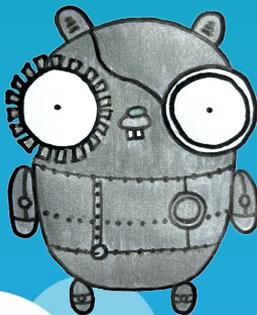




FOSDEM Go Devroom, February 2022

Fun with Finite Automata

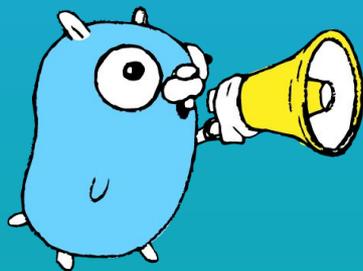


Bryan Boreham

Grafana Labs

 @bboreham

Overview



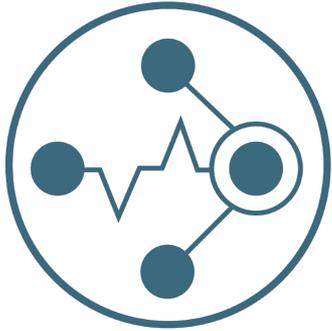
Who am I?

Why am I talking about this?

Automa-what?

How do regexps work?

Where did I speed them up?



cortex



loki



Tempo

This is where the story begins...



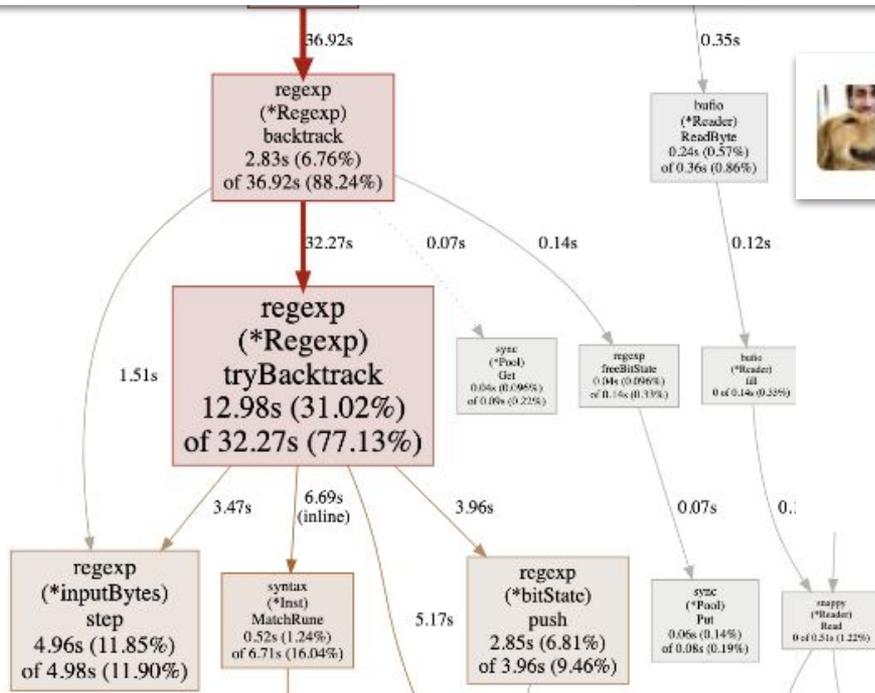
ewelch 17 3 months ago

someone running queries which spun the querier CPU cores to the max



Cyril 3 months ago

The regex is `|~ "(?i).*Finished.(?i).*heartbeat.*"`



"A **Regular Expression**, or Regexp, is used to express a search pattern in text" - *Wikipedia*

Example:

```
qui.*fox
```

Matches:

"The quick brown fox jumps
over the lazy dog"



.	any character
[xyz]	one of a set
xy	x followed by y
x y	x or y
x*	zero or more x
x+	one or more x
x?	zero or one x
^	beginning of text
\$	end of text
(re)	capturing group
(?i)	case-insensitive

Regular
expressions
are built out of
simple blocks.

