



This deck has been modified from its original version. It has been formatted to fit your PDF reader

The 'live' version may have a few more slides and some dynamic content

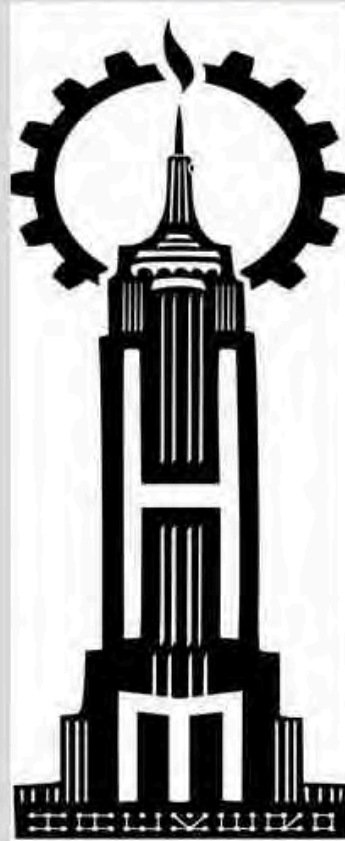
# RE: Exploring Regular Expression Denial of Service



# Places I've Hacked @



**NYCRESISTOR**



# About Me

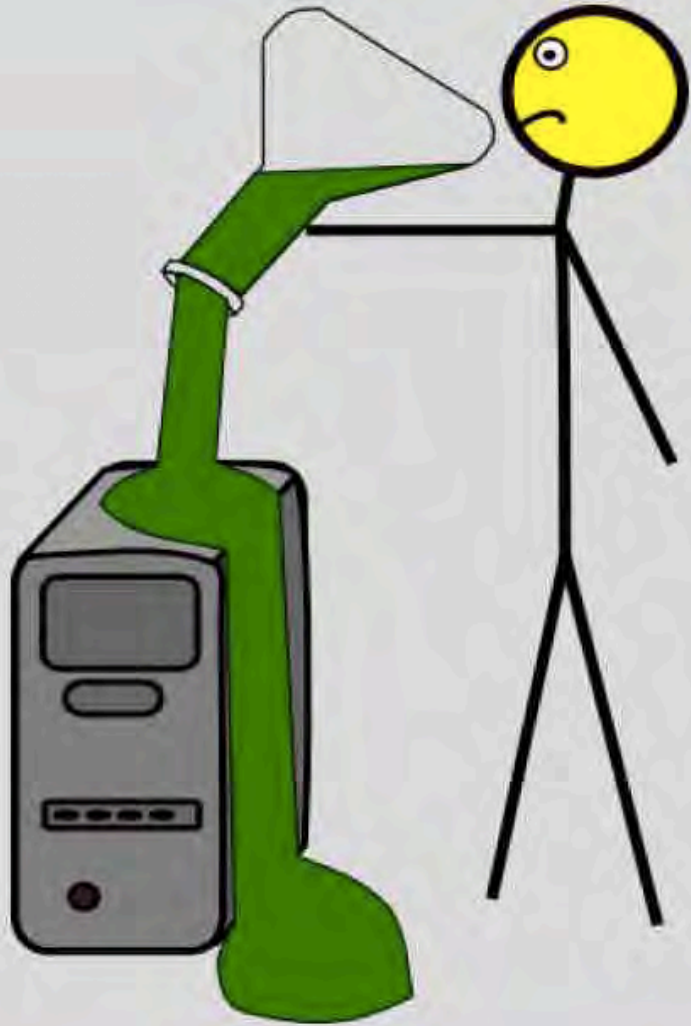
I'm Eric (XlogicX) - Not a "Security Researcher"

email: [no.axiom@gmail.com](mailto:no.axiom@gmail.com)

blog: [xlogicx.net](http://xlogicx.net)

github: <https://github.com/XlogicX>

RE Benchmark: <https://github.com/XlogicX/8ball/blob/master/benchrexes.pl>



## **RE:TL;DR**

**a{5,15}** iz 5 to 15 a's

**a?** iz a{0,1}

**a+** iz a{1,inf}

**a\*** iz a{0,inf}

**.** iz 1 char wildcard

**(afT)** iz 'a' or 'f' or 'T'

**[afT]** iz 'a' or 'f' or 'T' and perform moar good

**[^afT]** iz not those a,f, and T's

**\d** iz [0123456789]

**\s** iz whitespace character

**\w** iz [a-zA-Z\d]

**\+** iz just a '+', iz not {1,inf}

**^** anchor start, **\$** anchor end

# Quantifiers

}

?

