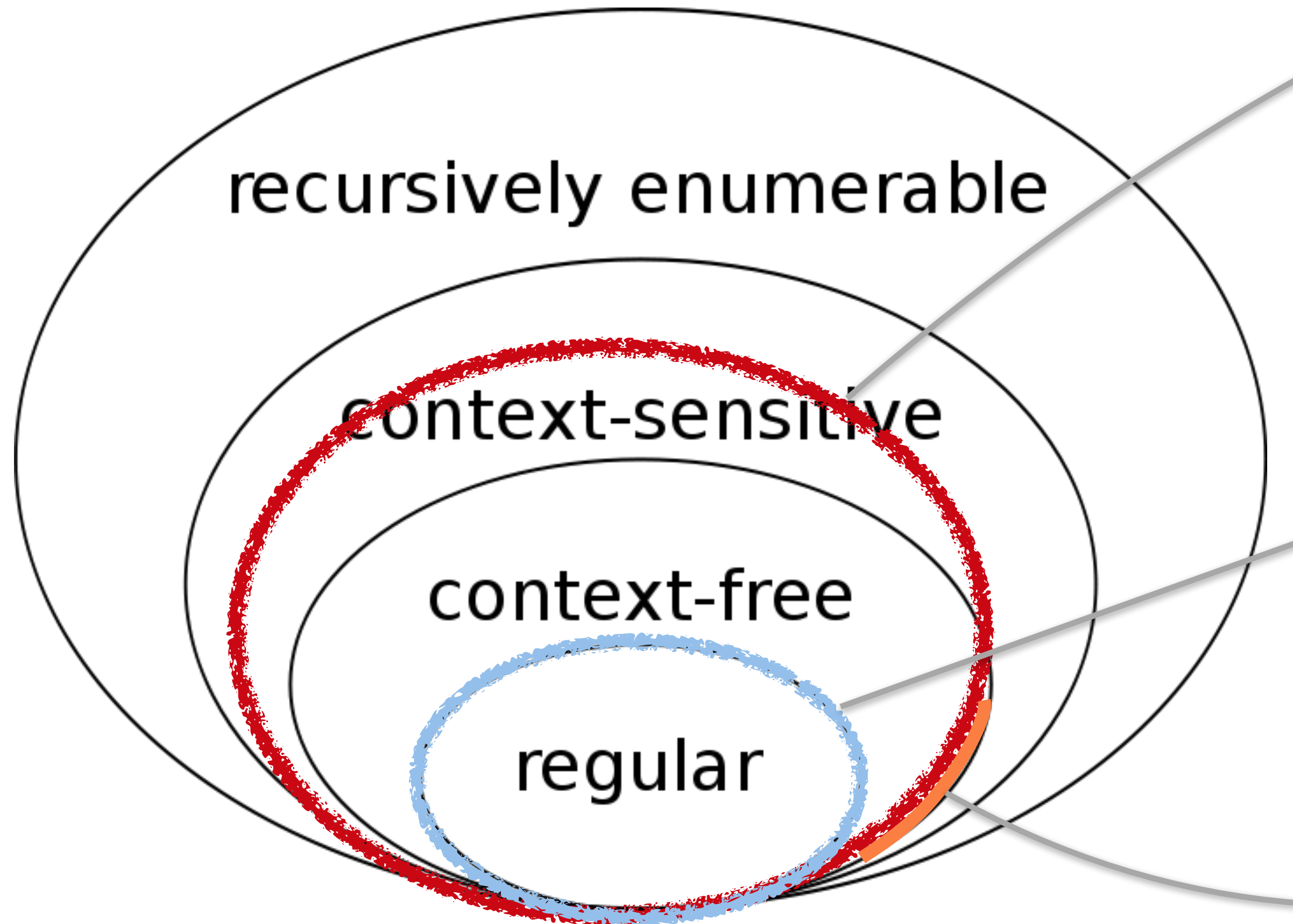
Motivated by these reasons, most scripting languages nowadays use pattern-matching tools that extend the original regular-expression formalism with a set of ad-hoc features, such as greedy repetitions, lazy repetitions, possessive repetitions, "longest match rule", lookahead, etc. These ad-hoc extensions bring their own set of problems, such as lack of a formal foundation and complex implementations.

In this paper, we propose the use of Parsing Expression Grammars (PEGs) as a basis for pattern matching. Following this proposal, we present LPEG, a pattern-matching tool based on PEGs for the Lua scripting language. LPEG unifies the ease of use of pattern-matching tools with the full expressive power of PEGs. Because of this expressive power, it can avoid the myriad of ad-hoc constructions present in several current pattern-matching tools. We also present a Parsing Machine that allows a small and efficient implementation of PEGs for pattern matching.

KEY WORDS: pattern matching, Parsing Expression Grammars, scripting languages

Formal basis

Chomsky hierarchy



**Rosie
Pattern
Language**
(and all PEG grammars)

**Regular
Expressions**
(strict)

**Open
Question:
PEG > CFG**

Parsing Expression Grammars: A Recognition-Based Syntactic Foundation

Bryan Ford
Massachusetts Institute of Technology
Cambridge, MA
baford@mit.edu

Abstract

For decades we have been using Chomsky's generative system of grammars, particularly context-free grammars (CFGs) and regular expressions (REs), to express the syntax of programming languages and protocols. The power of generative grammars to express ambiguity is crucial to their original purpose of modelling natural languages, but this very power makes it unnecessarily difficult both to express and to parse machine-oriented languages using CFGs. Parsing Expression Grammars (PEGs) provide an alternative, recognition-based formal foundation for describing machine-oriented syntax, which solves the ambiguity problem by not introducing ambiguity in the first place. Where CFGs express nondeterministic choice between alternatives, PEGs instead use *prioritized choice*. PEGs address frequently felt expressiveness limitations of CFGs and REs, simplifying syntax definitions and making it unnecessary to separate their lexical and hierarchical components. A linear-time parser can be built for any PEG, avoiding both the complexity and fickleness of LR parsers and the inefficiency of generalized CFG parsing. While PEGs provide a rich set of operators for constructing grammars, they are reducible to two minimal recognition schemas developed around 1970, TS/TDPL and gTS/GTDPL, which are here proven equivalent in effective recognition power.

1 Introduction

Most language syntax theory and practice is based on *generative* systems, such as regular expressions and context-free grammars, in which a language is defined formally by a set of rules applied recursively to generate strings of the language. A *recognition-based* system, in contrast, defines a language in terms of rules or predicates that decide whether or not a given string is in the language. Simple languages can be expressed easily in either paradigm. For example, $\{s \in a^* \mid s = (aa)^n\}$ is a generative definition of a trivial language over a unary character set, whose strings are "constructed" by concatenating pairs of a's. In contrast, $\{s \in a^* \mid (|s| \bmod 2 = 0)\}$ is a recognition-based definition of the same language, in which a string of a's is "accepted" if its length is even.

While most language theory adopts the generative paradigm, most practical language applications in computer science involve the recognition and structural decomposition, or *parsing*, of strings. Bridging the gap from generative definitions to practical recognizers is the purpose of our ever-expanding library of parsing algorithms with diverse capabilities and trade-offs [9].