A lot of simple blocks...



kinds of single-character expressions	examples
any character, possibly including newline	<code>.</code>
character class	<code>[xyz]</code>
negated character class	<code>[^xyz]</code>
Perl character class (link)	<code>\d</code>
negated Perl character class	<code>\D</code>
ASCII character class (link)	<code>[:alpha:]</code>
negated ASCII character class	<code>[^alpha:]</code>
Unicode character class (one-letter name)	<code>\pN</code>
Unicode character class	<code>\p{Greek}</code>
negated Unicode character class (1-letter)	<code>\PN</code>
negated Unicode character class	<code>\P{Greek}</code>

	Repetitions
<code>x*</code>	zero or more <code>x</code> , prefer more
<code>x+</code>	one or more <code>x</code> , prefer more
<code>x?</code>	zero or one <code>x</code> , prefer one
<code>x{n,m}</code>	<code>n</code> or <code>n+1</code> or ... or <code>m</code> <code>x</code> , prefer more
<code>x{n,}</code>	<code>n</code> or more <code>x</code> , prefer more
<code>x{n}</code>	exactly <code>n</code> <code>x</code>
<code>x*?</code>	zero or more <code>x</code> , prefer fewer
<code>x+?</code>	one or more <code>x</code> , prefer fewer
<code>x??</code>	zero or one <code>x</code> , prefer zero
<code>x{n,m}?</code>	<code>n</code> or <code>n+1</code> or ... or <code>m</code> <code>x</code> , prefer fewer
<code>x{n,}?</code>	<code>n</code> or more <code>x</code> , prefer fewer
<code>x{n}</code>	exactly <code>n</code> <code>x</code>

Character classes	
<code>\d</code>	digits ($\equiv [0-9]$)
<code>\D</code>	not digits ($\equiv [^0-9]$)
<code>\s</code>	whitespace ($\equiv [\t\n\f\r]$)
<code>\S</code>	not whitespace ($\equiv [^\t\n\f\r]$)
<code>\w</code>	word characters ($\equiv [0-9A-Za-z_]$)
<code>\W</code>	not word characters ($\equiv [^0-9A-Za-z_]$)

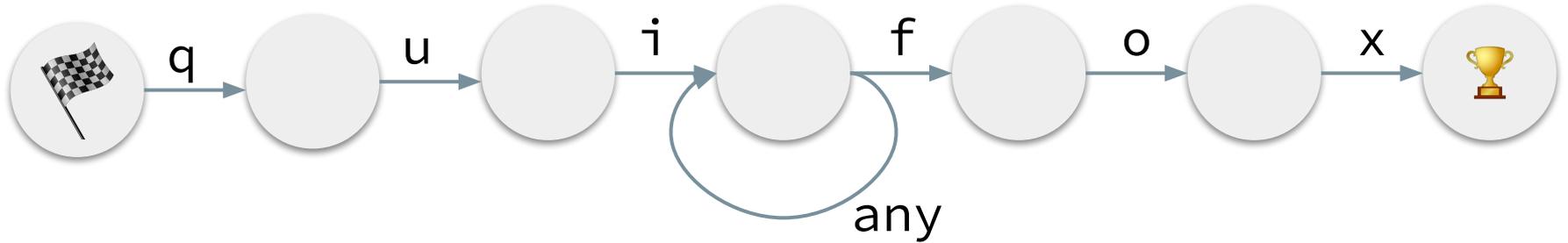
Composites	
<code>xy</code>	<code>x</code> followed by <code>y</code>
<code>x y</code>	<code>x</code> or <code>y</code> (prefer <code>x</code>)

Empty strings	
<code>^</code>	at beginning of text or line (<code>m=true</code>)
<code>\$</code>	at end of text (like <code>\z</code> not <code>\Z</code>) or line (<code>m=true</code>)
<code>\A</code>	at beginning of text
<code>\b</code>	at ASCII word boundary (<code>\w</code> on one side and <code>\W</code> , <code>\s</code> or <code>\S</code> on the other)
<code>\B</code>	not at ASCII word boundary
<code>\z</code>	at end of text