+

\*

# Freeform Range {}'s

Problem: We are looking for 5-15 instances of 'x'

Regex: `x{5,15}`



# Useful Aliases of {}'s

? = {0,1}

+ = {1,infinity}

\* = {0,infinity}

Problem: We are looking for 1 or more x's

Regex: x+

# Character Classes []'s

A set of specific characters

Problem: Looking for a string of 5-8 3's, 5's, or 9's

Regex: `[359]{5,8}`

String: This number (**5935935**) would make this string match

# Negative Space

Negating a set of specific characters

Problem: grab all text until a comma is reached

Regex: `[^,]+`

String: **All this text would be matched**, but not this

# Useful Character Class Aliases

`\s` = whitespace; space, tab, newline

`\d` = numbers

`\w` = alphanumeric

`\S`, `\D`, and `\W` = `[^\s]`, `[^\d]`, and `[^\w]`

`.` = any character but a newline

# Grouping and Alternation

Problem: Looking for words good, bad, or evil

Regex: `good|bad|evil`

String: This non-evil sentence would match

# Grouping and Alternation

Problem: Looking for words good, bad, or evil 3-4 times in a row

Regex: `(good | bad | evil){3,4}`

String: This text would make this string match:

**badgoodevilbad**

# Anchoring ^

Problem: Match string if it starts with the word "anchor"

Regex: ^anchor

Matching string: **anchor** is an anchor

Non-matching string: boat is an anchor

# Anchoring \$

Problem: Match string if it ends with the word "anchor"

Regex: anchor\$

Matching string: anchor is an **anchor**

Non-matching string: anchor is an a boat



# Escaping

Problem: You want to search for 3-6 \$'s

Regex: `\${3,6}`

WAT?: Because \$ is regex character (anchor); must be escaped with '\'

# Greediness/Laziness

String: <script>not really</script>just  
text<script>still...</script>moar text

Regex: <script>.+</script>

Matches: <script>not really</script>just  
text<script>still...</script>

Regex: <script>.+?</script>

Matches: <script>not really</script>

# RE:DoS

RE:respect

"Some people, when confronted with a problem, think 'I know, I'll use regular expressions.' Now they have two problems." - Jamie Zawinski

# RE:Recon

There are two popular regular expression algorithms; Deterministic Finite Automata (DFA) and Non-Deterministic Finite Automata (NFA)

Abuse of these two engines require different strategies

# DFA vs NFA

Some easily testable behavior differences in the two engines are:

- Longest alternation vs first alternation

- Laziness

- Possessiveness

# Greps

The following examples will use 2 greps:

`grep -P (NFA)`

`egrep (DFA)`

# Longest Alternation

NFA engines favor the first alternative  
DFA engines favor the longest alternative

```
xlogicx@xlogicxMVM ~ $ echo "ab" | grep -P --only-matching "a|ab"  
a  
xlogicx@xlogicxMVM ~ $ echo "ab" | egrep --only-matching "a|ab"  
ab
```

# Laziness

DFA engines can't be lazy

```
xlogicx@xlogicxMVM ~ $ echo "ababa" | grep -P --only-matching "a.+?a"  
aba  
xlogicx@xlogicxMVM ~ $ echo "ababa" | egrep --only-matching "a.+?a"  
ababa
```



# Possesiveness

NFA - greediness favors lexemes in order  
DFA - priority on longest overall match

```
xlogicx@xlogicxMVM ~ $ echo "abc" | grep -P --only-matching "a(b)?(bc)?"  
ab  
xlogicx@xlogicxMVM ~ $ echo "abc" | egrep --only-matching "a(b)?(bc)?"  
abc
```

