

**P L A Y H A P P Y S O U N D
F O R P L A Y E R T O
E N J O Y**

UNDERSTANDING CATASTROPHIC REGEX BACKTRACKING



Regex is essentially a pattern that search engines use to match strings, enabling for find, find and replace activities or input validation. Regex originated back in 1951 blah blah history, but because heavily used in the late 60s for pattern matching and lexical analysis.

Who knows What GREP stands for?


Globally search for a Regular Expression and Print matching lines.

How many people know regex? Of those, who has had to learn regex the same amount of times they've had to use it.

Often people say that you started with a problem and decided to use regex, now you have two problems. This talk will show just how true that is.




Lets take an example of a simple pattern we want to match – an IP Address. Could be many different reasons.



`\d+ \. \d+ \. \d+ \. \d+`

IP ADDRESS



$\backslash d+ \backslash . \backslash d+ \backslash . \backslash d+ \backslash . \backslash d+$
192.168.0.1

IP ADDRESS

`\d+\.\d+\.\d+\.\d+`

`192.168.0.2222222222222222`

IP ADDRESS

HAPPY SOUND

7

Unfortunately we can see the pattern doesn't properly validate what we're looking for.



`(?:\d{1,3}\.){3}\d{1,3}`

IP ADDRESS



(?:\d{1,3}\.){3}\d{1,3}


192.168.400.1

IP ADDRESS

HAPPY SOUND

9

Still it leaves room for improvement



(?:1?(?:1?\d{1,2}|2[0-4]\d|25[0-5])\.){3}(?:1?\d{1,2}|2[0-4]\d|25[0-5])

IP ADDRESS

2001:DB8:85A3:8D3:1319:8A2E:370:734
8

IPV6 ADDRESS

HAPPY SOUND

11

((?:?:1?(1?\d{1,2}|2[0-4]\d|25[0-5])\.(1??\d{1,2}|2[0-4]\d|25[0-5])\.(1?\d{1,2}|2[0-4]\d|25[0-5])\.(1?\d{1,2}|2[0-4]\d|25[0-5]))|(?:(?:1?((?:?:[0-9A-Za-f]{1,4}:){7}[0-9A-Za-f]{1,4}|(?:?:[0-9A-Za-f]{1,4}:){0,7}[0-9A-Za-f]{1,4}\z)|(((0-9A-Za-f){1,4}:){1,7}|:)((:[0-9A-Za-f]{1,4}){1,7}|:))))))

SMILEYFACE GUNEMOJI

BACKTRACKING

A\d+	A	1	2	3	4	
A\d+	A	1	2	3	4	

H A P P Y S O U N D

13

Here's another example of how backtracking works even though there 'appears' to be a match.

Progressing through the first line the regex fails because the . can't match anything.

The regex backtracks to A123, matches 0 non-digits, then matches 4 as the .

BACKTRACKING

/\d\d+\\$/g	0	4	0	8	-	0	1	7	-	8	4	2
/\d\d+\\$/g	0	4	0	8	-	0	1	7	-	8	4	2
/\d\d+\\$/g	0	4	0	8	-	0	1	7	-	8	4	2
/\d\d+\\$/g	0	4	0	8	-	0	1	7	-	8	4	2
/\d\d+\\$/g	0	4	0	8	-	0	1	7	-	8	4	2
/\d\d+\\$/g	0	4	0	8	-	0	1	7	-	8	4	2
/\d\d+\\$/g	0	4	0	8	-	0	1	7	-	8	4	2

H A P P Y S O U N D

14

Very simplified vision of how backtracking works. What we're looking for is the last four digits of a phone number – Verification questions?

As the regex progresses through the phone number it has to backtrack several times as the match fails.