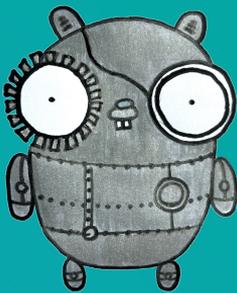
Grouping	
<code>(re)</code>	numbered capturing group (submatch)
<code>(?P<name>re)</code>	named & numbered capturing group (submatch)
<code>(?:re)</code>	non-capturing group
<code>(?flags)</code>	set flags within current group; non-capturing
<code>(?flags:re)</code>	set flags during re; non-capturing

Flags	
<code>i</code>	case-insensitive (default false)
<code>m</code>	multi-line mode: <code>^</code> and <code>\$</code> match begin/end line (default false)
<code>s</code>	let <code>.</code> match <code>\n</code> (default false)
<code>U</code>	ungreedy: swap meaning of <code>x*</code> and <code>x*?</code> , <code>x+</code> and <code>x+?</code> , etc

qui.*fox



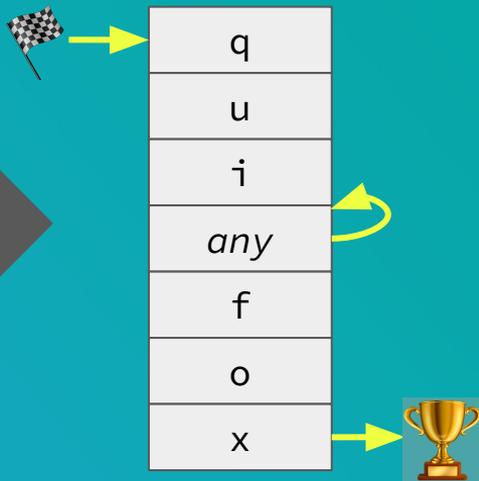
Nondeterministic
Finite
Automaton



Regex compilation in Go

qui.*fox

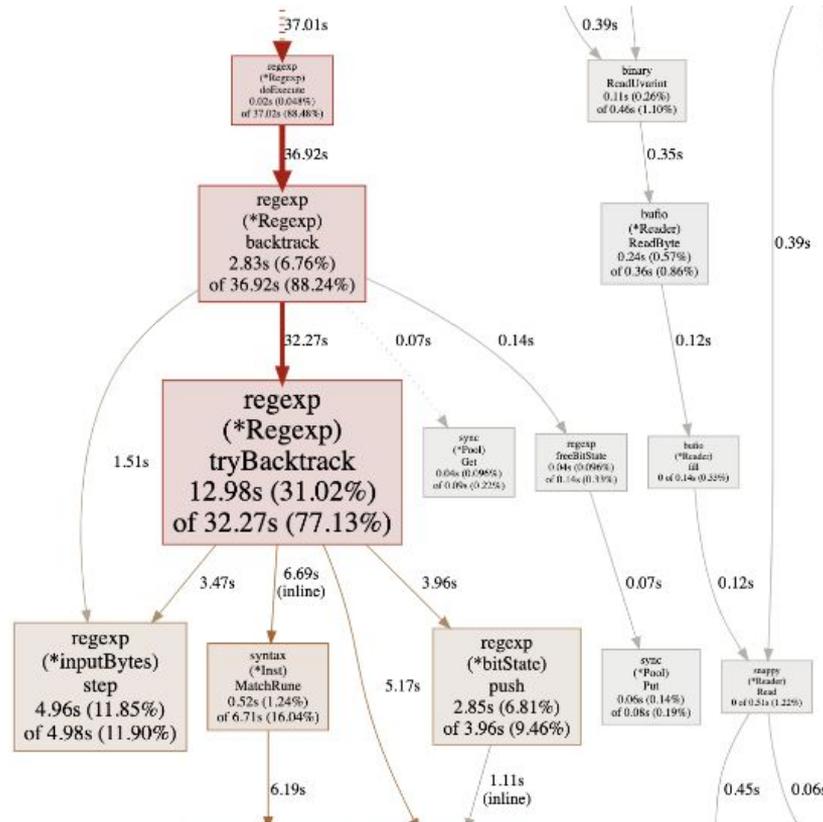
compile



This is the data structure for a compiled regexp:

```
type Prog struct {  
    Inst []Inst  
    Start int  
    NumCap int  
}  
  
type Inst struct {  
    Op InstOp  
    Out uint32  
    Arg uint32  
    Rune []rune  
}  
  
type InstOp uint8  
    InstAlt  
    InstAltMatch  
    InstCapture  
    InstEmptyWidth  
    InstMatch  
    InstFail  
    InstNop  
    InstRune  
    InstRune1  
    InstRuneAny  
    InstRuneAnyNotNL
```