Chomsky's generative system of grammars, from which the ubiquitous



A Text Pattern-Matching Tool based on Parsing Expression Grammars

Roberto Ierusalimsky¹

¹ PUC-Rio, Brazil

This is a preprint of an article accepted for publication in Software: Practice and Experience; Copyright 2008 by John Wiley and Sons.

SUMMARY

Current text pattern-matching tools are based on regular expressions. However, pure regular expressions have proven too weak a formalism for the task: many interesting patterns either are difficult to describe or cannot be described by regular expressions. Moreover, the inherent non-determinism of regular expressions does not fit the need to capture specific parts of a match.

Motivated by these reasons, most scripting languages nowadays use pattern-matching tools that extend the original regular-expression formalism with a set of ad-hoc features, such as greedy repetitions, lazy repetitions, possessive repetitions, "longest match rule", lookahead, etc. These ad-hoc extensions bring their own set of problems, such as lack of a formal foundation and complex implementations.

In this paper, we propose the use of Parsing Expression Grammars (PEGs) as a basis for pattern matching. Following this proposal, we present LPEG, a pattern-matching tool based on PEGs for the Lua scripting language. LPEG unifies the ease of use of pattern-matching tools with the full expressive power of PEGs. Because of this expressive power, it can avoid the myriad of ad-hoc constructions present in several current pattern-matching tools. We also present a Parsing Machine that allows a small and efficient implementation of PEGs for pattern matching.

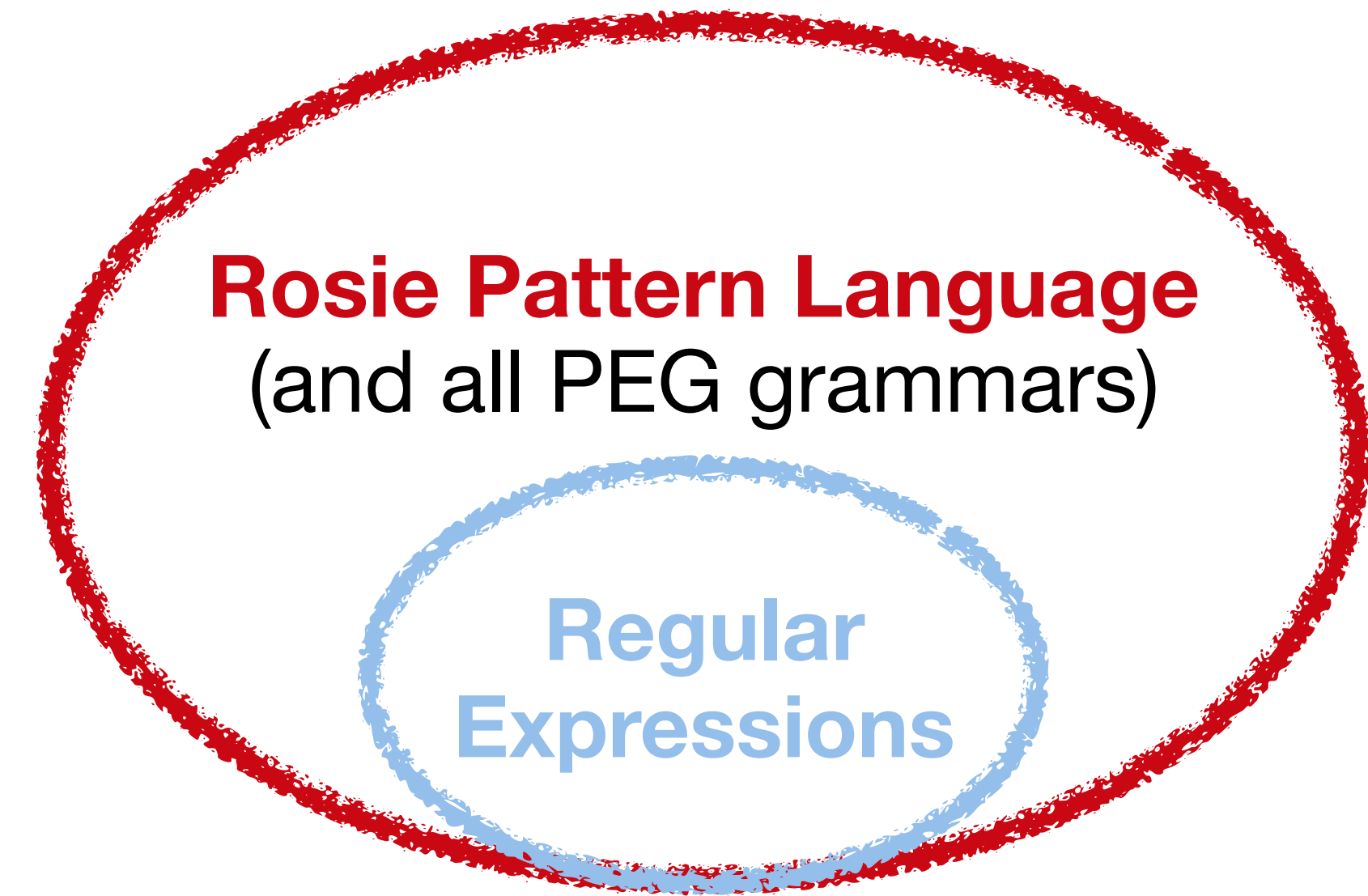
KEY WORDS: pattern matching, Parsing Expression Grammars, scripting languages