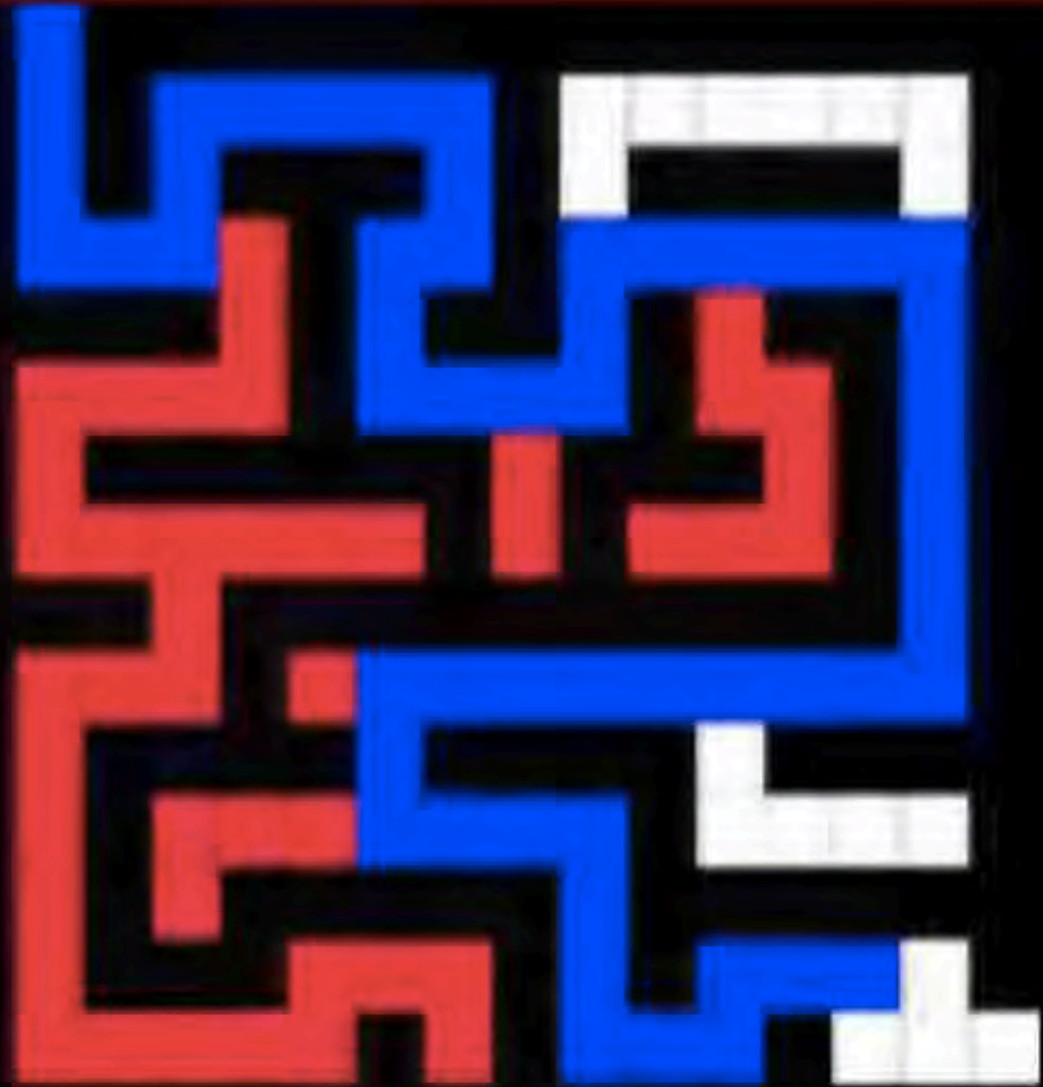
# NFA vs DFA part I

We will use a labyrinth metaphore to start our discussion of differences between these two engines

It should be noted that this metephore is imperfect; it is only intended to visualize a couple of concepts:

- NFAs do a lot of backtracking

- DFAs process options in parallel





# NFA vs DFA part II

The following 'animations' are also intended to be a visual aid (not an exact representation)

NFAs follow the expression

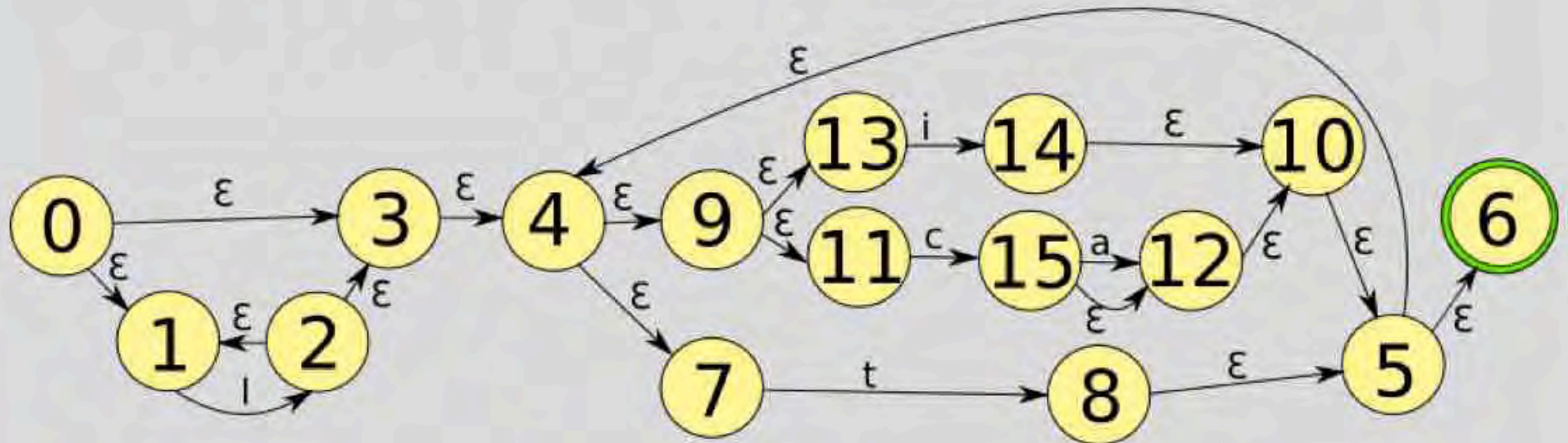
DFAs follow the string

(Click for each frame)

l\*(t|ca?|**i**) +

app**lication**on

# NFA State Diagram



Compiling REx "l\*(t|ca?|i)+"

Final program:

```
1: STAR (4)
2:   EXACT <l> (0)
4:  CURLYX[0] {1,32767} (24)
6:   OPEN1 (8)
8:     TRIE-EXACT[ci] (21)
      <t> (21)
      <c> (14)
14:      CURLY {0,1} (21)
16:      EXACT <a> (0)
      <i> (21)
21:    CLOSE1 (23)
23:  WHILEM[1/1] (0)
24:  NOTHING (25)
25:  END (0)
```