Back to the profile



```
func BenchmarkRegexp(b *testing.B) {  
    b.StopTimer()  
    x := bytes.Repeat([]byte("x"), 50)  
    re := MustCompile("(?i).*foo.*bar.*")  
    b.StartTimer()  
    for i := 0; i < b.N; i++ {  
        re.Match(x)  
    }  
}
```

```
go test -run xxx -bench Regexp -cpuprofile cprof
```

Benchmark results

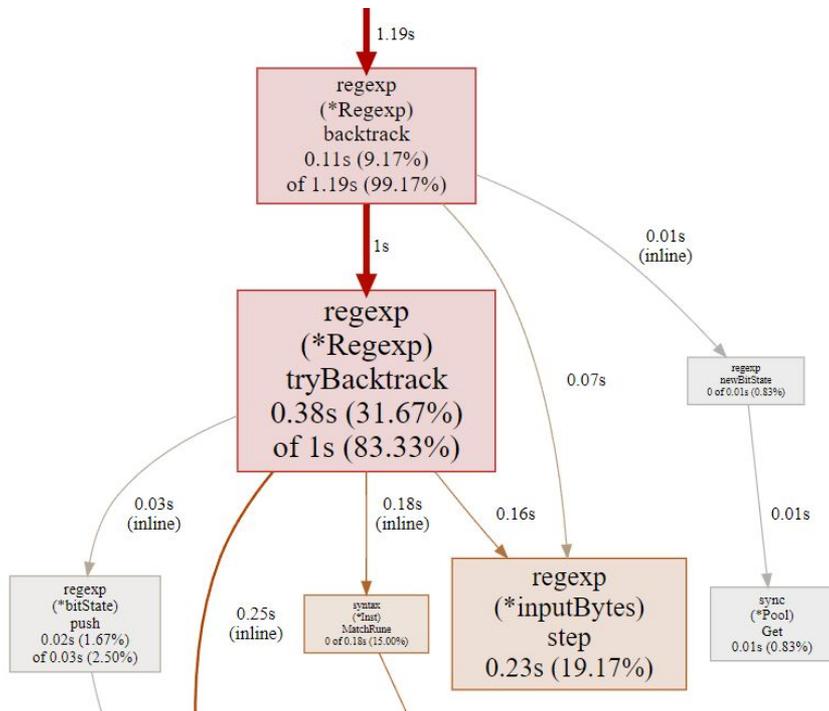


```
$ go test -run xxx -bench Regexp -cpuprofile cprof
goos: linux
goarch: amd64
pkg: github.com/grafana/regexp
cpu: Intel(R) Core(TM) i5-9400F CPU @ 2.90GHz
BenchmarkRegexp-4          283477          4123 ns/op
PASS
ok      regexp         1.415s
```