Comparison to regex: RPL is based on a different formalism

Comparison to regex: RPL is based on a different formalism

Parsing Expression Grammars

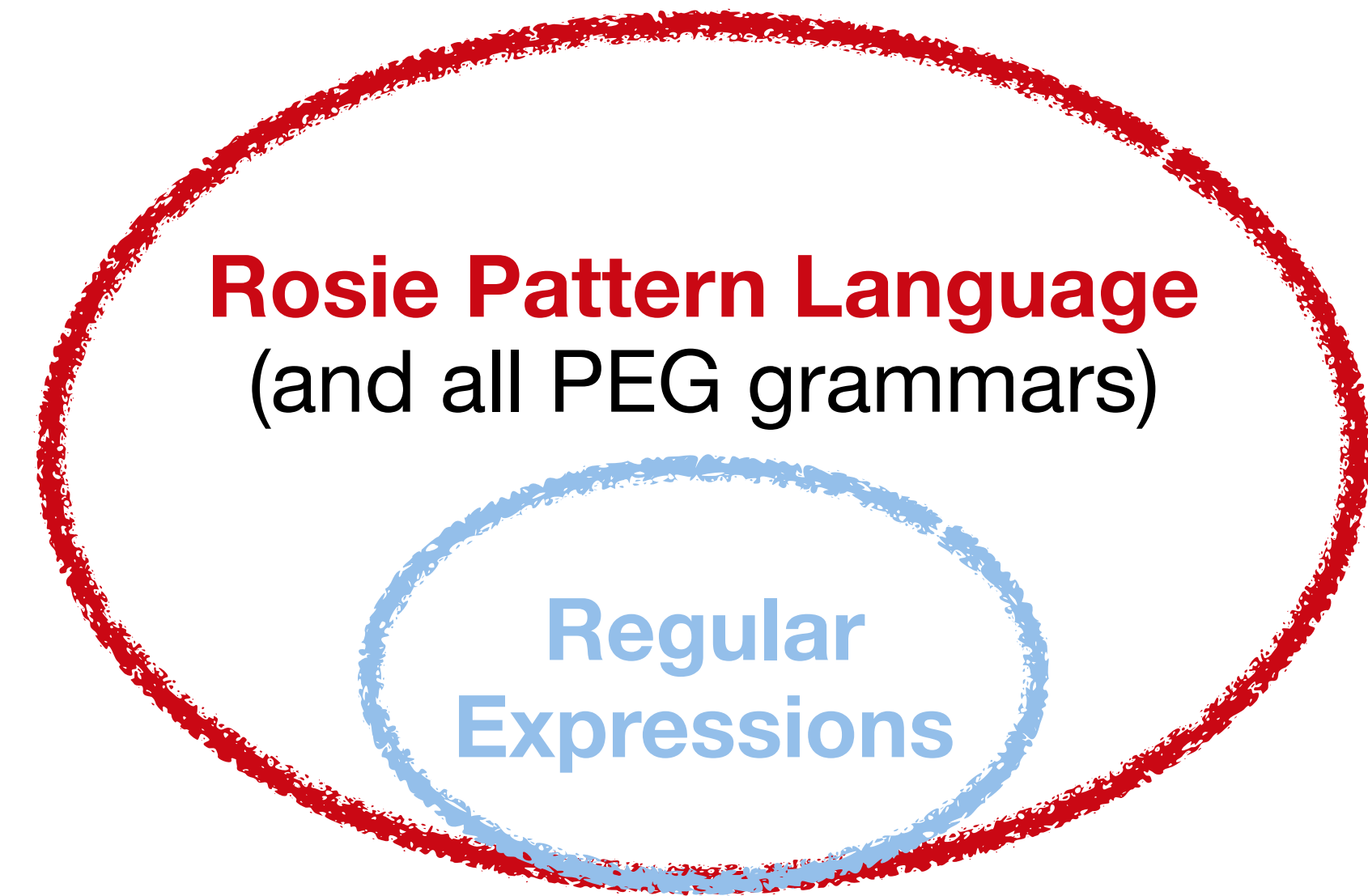
- Strictly more powerful than regular expressions



Comparison to regex: RPL is based on a different formalism

Parsing Expression Grammars

- Strictly more powerful than regular expressions
- Supports recursive pattern definitions



```
grammar
  bal = { "(" bal? ")" }+
end
```


Comparison to regex: RPL is based on a different formalism

Parsing Expression Grammars

- Strictly more powerful than regular expressions
- Supports recursive pattern definitions



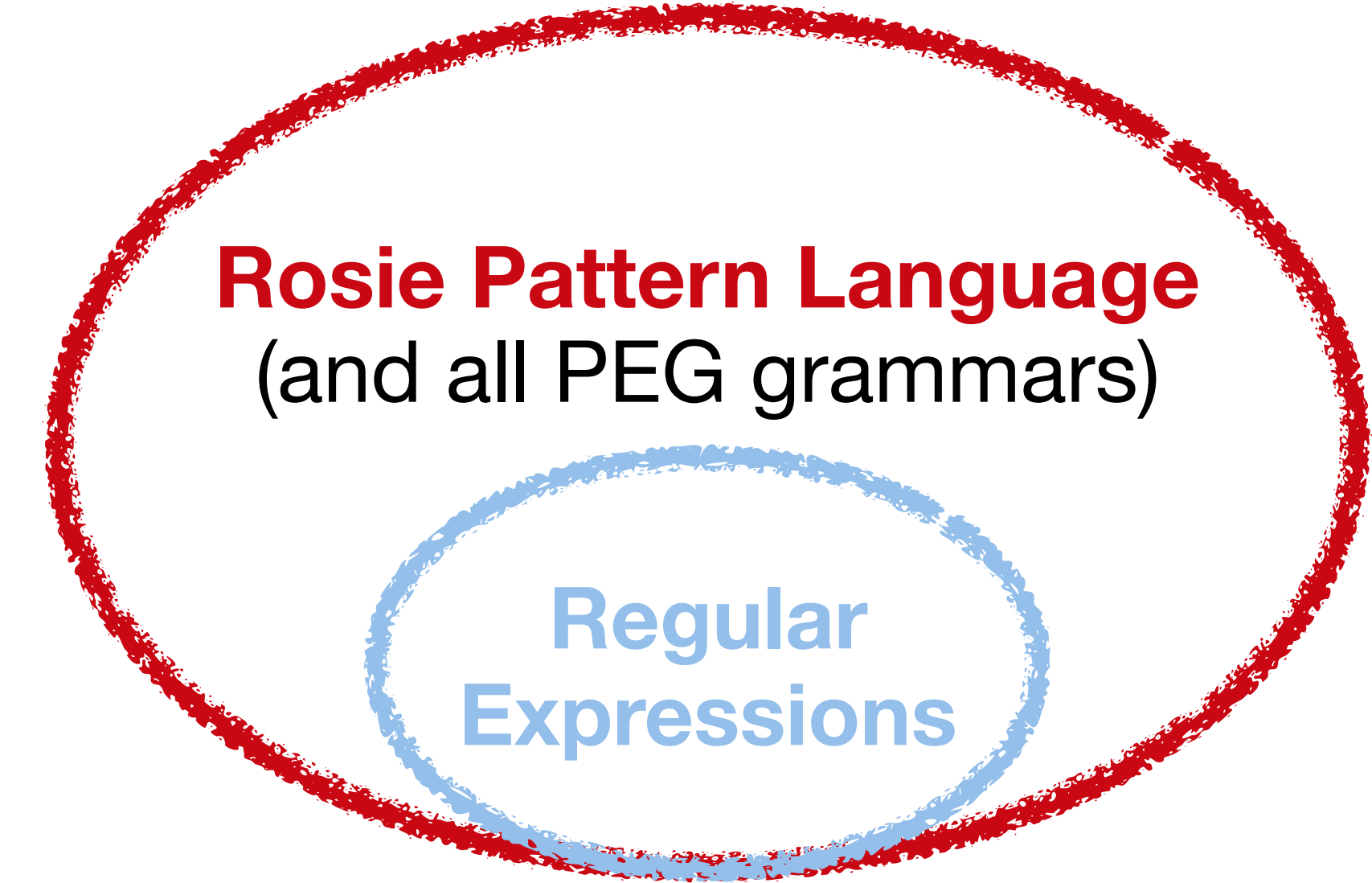
```
grammar  
  bal = { "(" bal? " " }+  
end
```

Perl: `(^\((?[-1])?\))+`

Comparison to regex: RPL is based on a different formalism

Parsing Expression Grammars

- Strictly more powerful than regular expressions
- Supports recursive pattern definitions
- *Packrat* implementation guarantees linear time