STACK OVERFLOW OUTAGE

JULY 20, 2016

1 4 4 4

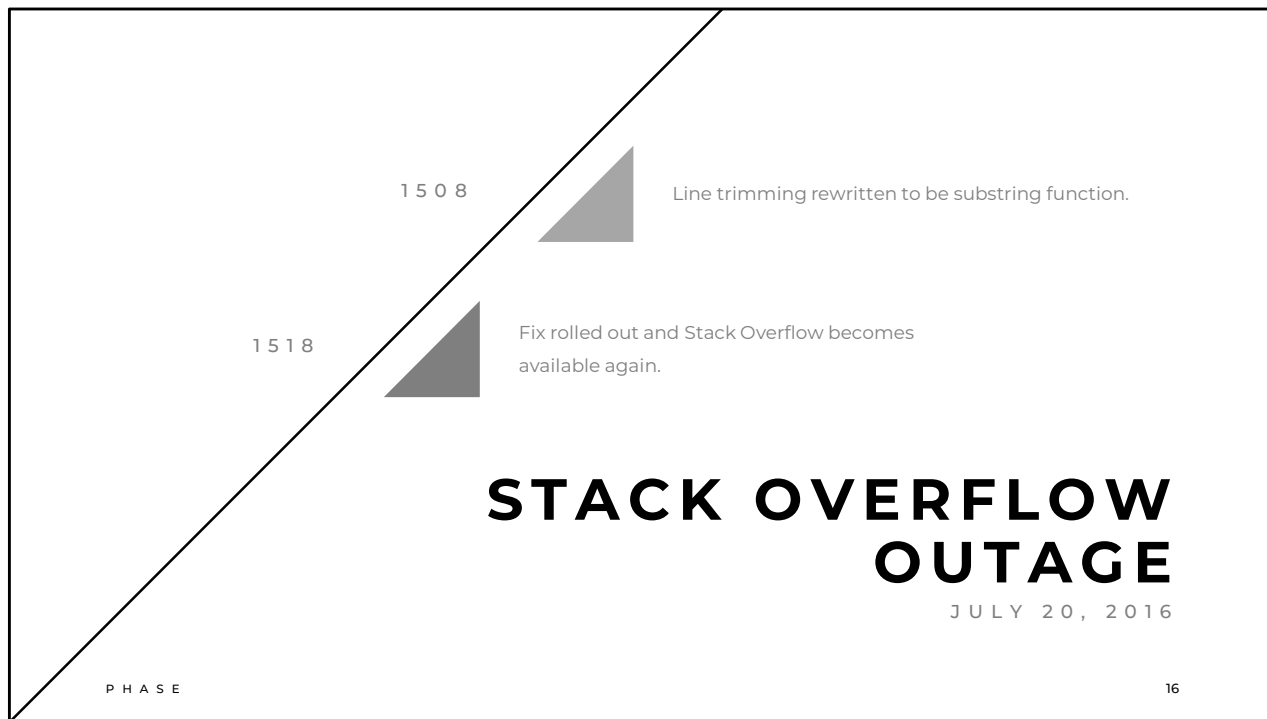
User nguyentrugkien submits a question: *"Join tiles in Corona SDK into one word for a Breakout game grid?"*