Profile of the benchmark



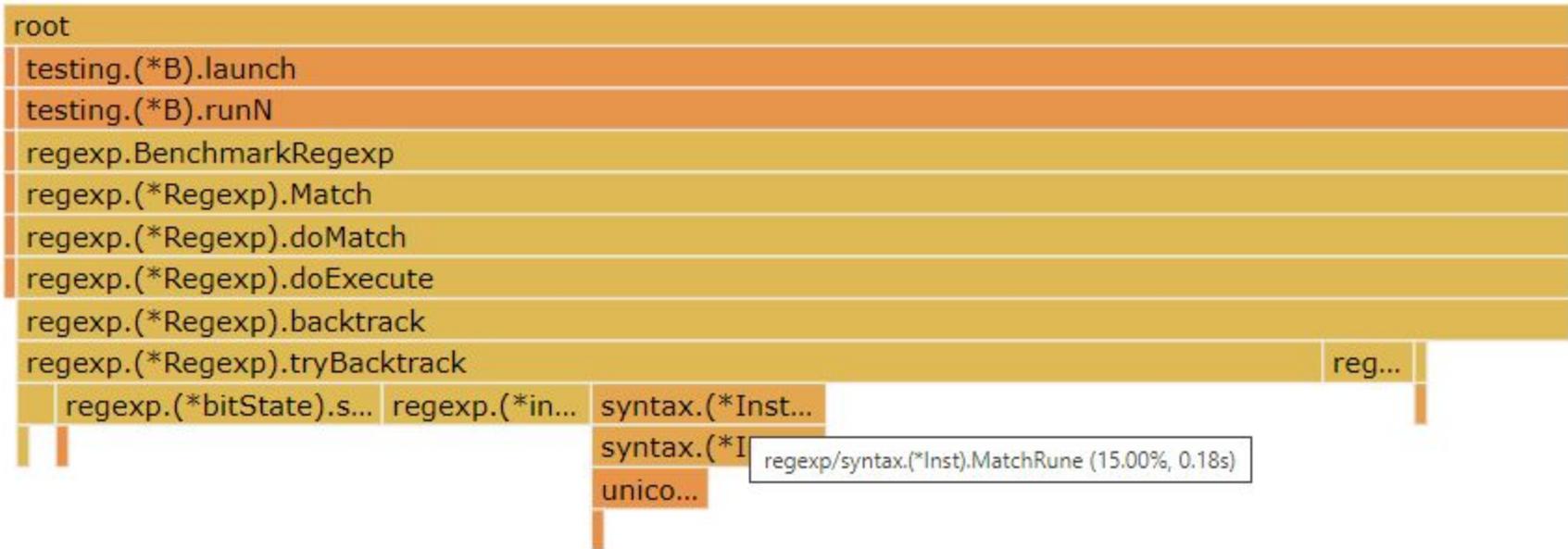
```
go tool pprof -http :8000 cprof
```



Profile of the benchmark (flame graph)



regexp/syntax.(*Inst).MatchRune (15.00%, 0.18s)



```
func BenchmarkRegexp(b *testing.B) {
    b.StopTimer()
    x := bytes.Repeat([]byte("x"), 50)
    re := MustCompile("(?i).*foo.*bar.*")
    b.StartTimer()
    for i := 0; i < b.N; i++ {
        re.Match(x)
    }
}
```

```
"(?i).*foo.*bar.*"
```

```
4123 ns/op
```

```
".*foo.*bar.*"
```

```
3413 ns/op
```

"foo.*bar"

91.25 ns/op

"foo"