Comparison to regex: RPL is based on a different formalism

Parsing Expression Grammars

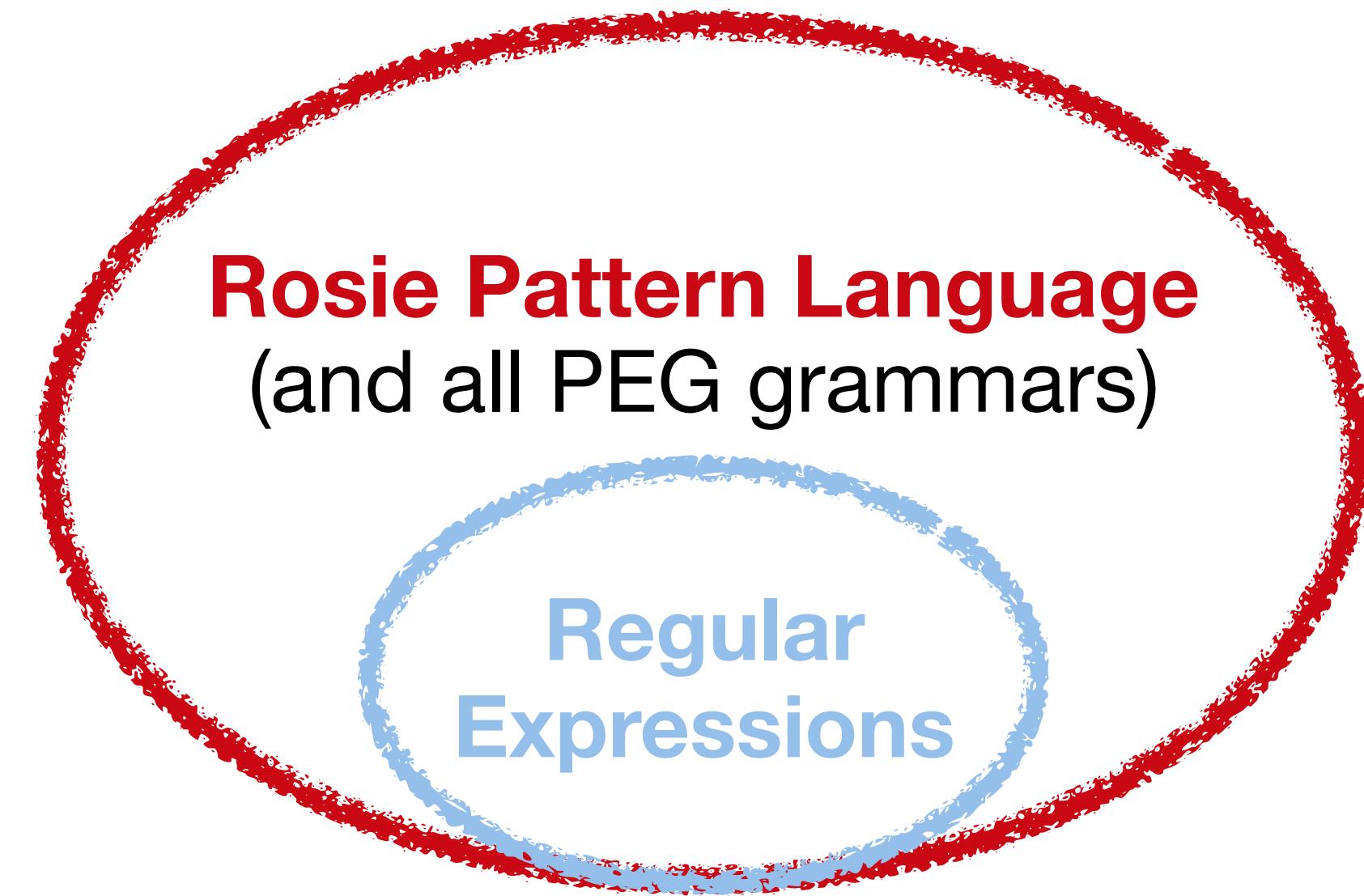
- Strictly more powerful than regular expressions
- Supports recursive pattern definitions
- *Packrat* implementation guarantees linear time
- Rosie uses a ***Matching VM*** implementation
 - Uses less space
 - Linear time for non-grammar, non-lookaround



Comparison to regex: RPL is based on a different formalism

Parsing Expression Grammars

- Strictly more powerful than regular expressions
- Supports recursive pattern definitions
- *Packrat* implementation guarantees linear time
- Rosie uses a ***Matching VM*** implementation
 - Uses less space
 - Linear time for non-grammar, non-lookaround
- Expressions are greedy and possessive