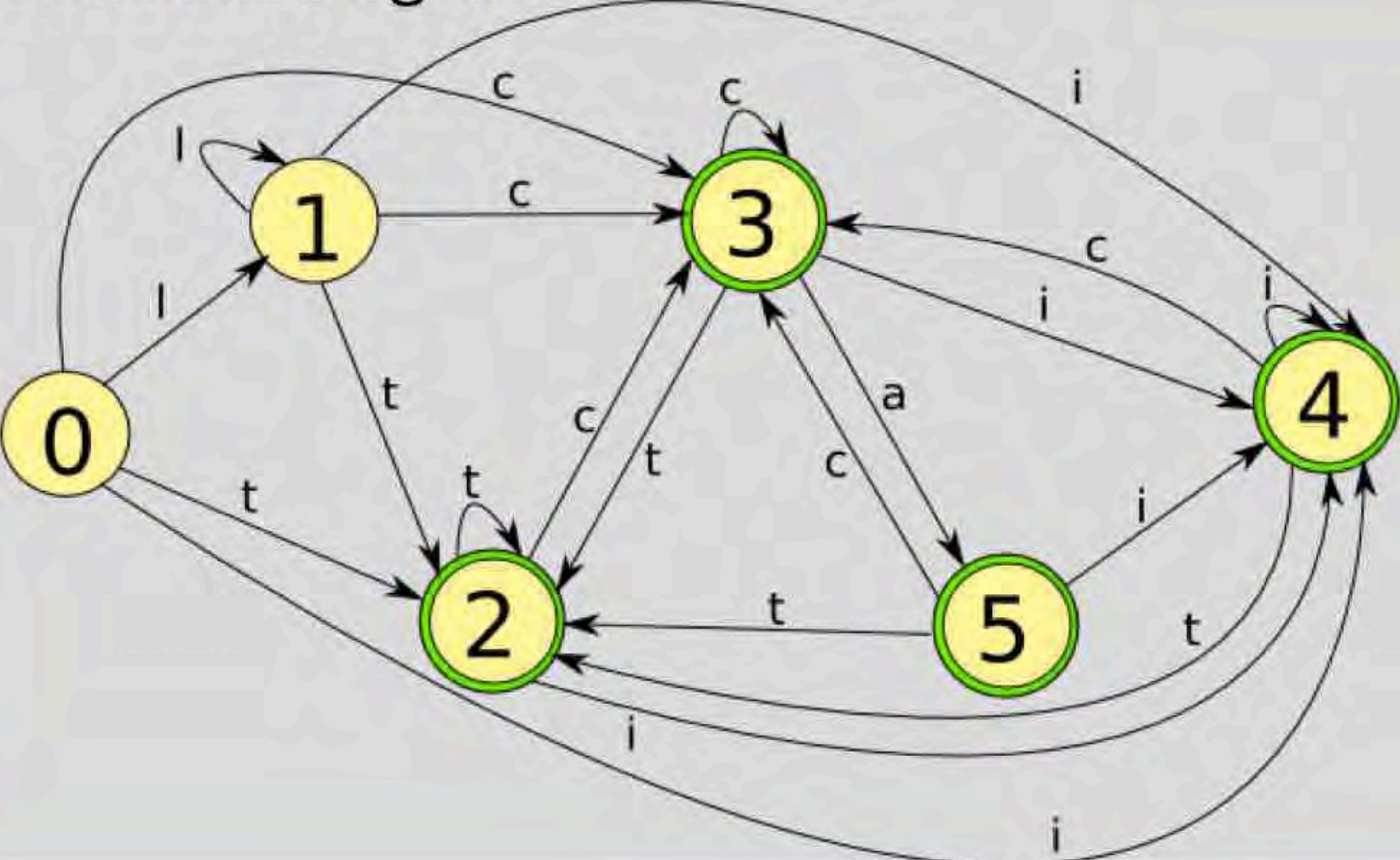
(Click for each frame)

l\*(t|ca?|**i**) +

app**lication**on



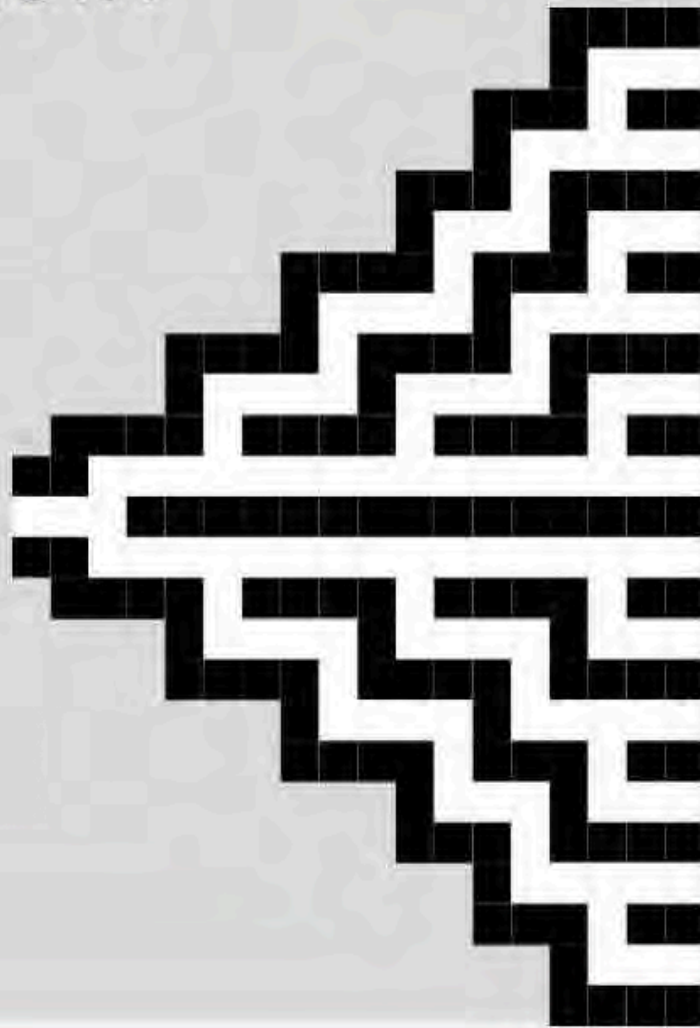
# DFA State Diagram



## DFA Table

	l	t	c	i	a	N/A	
State 0	1	2	3	4	-1	-1	
State 1	1	2	3	4	-1	-1	
State 2	-1	2	3	4	-1	-1	Match
State 3	-1	2	3	4	5	-1	Match
State 4	-1	2	3	4	-1	-1	Match
State 5	-1	2	3	4	-1	-1	Match