90.85 ns/op

Regex



One-pass

Backtrack

NFA

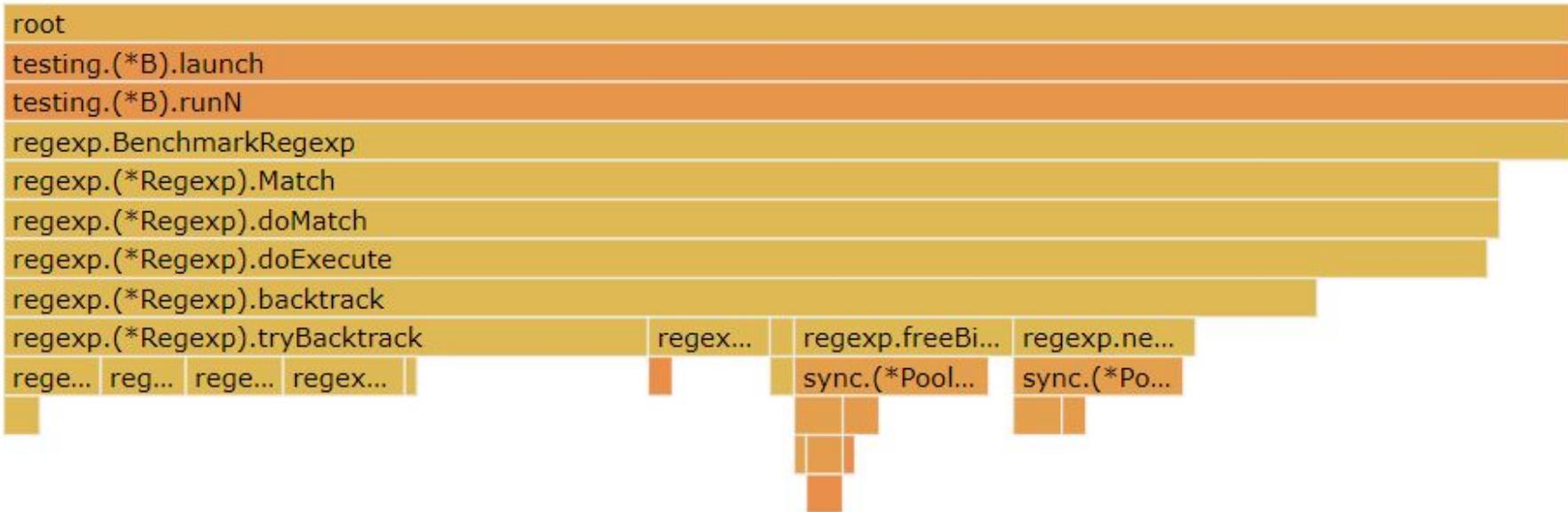
Try to force one-pass



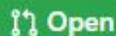
"^foo"

115.1 ns/op

Profile of benchmark "^foo"



regexp: allow patterns with no alternates to be one-pass #48748



bboreham wants to merge 3 commits into `golang:master` from `bboreham:more-onepass`

Conversation 8

Commits 3

Checks 0

Files changed 2



bboreham commented on 2 Oct 2021 • edited

Contributor



Check whether a regex has any 'alt' instructions before rejecting it as one-pass. Previously `^abc` would run the backtrack matcher.

I tried to make the comment match what the code does now.

Updates #21463

Reviewers

No reviews

Still in progress? Cc

Assignees

No one assigned

Active ☆ 353711 regexp: allow patterns with no alternates to be one-pass

EDIT

Change Info

SHOW ALL

REPLY

- Owner: GerritBot
- Author: Bryan Boreham
- Reviewers: Russ Cox
- CC: Bryan Boreham, Brad Fitzpatrick, Michael Matlo..., Gopher Robot
- Repo | Branch: go | master
- Topic:

Submit requirements

- Code-Review: No votes
- Untrusted: Gopher Robot

Other labels

SHOW LESS

- Run-TryBot: No votes
- Trust: No votes
- TryBot-Result: No votes

regexp: allow patterns with no alternates to be one-pass

Check whether a regex has any 'alt' instructions before rejecting it as one-pass. Previously '^abc' would run the backtrack matcher.

I tried to make the comment match what the code does now.

Updates #21463

...

name	old time/op	new time/op	delta
Find-8 (p=0.500 n=5+5)	167ns ± 1%	170ns ± 3%	~
FindAllNoMatches-8 (p=0.095 n=5+5)	88.8ns ± 5%	87.3ns ± 0%	~

SHOW ALL

EDIT

Checks: No results

Comments: 5 resolved

```
469 469 » // onepass regexp is anchored
470 470 » if prog.Inst[prog.Start].Op != syntax.InstEmptyWidth ||
471 471 »     » syntax.EmptyOp(prog.Inst[prog.Start].Arg)&syntax.EmptyBeginText != syntax.EmptyBegin
Text {
472 472 »     » return nil
473 473 » }
474 474 » // every instruction leading to InstMatch must be EmptyEndText
474 474 » hasAlt := false
475 475 » for _, inst := range prog.Inst {
476 476 »     » if inst.Op == syntax.InstAlt || inst.Op == syntax.InstAltMatch {
477 477 »         »     » hasAlt = true
478 478 »         »     » break
479 479 »         »     }
480 480 »     }
481 481 » // If we have alternates, every instruction leading to InstMatch must be EmptyEndText.
482 482 » // Also, any match on empty text must be $.
475 483 » for _, inst := range prog.Inst {
476 484 »     » opOut := prog.Inst[inst.Out].Op
477 485 »     » switch inst.Op {
478 486 »     »     » default:
479 487 »         »         » if opOut == syntax.InstMatch {
487 487 »         »         » if opOut == syntax.InstMatch && hasAlt {
480 488 »         »         »     »     » return nil
481 489 »         »         »     }
482 490 »         »         » case syntax.InstAlt, syntax.InstAltMatch:
```