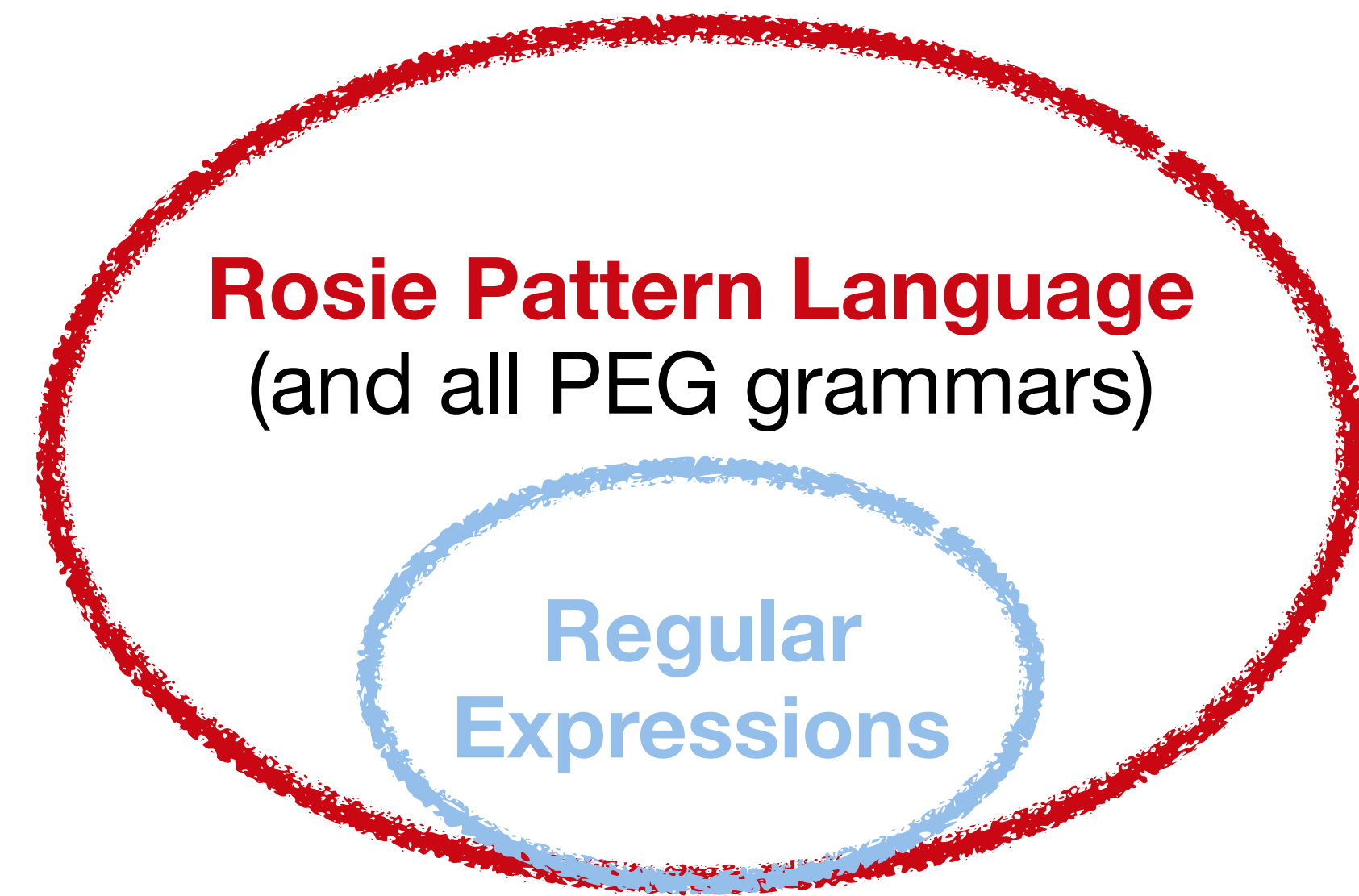
`.* "x"` *always fails!*

`{!"x" .}* "x"`

Comparison to regex: RPL is based on a different formalism

Parsing Expression Grammars

- Strictly more powerful than regular expressions
- Supports recursive pattern definitions
- *Packrat* implementation guarantees linear time
- Rosie uses a *Matching VM* implementation
 - Uses less space
 - Linear time for non-grammar, non-lookaround
- Expressions are greedy and possessive



`.* "x"` *always fails!*

`{!"x" .}* "x"`

`find:"x"`

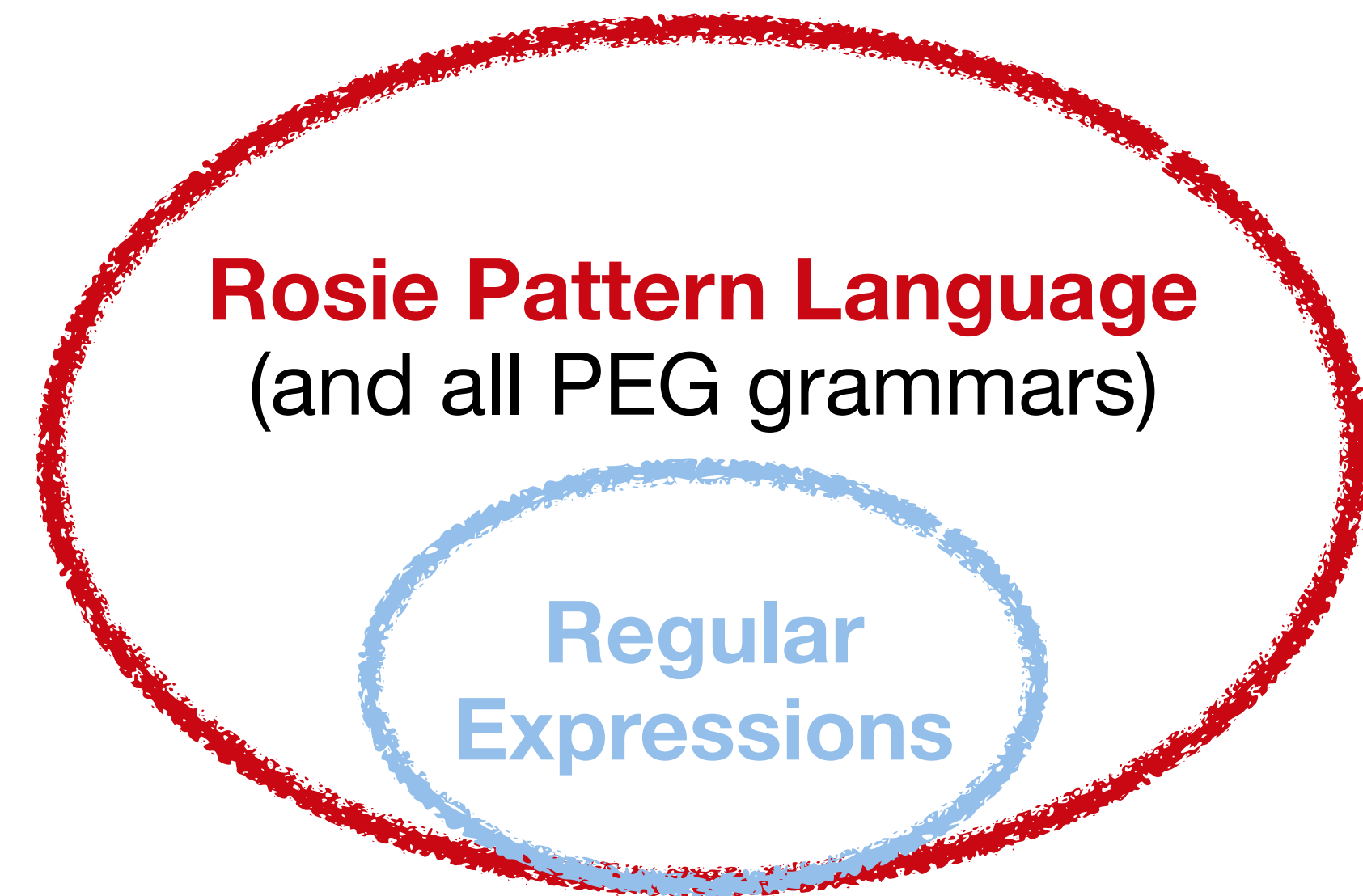
Comparison to regex: RPL is based on a different formalism

Parsing Expression Grammars

- Strictly more powerful than regular expressions
- Supports recursive pattern definitions
- *Packrat* implementation guarantees linear time
- Rosie uses a *Matching VM* implementation
 - Uses less space
 - Linear time for non-grammar, non-lookaround
- Expressions are greedy and possessive

Automated conversion of regex to RPL

- A practical implementation is underway



`. * "x"` *always fails!*

`{!"x" .} * "x"`

`find: "x"`

A “little language” built like a big language

1. Parser is Rosie itself (self-hosted)
2. Macro expansion phase
 - Today, macros written in Lua
 - Future?
3. Conversion to tree representation
4. Optimizations
 - Inlining (always possible with pure functions!)
 - Common sub-expression elimination
 - Lots of small opportunities also
5. Code generation
 - Further optimizations, e.g. peephole

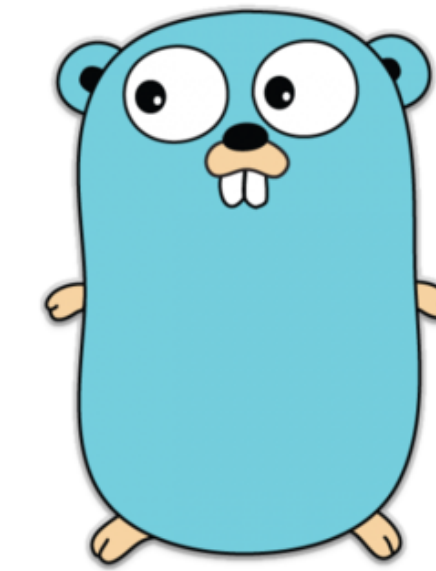


Using Rosie in programs

Today:



Haskell



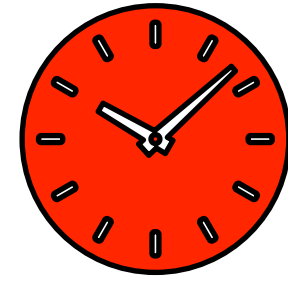
Go

Once and future:

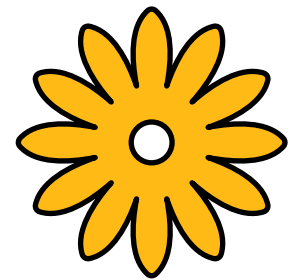


Clojure

Join the Rosie community!