# Abusing DFA



# Consider the State Table

State Table uses memory

Increase memory usage by:

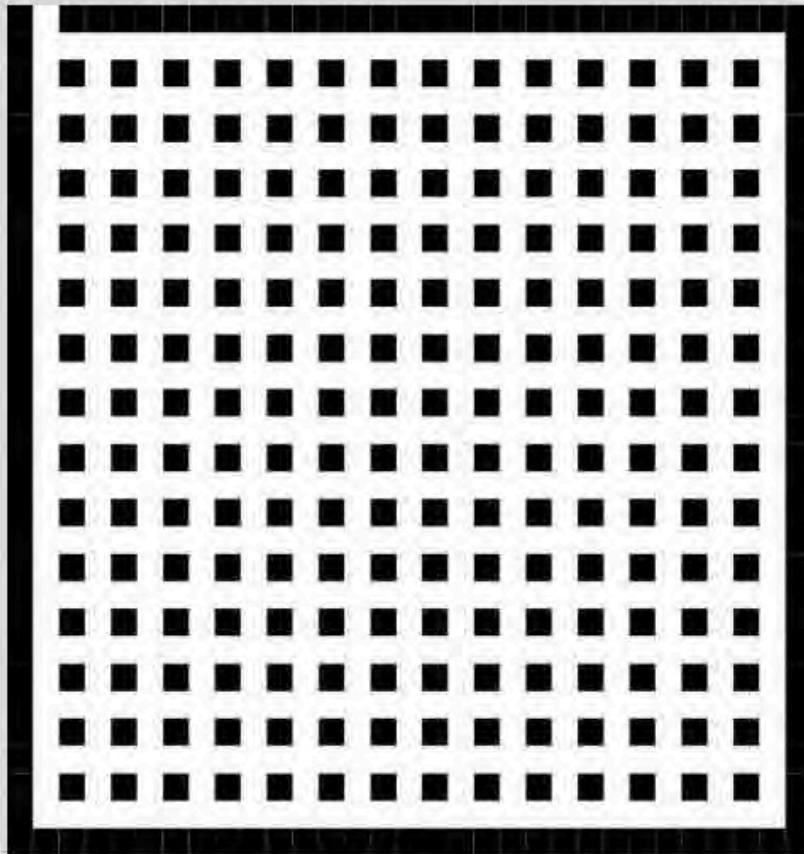
- Adding more states (very effective)

- Adding more lexemes

PoC: DFA State Flood

$((a_{0,75})_{0,75})_{0,75}?$

# Abusing NFA





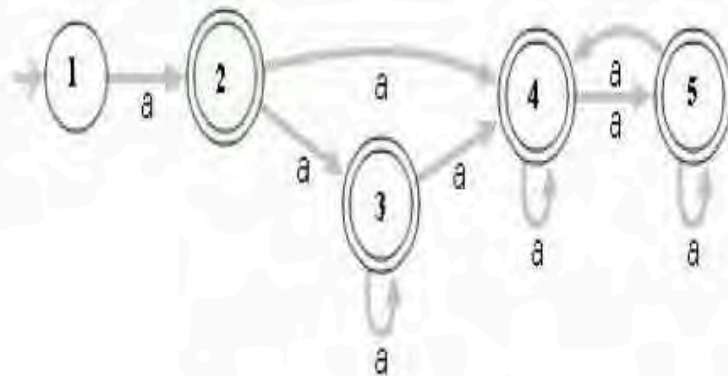
# Considerations

DFA tries **everything** until it finds a match

Until it finds a match...

# OWASP example

For example, the Regex  $^(a+)$$  is represented by the following NFA:



For the input  $aaaaX$  there are 16 possible paths in the above graph. But for  $aaaaaaaaaaaaaaaaX$  there are 65536 possible paths, and the number is double for each additional  $a$ . This is an extreme case where the naïve algorithm is problematic, because it must pass on many many paths, and then fail.

# Optimization

OWASP ReDoS:  $^(a+)+\$$



# Optimization

Optimized: ^a+\$



# Some Simple C

```
1  #include <stdio.h>
2
3  int main() {
4      if (5 > 0) {
5          printf("True\n");
6      } else {
7          printf("False\n");
8      }
9  }
```

# What Actually Happens

```
00000000:0040052d 55          push rbp
00000000:0040052e 48 89 e5    mov rbp, rbp
00000000:00400531 bf c4 05 40 00 mov edi, 0x004005c4
00000000:00400536 e8 d5 fe ff ff call 0x0000000000400410
00000000:0040053b 5d          pop rbp
00000000:0040053c c3          ret
00000000:0040053d 8f 1f 00    nop dword ptr [rax]
00000000:00400540 41 57          push r15
00000000:00400542 41 89 ff    mov r15d, edi
00000000:00400545 41 56          push r14
00000000:00400547 49 89 f6    mov r14, rsi
00000000:0040054a 43 55          push r13
00000000:0040054c 49 89 d5    mov r13, rdx
```