Now this case is much faster



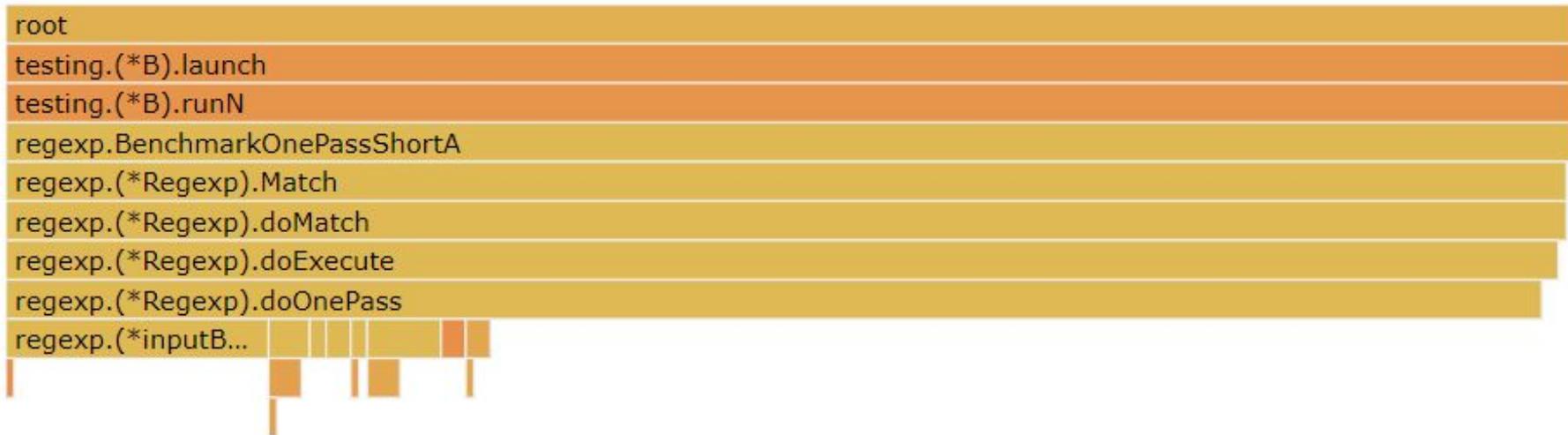
"^foo"

78.09 ns/op

```
func BenchmarkOnePassShortA(b *testing.B) {
    b.StopTimer()
    x := []byte("abcdddddeeeededd")
    re := MustCompile("^bc(d|e)*$")
    b.StartTimer()
    for i := 0; i < b.N; i++ {
        re.Match(x)
    }
}
```

730.6 ns/op

Profile of BenchmarkOnePassShortA



```
regexp.(*Regexp).doOnePass in regexp/exec.go
 440ms    720ms (flat, cum) 98.63% of Total
    .      .    398: func (re *Regexp) doOnePass ...
□
    .      .    444:     for {
170ms    170ms  445:         inst = re.onepass.Inst[pc]
 40ms    40ms  446:         pc = int(inst.Out)
 40ms    40ms  447:         switch inst.Op {
    .      .    448:         default:
```

□

Simple change



[355789](#) regexp: avoid copying each instruction executed [src/regexp/exec.go](#) File 3 of 3 | [Prev](#)

Base [browse](#) → Patchset 1 [browse](#) | DOWNLOAD [SHOW BLAME](#) | EDIT | Diff view:

```
440 440 »      »      r1, width1 = i.step(pos + width)
441 441 »      »      flag = i.context(pos)
442 442 »      »      pc = int(re.prefixEnd)
443 443 »      »      }
444 444 »      »      for {
445 445 »      »      inst = re.onepass.Inst[pc]
445 445 »      »      inst = &re.onepass.Inst[pc]
446 446 »      »      pc = int(inst.Out)
447 447 »      »      }
```