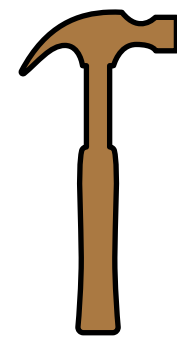


```
make ;  
make install (optional)
```



Contribute Patterns

- Domain-specific
- Authoritative
 - E.g. from RFC
- Non-English patterns!
- “Looks like” (recognizers)
- Byte-encoded data?



Write Tools

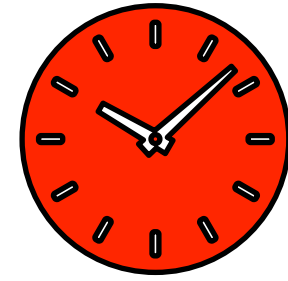
- Package info
- Better trace (compact)
- Linter
- Notebook (Jupyter?)
- Integrations
 - scikit-learn
 - Spark



Implement features

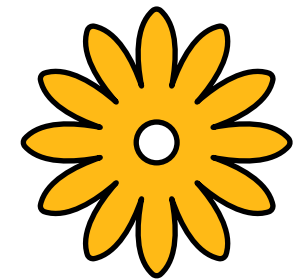
- Optimizations
- Language-specific libs
 - Improve or create
 - Python, R, Go, Java, ...
- User-written extensions
 - Output encoders
 - Macros
 - Character sets

Join the Rosie community!



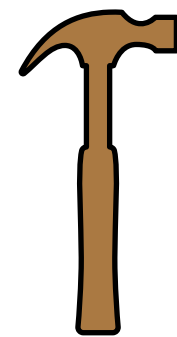
```
make ;  
make install (optional)
```

Or: brew install rosie
Also: pip install rosie



Contribute Patterns

- Domain-specific
- Authoritative
 - E.g. from RFC
- Non-English patterns!
- “Looks like” (recognizers)
- Byte-encoded data?



Write Tools

- Package info
- Better trace (compact)
- Linter
- Notebook (Jupyter?)
- Integrations
 - scikit-learn
 - Spark



Implement features

- Optimizations
- Language-specific libs
 - Improve or create
 - Python, R, Go, Java, ...
- User-written extensions
 - Output encoders
 - Macros
 - Character sets

Thank you!

Rosie Pattern Language

■ Pattern libraries

- Standard library, including full Unicode (UTF-8) support
- Community libraries (e.g. GitHub)
- User libraries

■ Output formats

- Colorized text for humans
- JSON for programs
- Full lines or just matches (like grep)
- And others...

■ Development tools

- Command line interface, read/eval/print loop
- Trace output
- Unit tests (automated)
- Packages (shareable)

■ Built for big data but makes a better grep

- Readable, maintainable
- Works well with git/diff, pipelines (unit tests), dependency mgmt

Rosie Pattern Language



Formal basis:

- ◆ Parser combinators
- ◆ Based on Parsing Exp. Grammars
- ◆ Good performance, often linear
- ◆ Not a “packrat” implementation

Additional slides follow...



Summary

Faster

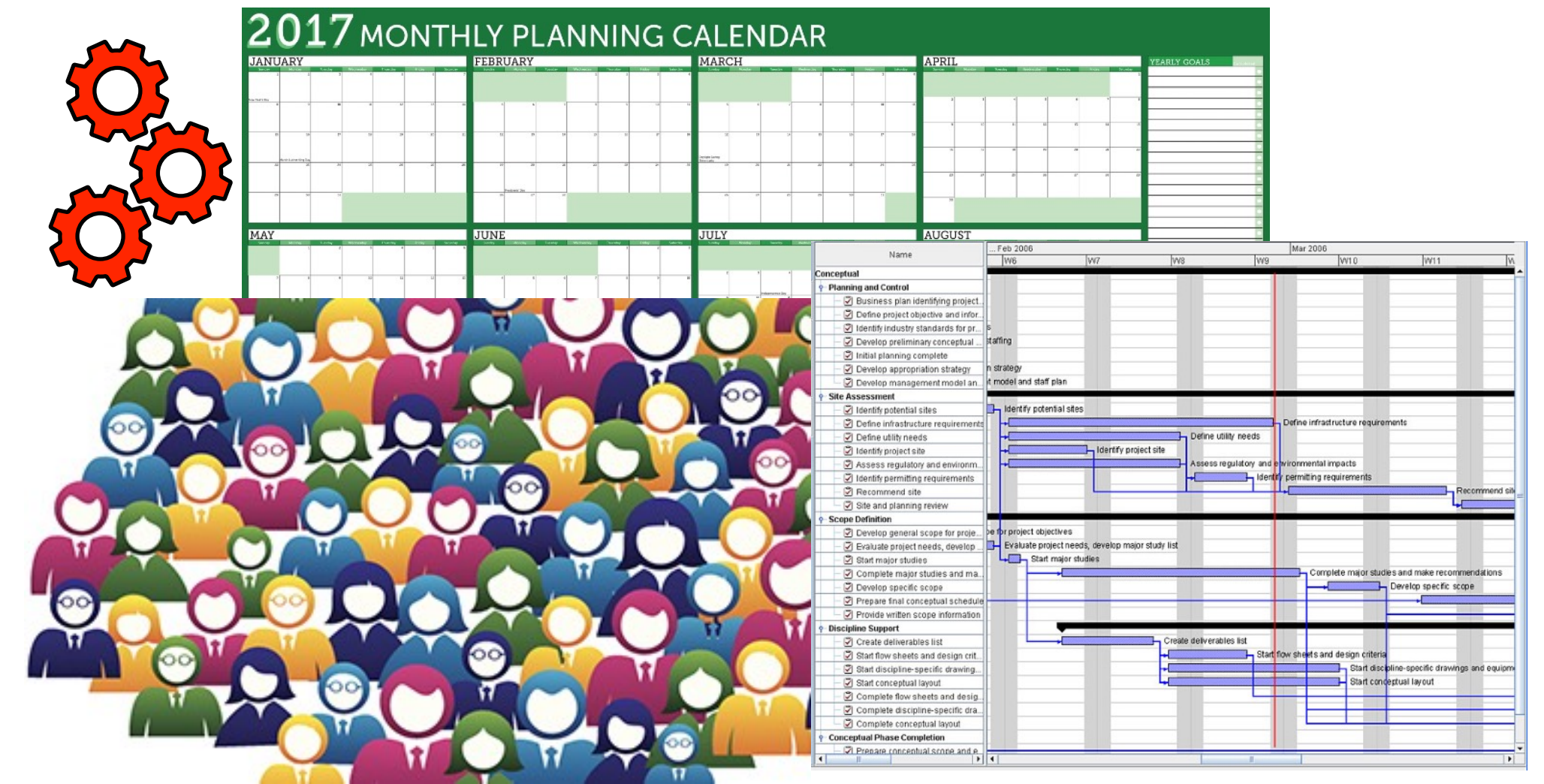
- ◆ Dev time:
 - ✓ library of patterns you don't have to write
 - ✓ new patterns composed of existing patterns
- ◆ Run time: matching performance very good