0x0000000000400410 = 0000000000400410

Registers

RSI:	00007fff544188a8
EDI:	00000000004005c4 ASCII "True"
R8:	00007f6778eae80
R9:	00007f6778ec4560
R10:	00007fff54418650

Bookmarks

Address	Comment
---------	---------

Add Del Clear

## Data Dump

+ 0000000000400000-0000000000401000

00000000:004005b6 f3 c3 00 00 48 83 ec 08 48 83 c4 00 c3 00 00 00 00	0A .H.L.H.A.A. ....
00000000:004005c8 01 00 02 00 54 72 75 65 06 00 00 00 01 1b 03 3b	...True.....
00000000:004005d0 30 00 00 00 05 00 00 00 34 fe ff ff 7c 00 00 00	0.....4py)...

## Stack

00007fff:544187c0	0000000000000000	.....
00007fff:544187c8	00007f6778b0fec5	Ap*xg... return to
00007fff:544187d0	0000000000000000	.....
00007fff:544187d8	00007fff544188a8	.ATy...
00007fff:544187e0	0000000100000000	.....
00007fff:544187e8	000000000040052d	-0.....

# Throw A Wrench In

```
1  #include <stdio.h>
2
3  int main() {
4      int derp = 5;
5      if (derp > 0) {
6          printf("True\n");
7      } else {
8          printf("False\n");
9      }
10 }
```

# What Happened

```
00000000:0040052d 55          push rbp
00000000:0040052e 48 89 e5    mov rbp, rsp
00000000:00400531 48 83 ec 10 sub rsp, 16
00000000:00400535 c7 45 fc 05 00 00 00 mov dword ptr [rbp-4], 5
00000000:0040053c 83 7d fc 00 cmp dword ptr [rbp-4], 0
00000000:00400540 7e 0c      jle 0x000000000040054e
00000000:00400542 bf e4 05 40 00 mov edi, 0x0040054
00000000:00400547 e8 c4 fe ff ff call 0x0000000000400410
00000000:0040054c eb 0a      jmp 0x0000000000400558
00000000:0040054e bf e9 05 40 00 mov edi, 0x004005e9
00000000:00400553 e8 b8 fe ff ff call 0x0000000000400410
00000000:00400558 c9        leave
00000000:00400559 c3        ret
```

rbp = 0000000000000000

ata Dump

0000000000400000-0000000000401000

```
00000000:004005d0 f3 c3 00 00 48 83 ec 08 48 83 c4 08 c3 00 00 00 00 0A H.I.H.A.A.
00000000:004005e0 01 00 02 00 54 72 75 65 00 46 61 6c 73 65 00 00 ... True False
```



# Wrench for Regex

Finding something else like '+'

+ means '1 or more'

{1,9001} is kind of like that

So:  $^(a\{1,9001\})\{1,9001}\$$

# PoC: NFA

```
XlogicXs-Air:8ball XlogicX$ /usr/bin/time perl -e 'if("aaaaaaaaaaaaaaaaaaaaaab"=~/(a{1,30}){1,30}$/){}'  
0.59 real 0.58 user 0.00 sys  
XlogicXs-Air:8ball XlogicX$ /usr/bin/time perl -e 'if("aaaaaaaaaaaaaaaaaaaaaab"=~/(a{1,30}){1,30}$/){}'  
1.20 real 1.19 user 0.00 sys  
XlogicXs-Air:8ball XlogicX$ /usr/bin/time perl -e 'if("aaaaaaaaaaaaaaaaaaaaaab"=~/(a{1,30}){1,30}$/){}'  
2.26 real 2.24 user 0.00 sys  
XlogicXs-Air:8ball XlogicX$ /usr/bin/time perl -e 'if("aaaaaaaaaaaaaaaaaaaaaab"=~/(a{1,30}){1,30}$/){}'  
4.56 real 4.54 user 0.01 sys  
XlogicXs-Air:8ball XlogicX$ /usr/bin/time perl -e 'if("aaaaaaaaaaaaaaaaaaaaaab"=~/(a{1,30}){1,30}$/){}'  
9.43 real 9.37 user 0.02 sys  
XlogicXs-Air:8ball XlogicX$ /usr/bin/time perl -e 'if("aaaaaaaaaaaaaaaaaaaaaab"=~/(a{1,30}){1,30}$/){}'  
18.31 real 18.23 user 0.04 sys  
XlogicXs-Air:8ball XlogicX$ /usr/bin/time perl -e 'if("aaaaaaaaaaaaaaaaaaaaaab"=~/(a{1,30}){1,30}$/){}'  
37.85 real 37.66 user 0.09 sys  
XlogicXs-Air:8ball XlogicX$ /usr/bin/time perl -e 'if("aaaaaaaaaaaaaaaaaaaaaab"=~/(a{1,30}){1,30}$/){}'  
74.26 real 73.91 user 0.17 sys  
XlogicXs-Air:8ball XlogicX$ /usr/bin/time perl -e 'if("aaaaaaaaaaaaaaaaaaaaaab"=~/(a{1,30}){1,30}$/){}'  
149.36 real 148.67 user 0.35 sys  
XlogicXs-Air:8ball XlogicX$ /usr/bin/time perl -e 'if("aaaaaaaaaaaaaaaaaaaaaab"=~/(a{1,30}){1,30}$/){}'  
295.81 real 294.48 user 0.68 sys  
XlogicXs-Air:8ball XlogicX$ /usr/bin/time perl -e 'if("aaaaaaaaaaaaaaaaaaaaaab"=~/(a{1,30}){1,30}$/){}'  
602.57 real 599.44 user 1.50 sys  
XlogicXs-Air:8ball XlogicX$ /usr/bin/time perl -e 'if("aaaaaaaaaaaaaaaaaaaaaab"=~/(a{1,30}){1,30}$/){}'  
1204.79 real 1198.21 user 3.09 sys  
XlogicXs-Air:8ball XlogicX$ /usr/bin/time perl -e 'if("aaaaaaaaaaaaaaaaaaaaaab"=~/(a{1,30}){1,30}$/){}'  
2398.76 real 2386.04 user 6.00 sys  
XlogicXs-Air:8ball XlogicX$ /usr/bin/time perl -e 'if("aaaaaaaaaaaaaaaaaaaaaab"=~/(a+)+$/){}'  
0.02 real 0.00 user 0.01 sys  
XlogicXs-Air:8ball XlogicX$ /usr/bin/time perl -e 'if("aaaaaaaaaaaaaaaaaaaaaab"=~/^(a+)+$/){}'  
0.02 real 0.00 user 0.01 sys
```