730.6 -> 633.0 ns/op

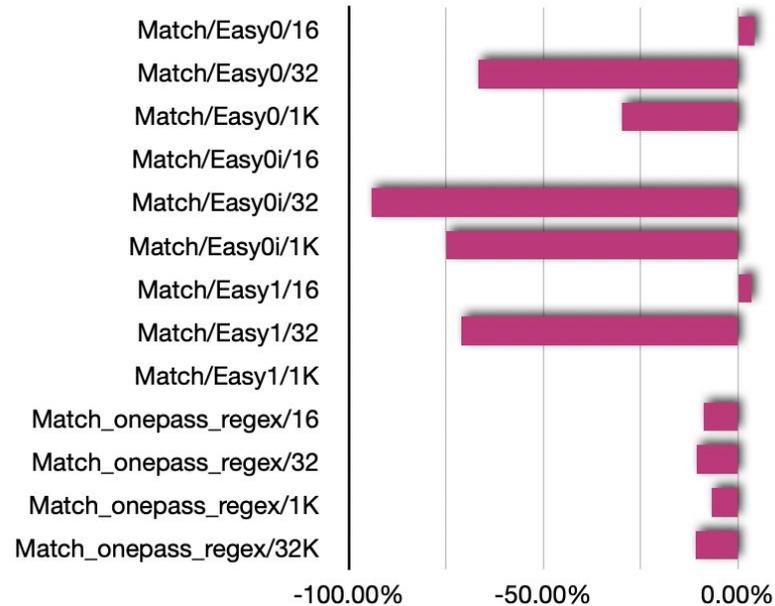
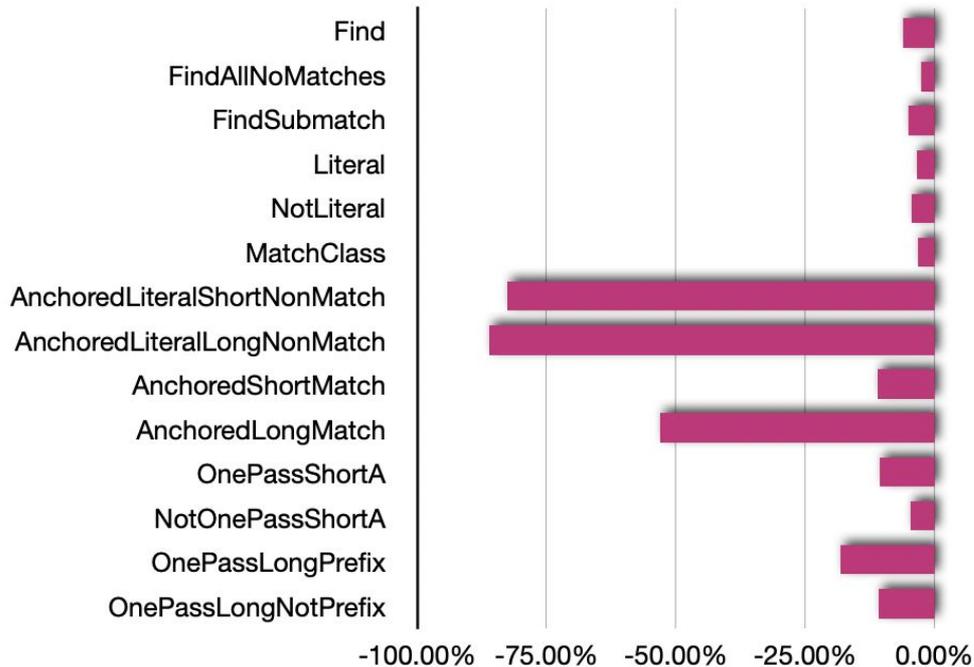


Brad Fitzpatrick

Wow, nice find!

- [353711](#): allow patterns with no alternates to be one-pass
- [354909](#): speed up onepass prefix check
- [355789](#): avoid copying each instruction executed
- [358756](#): handle prefix string with fold-case
- [377294](#): allow prefix string anchored at beginning

Overall Score



<https://github.com/grafana/regexp/tree/speedup>

<https://github.com/google/re2/wiki/Syntax>



GopherBot automaton by Mary Pollard.

Other gophers and the original Gopher design by Renee French.