Better

- ◆ shareable libraries
- ◆ conformance to RFCs
- ◆ readable syntax, and strict semantics (and no flags)
- ◆ plays well with DevOps tools (git/diff, package management, unit tests)

Cheaper

- ◆ ROI in reduced development and maintenance costs
- ◆ And, it's free open source software (MIT license)



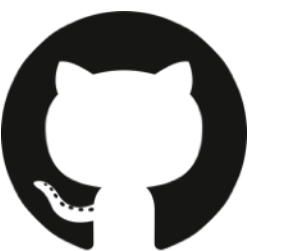
Jenkins



git



Travis CI



1. Mining source code repositories

- “Micro-grammar” approach:

How to build static checking systems using orders of magnitude less code by Brown, Nötzli, Engler

- NCSU students: 6 features x 10 languages

Features →	Comments	Dependencies	Class / Struct Defs	Function Defs	Error Handling	String Literals	Function Bodies	Class Bodies
Languages ↓								
Java	✓	✓	✓	✓	✓	✓		
C	✓	✓	✓	✓	✓	✓		
C++	✓	✓	✓	✓	✓	✓		
C#	✓	✓	✓	✓	✓	✓		
Python	✓	✓	✓	✓	✓	✓		
JScript	✓	✓	✓	✓	✓	✓		
Ruby	✓	✓	✓	✓	✓	✓		
R	✓	✓	✗	✓	✓	✓		
Go	✓	✓	✓	✓	✓	✓		
Bash	✓	✗	✗	✓	✓	✓		
VB	✓	✓	✓	✓	✓	✓		

2. Application log processing (streaming or batch)

```
$ tail -n 3 /var/log/system.log | rosie match all.things
Jul 29 09:48:58 Jamies-Compabler com.apple.xpc.launchd[1] (com.apple.quicklook): Service only ran for 0
Jul 29 09:48:59 Jamies-Compabler com.apple.xpc.launchd[1] (com.apple.quicklook[91387]): Endpoint has be
Jul 29 09:48:59 Jamies-Compabler kcm[91389]: DEPRECATED USE in libdispatch client: Setting timer interv
```


1. Mining source code repositories

- “Micro-grammar” approach:

How to build static checking systems using orders of magnitude less code by Brown, Nötzli, Engler

- NCSU students: 6 features x 10 languages

Features → Languages ↓	Comments	Dependencies	Class / Struct Defs	Function Defs	Error Handling	String Literals	Function Bodies	Class Bodies
Java	✓	✓	✓	✓	✓	✓		
C	✓	✓	✓	✓	✓	✓		
C++	✓	✓	✓	✓	✓	✓		
C#	✓	✓	✓	✓	✓	✓		
Python	✓	✓	✓	✓	✓	✓		
JScript	✓	✓	✓	✓	✓	✓		
Ruby	✓	✓	✓	✓	✓	✓		
R	✓	✓	✗	✓	✓	✓		
Go	✓	✓	✓	✓	✓	✓		
Bash	✓	✗	✗	✓	✓	✓		
VB	✓	✓	✓	✓	✓	✓		

2. Application log processing (streaming or batch)

```
$ tail -n 3 /var/log/system.log | rosie match all.things
Jul 29 09:48:58 Jamies-Compabler com.apple.xpc.launchd[1] (com.apple.quicklook): Service only ran for 0
Jul 29 09:48:59 Jamies-Compabler com.apple.xpc.launchd[1] (com.apple.quicklook[91387]): Endpoint has be
Jul 29 09:48:59 Jamies-Compabler kcm[91389]: DEPRECATED USE in libdispatch client: Setting timer interv
```

3. Secure engineering principle: *Parse everything!*

The most critical risk in every OWASP report since 2003: **Injection attacks (unvalidated input)**

Best practice: Whitelist valid input, which requires parsing every input

RPL: Familiar concepts (and syntax)

RPL expression	Matches
<code>pat*</code>	Zero or more copies of <code>pat</code>
<code>pat+</code>	One or more copies of <code>pat</code>
<code>pat?</code>	Zero or one copies of <code>pat</code>
<code>pat{n}</code>	Exactly <code>n</code> copies of <code>pat</code>

RPL expression	Meaning
<code>[:name:]</code>	Named character set (see <i>note [a]</i>)
<code>[:^name:]</code>	Complement of a named character set
<code>[x-y]</code>	Range of characters, from <code>x</code> to <code>y</code> (see <i>note [b]</i>)
<code>[^x-y]</code>	Complement of a character range
<code>[...]</code>	List of characters (in place of <code>...</code>)
<code>[^...]</code>	Complement of the character list <code>...</code>
<code>[cs1 cs2 ...]</code>	Union of character sets <code>cs1</code> , <code>cs2</code> , etc. (E.g. <code>[[a-f][0-9]]</code>)
<code>[^ cs1 cs2 ...]</code>	Complement of a union of character sets

RPL expression	Meaning
<code>> pat</code>	Look ahead at <code>pat</code> (predicate: consumes no input)
<code>< pat</code>	Look behind at <code>pat</code> (predicate: consumes no input)
<code>!pat</code>	Not <code>pat</code> , i.e. not looking at <code>pat</code> . Same as <code>!>pat</code> .

RPL expression	Meaning
<code>p / q</code>	Ordered choice: match <code>p</code> , and <code>p</code> fails, match <code>q</code>

RPL: Familiar concepts (and syntax)

RPL expression	Matches
<code>pat*</code>	Zero or more copies of <code>pat</code>
<code>pat+</code>	One or more copies of <code>pat</code>
<code>pat?</code>	Zero or one copies of <code>pat</code>
<code>pat{n}</code>	Exactly <code>n</code> copies of <code>pat</code>

RPL expression	Meaning
<code>[:name:]</code>	Named character set (see <i>note [a]</i>)
<code>[:^name:]</code>	Complement of a named character set
<code>[x-y]</code>	Range of characters, from <code>x</code> to <code>y</code> (see <i>note [b]</i>)
<code>[^x-y]</code>	Complement of a character range
<code>[...]</code>	List of characters (in place of <code>...</code>)
<code>[^...]</code>	Complement of the character list <code>...</code>
<code>[cs1 cs2 ...]</code>	Union of character sets <code>cs1</code> , <code>cs2</code> , etc. (E.g. <code>[[a-f][0-9]]</code>)
<code>[^ cs1 cs2 ...]</code>	Complement of a union of character sets

RPL expression	Meaning
<code>> pat</code>	Look ahead at <code>pat</code> (predicate: consumes no input)
<code>< pat</code>	Look behind at <code>pat</code> (predicate: consumes no input)
<code>!pat</code>	Not <code>pat</code> , i.e. not looking at <code>pat</code> . Same as <code>!>pat</code> .