# PoC: Results Digest

34 characters takes 40 minutes (on my Air)

Every additional 'a' doubles the time to complete

So 50 characters of this would take roughly 10 years

# Automation for DFA's

Because of the way the DFA engine works, only the expression matters; not the string

As a review, we are attacking memory

I use RE2 and slowly starve it of resources to see where things fall (for each expression)

# Automation for NFA's

The regular expression by itself is not inefficient

The expression may, however, be vulnerable to strings that make it go super-linear

# Strategies for a TERRIBLE string

Match as much as possible, while failing at the end  
(long-circuit attack)

When a quantifier is found, max it out

When alternations are seen, pick the last one







All of This is a  
Challenge to  
Automate...

# Challenge Accepted!

Before getting into benchmarking (the good stuff)

Let's look at what strings the automation produces

(Because it's unintentionally funny)

hi

ha

FYI: Green matches, Red does not...

ha

h1

$a_{\{1,15\}}h!$

aaaaaaaaaaaaaha

(ab | cd | ef | gh | ij | kl | mn |  
op | qr | st | uv | wx | yz)d1

yzda

h+i

hhhhhhhhhhhhhhhhhh

hhhhhhhhhhhhhhhhhh

hhhhhhhhhhhhhhhhhh

hha

$x(a\{1,10\})\{1,10\}$

x1



$x(a\{1,10\})\{1,10\}\$$

$xaaaaaaaaaaaaaaaa$

$aaaaaaaaaaaaaaaa$

$aaaaaaaaaaaaaaaa$

$aaaaaaaaaaaaaaaa$

$aaaaaaaaaaaaaaaa1$

# Benchrexes.pl

Because I can't think of a clever name for it

- Feed the script a txt file with regular expressions (1 per line)
- It will output a csv of benchmarks
- It benchmarks for both time and memory use

# Results

Now lets look at some interesting real world results from

- Emerging Threats IDS rules
- [regexlib.com](http://regexlib.com)

Tested in x64 Linux Mint 17 VM @ 1.7GHz (i7) with 1 GB of RAM

# Regexlib.com/Memory

Most Complete URL Validator

```
^((http(s){0,1}\:VV){0,1}([a-z|A-Z|0-9|\.\|\-|_]){4,255}(\:\d{1,5})  
{0,1}){0,1}((V([a-z|A-Z|0-9|\.\|\-|_]|\%[A-F|a-f|0-9]{2})){1,255})  
{1,255}V{0,1}){0,1}(|V{0,1}\?[a-z|A-Z|0-9|\.\|\-|_]{1,255}\=([a-  
z|A-Z|0-9|\.\|\-|_|+|\:]\%[A-F|a-f|0-9]{2}|\&[a-z|A-Z]{2,12}  
\;){0,255}){0,1}((\&[a-z|A-Z|0-9|\.\|\-|_]{1,255}\=([a-z|A-Z|0-9|  
\.\|\-|_|+|\:]\%[A-F|a-f|0-9]{2}|\&[a-z|A-Z]{2,12}\;){0,255})  
{0,255})(V{0,1}|\#[a-z|A-Z|0-9|\.\|\-|_|+|\:]\%[A-F|a-f|0-9]{2}|\&  
[a-z|A-Z]{2,12}\;){0,255})$
```

149 MB

# Regexlib.com/Memory

Validates a long (windows) filename

```
^[^\\.\/*?\"<>\\|]{1}[^\\V:\/*?\"<>\\|]{0,254}$
```