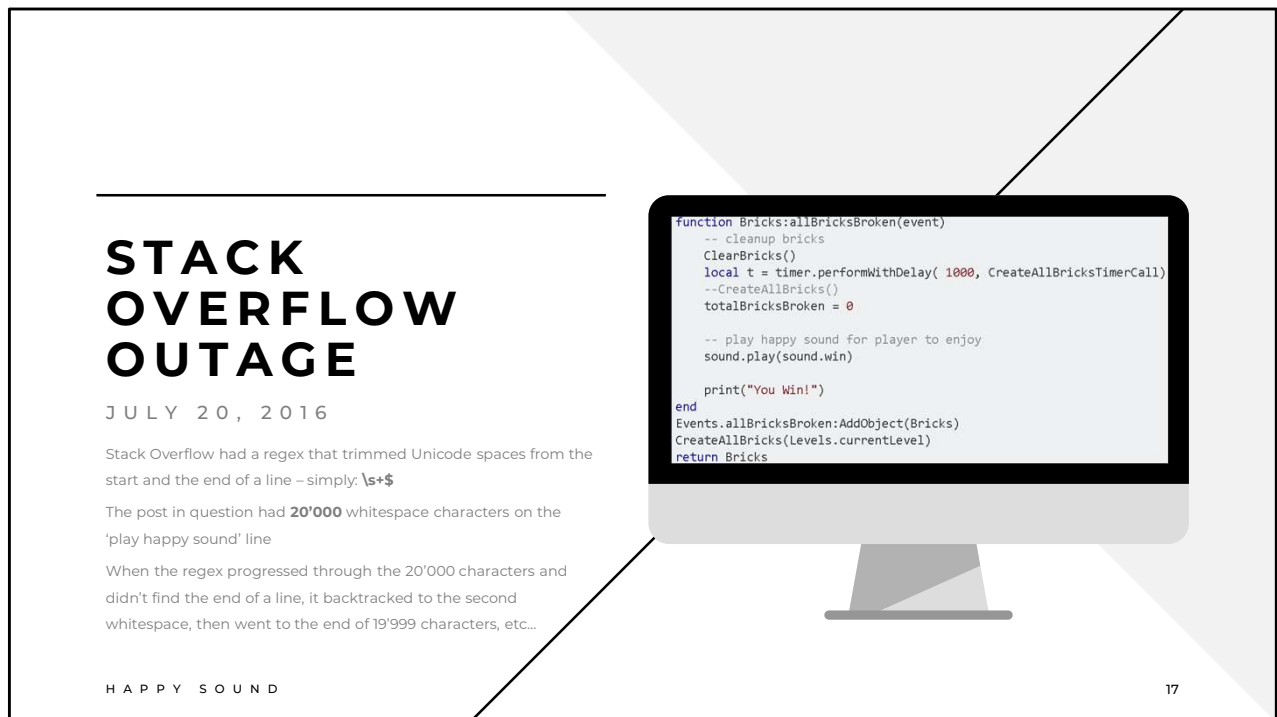
1 4 5 4

Memory dump leads Stack Overflow team to the specific regex and the 'Front Page' (last 3000 posts). Reran all posts against the regex to find which post took the longest to run.

HAPPY SOUND

15





So what caused the stackoverflow outage?

They had a regex element that trimmed Unicode spaces from the start or end of a line. To break a long regex down to the simple element that caused the problem `\s$`


Somehow, when the user copied part of their code to stackoverflow, they managed to include 20'000 whitespaces on the line “– play happy sound for player to enjoy”

The regex would start at the beginning, progress through all 20'000 characters, then fail as the line didn't end.

It would backtrack, and start at the second character, for all 19'999 characters, then fail as the line didn't end.

It would backtrack and start at the third character, for all 19'998 characters, then ...

```
any any any -> any any (msg:"99_GLOBAL_0365_Azure_ActiveDirectory";adsid
:586;content:"Workload";;content:"\AzureActiveDirectory";json:pcr:
"\x22Workload\x22\x3a\x22\x22AzureActiveDirectory\x22";pcr:$(Operation);
"(?<TrimmedOperation>.*(?<\x2e5))";pcr:
0Home\x22\x3a\x22UserAgent\x22\x2c\x22Value\x22\x3a\x22(?<Pcu>{["\x22]+)"
;pcr:$(group);"targetName\x22\x2c\x22Value\x22\x3a\x22{["\x22]+}\x22";
pcr:$(RequestType);
"\x7b\x22Name\x22\x3a\x22RequestType\x22\x2c\x22Value\x22\x3a\x22(.*)\x2
2\x7d.";pcr:$(UserAuthenticationMethod);
"\x7b\x22Name\x22\x3a\x22UserAuthenticationMethod\x22\x2c\x22Value\x22\x3
a\x22(.*)\x22\x7d.";pcr:$(ResultStatusDetail);
"\x7b\x22Name\x22\x3a\x22ResultStatusDetail\x22\x2c\x22Value\x22\x3a\x22(
.*)\x22\x7d.";pcr:$(GroupName);
"\x7b\x22Name\x22\x3a\x22Group_DisplayName\x22\x2c\x22NewValue\x22\x3a\x2
2(.*)\x22\x2c";pcr:$(GroupChangeTarget);
"\x7b\x22Name\x22\x3a\x22targetUPN|targetName\x22\x2c\x22Value\x22\x3a\x2
2(.*)\x22\x7d.";var@Access_Privileges:$(UserType);var@action:$(
ResultStatus);var@AppID:$(Workload);var@Authentication_Type:$(
UserAuthenticationMethod);var@CommandID:$(RecordType);var@
Destination_Logon_ID:$(GroupChangeTarget);var@firsttime,lasttime:$(
CreationTime);var@ObjectID:$(GroupName);var@ObjectID:$(group);
var@Request_Type:$(RequestType);var@Security_ID:$(UserKey);var@
Session_Status:$(ResultStatusDetail);var@sigdesc,sigid:$(Workload)
"$TrimmedOperation);var@src_ip:$(ClientIP);var@src_ip:$(
ActorIpAddress);var@UserIDSrc:$(UserId);var@Status:$(LogonError);
var@User_Agent:$(ue);map@actions:$(success="0",fail="
0";failed="0";failed="0";map@CommandID:"1"=1 - Exchange Admin
Audit","2"=2 - Single Mailbox Audit","3"=3 - Multiple Mailbox Audit",
"4"=4 - Assigning Permissions","6"=6 - File or Folder Operation","8"=
"8 - Admin Operation Performed","9"=9 - OrgId Logon Event","10"=10 -
Security Cmdlet Event","11"=11 - Data Loss Protection Event","12"=12
- Sway Event","14"=14 - Sharing Event","15"=15 - Secure Token Service
Logon Event","18"=18 - Security and Compliance Center Event","20"=20
- Power BI Event","22"=22 - Yammer Event","25"=25 - Microsoft Teams
Event","26"=26 - Microsoft Teams Event","27"=27 - Microsoft Teams
Event"; map@Access_Privileges:"0"=0 - Regular User","2"=2 - Office
365 Admin","3"=3 - Datacenter Admin or System Account","4"=4 - System
Account","5"=5 - An Application","6"=6 - Service Principal";
map@Attribute_Type:"0"=0 - Account Login Event","1"=1 - Application
Security Event");
HAPPY SOUND
```



CATASTROPHIC BACKTRACKING

Catastrophic backtracking often occurs when nested quantifiers are repeated, or alternated tokens exist within a group that itself is repeated.