RPL expression	Meaning
<code>p / q</code>	Ordered choice: match <code>p</code> , and <code>p</code> fails, match <code>q</code>

Differences

- Always greedy
- Always possessive
- Choices are ordered

RPL: Familiar concepts (and syntax)

RPL expression	Matches
<code>pat*</code>	Zero or more copies of <code>pat</code>
<code>pat+</code>	One or more copies of <code>pat</code>
<code>pat?</code>	Zero or one copies of <code>pat</code>
<code>pat{n}</code>	Exactly <code>n</code> copies of <code>pat</code>

RPL expression	Meaning
<code>[:name:]</code>	Named character set (see <i>note [a]</i>)
<code>[:^name:]</code>	Complement of a named character set
<code>[x-y]</code>	Range of characters, from <code>x</code> to <code>y</code> (see <i>note [b]</i>)
<code>[^x-y]</code>	Complement of a character range
<code>[...]</code>	List of characters (in place of <code>...</code>)
<code>[^...]</code>	Complement of the character list <code>...</code>
<code>[cs1 cs2 ...]</code>	Union of character sets <code>cs1</code> , <code>cs2</code> , etc. (E.g. <code>[[a-f][0-9]]</code>)
<code>[^ cs1 cs2 ...]</code>	Complement of a union of character sets

RPL expression	Meaning
<code>> pat</code>	Look ahead at <code>pat</code> (predicate: consumes no input)
<code>< pat</code>	Look behind at <code>pat</code> (predicate: consumes no input)
<code>!pat</code>	Not <code>pat</code> , i.e. not looking at <code>pat</code> . Same as <code>!>pat</code> .

RPL expression	Meaning
<code>p / q</code>	Ordered choice: match <code>p</code> , and <code>p</code> fails, match <code>q</code>

Differences

- Always greedy
- Always possessive
- Choices are ordered

To regex, and beyond!

- Structured output, not flat
- Sane syntax
- A few key concepts
- Auto tokenization
- Package system
- Macros
- Unit tests
- REPL
- Trace output

Patterns in the standard library (v1.0.0)

■ Collections

- `net.any`, `date.any`, etc.
- `all.things`

■ Commonly needed

- int, float, hex, and other numbers
- several kinds of identifiers
- path names for Unix and Windows
- GUIDs

■ Network patterns

- ip address (v4, v6, mixed), domain name, email address, url, URI, MAC, HTTP

■ Timestamps

- RFC3339, RFC2822, and more than a dozen other common formats

■ CSV data

- delimiters: `,` `;` `|`
- quoted fields: `"foo"` or `'bar'`
- escapes: `""` or `\` or `\"`

■ JSON data

- full parse
- match nested and balanced `}` `[]`

■ Source code features

- 10 popular languages

■ De-structuring

- E.g. `"CSC316"` ==> `"CSC"`, `"316"`
- E.g. `"(1.2, 3.77, 0)"` ==> `"1.2"`, `"3.77"`, `"0"`

■ Log files

- syslog constituents (covers most log files)
- Java exceptions, Python tracebacks

Community

There are faster parsers for formats like JSON and CSV!

- Why use Rosie to parse JSON or CSV when there are special-purpose solutions for those that are very fast?
- Because you'll find those formats embedded into:
 - Semi-structured data, e.g. JSON-formatted log messages
 - Unstructured data, e.g. CSV as part of a larger piece of text/doc
- And in those cases, you can either separate the input and use different parsers on each part, or you can use one parser for the whole thing
- It comes down to volume, perhaps:
 - Specialized tools will run faster, and you'll need them if the volume of data in that format is high.
 - Otherwise, the “Swiss Army knife” approach may be better

The formal basis of RPL

- Rosie's operators are parser combinators
 - Based on Parsing Expression Grammars
 - Not CFG (slow!) or regex (limited!)
 - Express all deterministic (unambiguous) CFLs
 - And some non-CFLs, e.g. $a^n b^n c^n$
 - Key advantage: can match recursively structured input
- PEGs [Ford, 2004]
 - “Scanner-less parsing”
 - Linear time matching (at space cost)
 - Languages recognized by PEGs are
 - A superset of regular languages
 - All languages recognized by LL(k) and LR(k) parsers
- LPEG library [Ierusalimschy, 2008]
 - ➔ Gives a space-efficient PEG matching algorithm
 - ➔ Linear time in input size (non-grammars, no look-around)

Parsing Expression Grammars: A Recognition-Based Syntactic Foundation

Bryan Ford
Massachusetts Institute of Technology
Cambridge, MA

A Text Pattern-Matching Tool based on Parsing Expression Grammars

Roberto Ierusalimschy¹

¹ PUC-Rio, Brazil

This is a preprint of an article accepted for publication in Software: Practice and Experience; Copyright 2008 by John Wiley and Sons.

SUMMARY

Current text pattern-matching tools are based on regular expressions. However, pure regular expressions have proven too weak a formalism for the task: many interesting patterns either are difficult to describe or cannot be described by regular expressions. Moreover, the inherent non-determinism of regular expressions does not fit the need to capture specific parts of a match.

Motivated by these reasons, most scripting languages nowadays use pattern-matching tools that extend the original regular-expression formalism with a set of ad-hoc features, such as greedy repetitions, lazy repetitions, possessive repetitions, “longest match rule”, lookahead, etc. These ad-hoc extensions bring their own set of problems, such as lack of a formal foundation and complex implementations.

In this paper, we propose the use of Parsing Expression Grammars (PEGs) as a basis for pattern matching. Following this proposal, we present LPEG, a pattern-matching tool based on PEGs for the Lua scripting language. LPEG unifies the ease of use of pattern-matching tools with the full expressive power of PEGs. Because of this expressive power, it can avoid the myriad of ad-hoc constructions present in several current pattern-matching tools. We also present a Parsing Machine that allows a small and efficient implementation of PEGs for pattern matching.

KEY WORDS: pattern matching, Parsing Expression Grammars, scripting languages

sed on *generative*
free grammars, in
f rules applied re-
recognition-based
of rules or predi-
s in the language.
er paradigm. For
fnition of a trivial
are “constructed”
| (*|s| mod 2 = 0*) }
guage, in which a

ve paradigm, most
ience involve the
arsing, of strings.
ractical recogniz-
of parsing algo-

which the ubiqui-
expressions (REs)
for modelling and
heir elegance and
nerative grammars
ell. The ability of
tant and powerful
power gets in the
anguages that are
iguity in CFGs is



Rosie's matching engine is an
enhanced version of LPEG

Rosie is self-hosting

- Rosie is a parser, and Rosie is used to parse Rosie Pattern Language
- About 115 lines of RPL (core version) to define the current RPL version
- Could support multiple versions of RPL, even different dialects
- Non-trivial user extensions to RPL can be enabled by:
 - Specifying RPL for the extension (to RPL)
 - Writing a compiler “plug-in” for the extension
 - The compiler plug-in interface has not yet been designed... *hint!*

Match non-blank, non-comment lines of RPL:

```
$ rosie match -o data '!{[:space:]*$} !{[:space:]* "--"}' rpl_1_2.rpl | wc -l  
115
```

Roadmap

“If you want to go fast, go alone.
If you want to go far, go together.”

“Proverb”

Roadmap



1/14

Roadmap

Pattern generation

Algorithmic, e.g. from static analysis
Statistical / ML

Extensibility

User-written macros
User-written output encoders

Compiler Optimizations

Common subexpression elimination
New vm instructions
Flow analysis

Command line/scripting convenience

Traverse directories
Follow links or not, etc.

Regex-to-rosie converter

Re-use existing regex
Give them unit tests
Debug them

Ahead of time compilation

Fast startup
Small matching run-time (~50Kb binary)