660 KB

# Regexlib.com/Time

Siteswap validator

```
^((0)*(1|(2|(3|(4|50)0)0)0)*|(01*(2|(3|(4|50)0)0))*|(00(1|20)*|(3|(4|50)0))*|(000(1|(2|30)0)*|(4|50))*|(0000(1|(2|(3|40)0)0)*5)*|(2|(5(3|50)0|(4|51)2|40)0|(3|(4|51)1|(52|(4|51)50)0)(501|(4|5050)0)*1|(3|50(2|40)0))*|(501|(4|5050)0|(1|(3|50(2|40)0)0)(2|(5(3|50)0|(4|51)2|40)0)*3|(4|51)1|(52|(4|51)50)0)*|(0(2|(3|(4|51)1|(52|(4|51)50)0)(501|(4|5050)0)*1*(5(3|50)0|(4|51)2|40)|(3|(4|51)1|(52|(4|51)50)0)(501|(4|5050)0)*3|50(2|40)))*(0(501)*(4|5050)0(501)*(1|(3|50(2|40)0)0)(2|(5(3|50)0|(4|51)2|40)0)|(3|(4|51)1|(501)*(1|(3|50(2|40)0)0))*52|(4|51)50|(3|(4|51)1|(501)*(4|5050))*|((01|0500)(40)*5|(0(2|40)0|(01|0500)(40)*(1|30))2|(5(3|50)0|(4|51)2|40)0|(3|(4|51)1|(52|(4|51)50)0)(40)*(1|30))*3|(4|51)1|(52|(4|51)50)0)(40)*5)*|((1|500)(40)*50|((2|40)0|(1|500)(40)*(1|30))2|5(3|50)00|(3|520)(40)*(1|30))*4|51|(3|520)(40)*50))*|(00(2|(4|51)20|3|(4|51)1|(52|(4|51)50)0)(501|(4|5050)0)*1|(3|50(2|40)0))*5(3|50)|(4|51)4|(3|(4|51)1|(52|(4|51)50)0)(501|(4|5050)0)*504))*|00(501|40)*505|00(501|40)*(1|(3|50(2|40)0)0)(2|(5(3|50)0|(4|51)2|40)0)|(3|(4|51)1|520)(501|40)*(1|(3|50(2|40)0)0))*((4|51)5|(3|(4|51)1|520)(501|40)*505))*|(((3|50)0|(1|(2|40)0)0|(11|(2|150)0)(501|(4|5050)0)*1|(3|50(2|40)0)0)(2|4(2|40)0|(3|41|4500)(501|(4|5050)0)*1|(3|50(2|40)0)0))*5)*|(000(2|(530|(4|51)2|40)0)|(3|(4|51)1|(52|(4|51)50)0)(501|(4|5050)0)*1|(3|50(2|40)0)0))*55)*|(3|55(2|50)0|5(3|51)(1|40)|(4|52|5(3|51)50)|51|(4|550)50)*2|530|(4|550)(1|40))*|(51|(4|550)50|(2|530|(4|550)(1|40))3|55(2|50)0|5(3|51)(1|40))*4|52|5(3|51)50))*|(50|51)*(4|550)|1|40|50(51)*2|530)|3|55(2|50)0|(4|52|51)*2|530))*5(3|51)|(4|52|51)*(4|550))*|(0(3|5(3|51)1|(4|52|5(3|51)50)(51|(4|550)50)*2|(4|550)1))*55(2|50)|5(3|51)4|(4|52|5(3|51)50)(51|(4|550)50)*53|(4|550)4))*|(0|51)*(4|550)5|0(51)*2|530|(4|550)(1|40))3|55(2|50)0|5(3|51)(1|40)|(4|52|51)*2|530|(4|550)(1|40))*|(5(3|51)5|(4|52|51)*4|550)5))*|(1|5050)(450)*5|(30|50(1|40)|(1|5050)(450)*2|4(1|40))3|55(2|50)0|5(3|51)(1|40)|(4|52|5(3|51)50)(450)*2|4(1|40))*4|52|5(3|51)50)(450)*5)|(050(51|450)*55|(0(1|40)|050(51|450)*2|530|4(1|40))3|55(2|50)0|5(3|51)(1|40)|(4|52|5(3|51)50)(51|450)*2|530|4(1|40))*|(4|52|5(3|51)50)(51|450)*55)*|((5(2|50)0|(3|51)(1|40)|(2|(3|51)50)(51|(4|550)50)*2|530|(4|550)(1|40)))3|4|51|(4|550)50)*2|530|(4|550)(1|40))*5)*|(((2|50)0|1(1|40)|150(51|(4|550)50)*2|530|(4|550)(1|40))3|53(1|40)|(4|52|5350)(51|(4|550)50)*2|530|(4|550)(1|40))*55)*|(00(3|5520|5(3|51)(1|40)|(4|52|5(3|51)50)(51|(4|550)50)*2|530|(4|550)(1|40))*555)*|(4|5(3|5(2|5(1|50))))*(3|5(2|5(1|50)))4*5)*|((2|5(1|50))(4|53)*55)*|(1|50)(4|5(3|52))*555)*|(0(4|5(3|5(2|51)))5555)*|(5)*$
```

.01 seconds

# IDS Rule/Memory

GPL SQL time\_zone buffer overflow attempt

```
TIME_ZONE\s*=\s*((\x27[^\x27]{1000,})|(\x22[^\x22]{1000,}))
```

325 KB

# IDS Rule/Time

ET ACTIVEX ImageShack Toolbar Remote Code Execution

```
<object\s*[\^>>]*\s*classid\s*=\s*[\x22\x27]\s*clsid  
\s*\x3a\s*{\s*DC922B67-FF61-455E-  
9D79-959925B6695C\s*}\s*(.*)\>
```

1.6 seconds





# Real Bad Regex (Mem)

This slide not yet populated; it will be for the live presentation

# Real Bad Regex (Time)

This slide not yet populated; it will be for the live presentation