Tldr; Catastrophic Backtracking is always the dev's fault

18

This is what's known as catastrophic backtracking. Why do I know so much about this? Had to deal with an outage myself courtesy of McAfee.

On the left-hand side here we can see an example of a parsing file that takes O365 Azure AD events for their SIEM product. Unfortunately I don't have the regex example available that caused this, but coming up hopefully you can take a guess what it looked like.

Catastrophic backtracking occurs when nested quantifiers are repeated, or groups themselves are repeated.

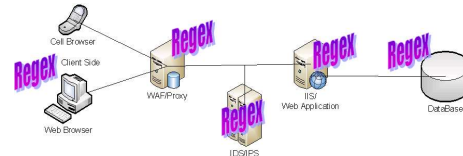
Catastrophic backtracking has been a problem since regex has been a thing. The solution to catastrophic backtracking has existed almost as long (first being identified and proposed in 1968). In more modern times, the .NET regex engine has been known to cause lots of problems (in the early 2000s, Hotmail devs facing 'catastrophic' issues... unfortunately, the blogs around that have been deleted).

WHAT DOES THIS MEAN FOR SECURITY

HAPPY SOUND

DENIAL OF SERVICE

Primarily, catastrophic regex backtracking can cause denial of service. This could be triggered against users, against edge infrastructure. Almost every path on the internet contains regex.



EXAMPLES OF EVIL REGEX

(a+)+ (a|aa)+
([a-zA-Z]+)* (a|a?)+

19

As far as security is concerned, the primary problem that regex is going to cause you is denial of service – just like with the example that we saw with stack overflow.

At basically every stage of browsing the web you'll come up against a regex *somewhere* in the backend of your ops.

Evil regex – regex that can cause catastrophic backtracking – are usually relatively simple. They are almost always groups that are being repeated, or groups that are being repeated having further repetition or alternation inside.

1, perl5.18

Starting regex match
|
V
/.*.*=.*/

|
V
'x=x'

[Visual of regex at 're.pl' line 6] [step: 0]

.*.*=.*

This is an example of a repeated regex group looking for 0 or more of any character.

The first .* acts greedily and matches the entire string. The second .* can't match anything so doesn't. Then we hit the literal character of = and the match fails.

The regex backtracks, returns to the first .* which matches x=, etc.

HAPPY SOUND

20

.*.*=.* is essentially zero or more of any character, followed by zero or more of any character, followed by a literal =, followed by zero or more of any character.

When looking at a simple match such as x=x it takes 23 steps.

The first step it takes is the .* matches the entire string. It can't match the second .*, it doesn't care, but when it gets to = it fails and backtracks.

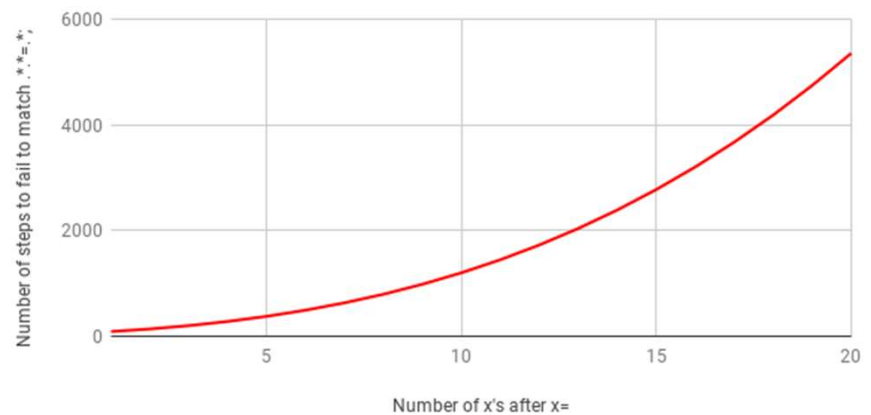
foo=bar;

Previously our strings matched. If we send partially matching strings into a poorly designed regex the backtracking becomes catastrophic.

Failing to match `x=x` takes 90 steps.

Failing to match 20 `x`'s takes 5353 steps.

Failing to match `.*.*=.*;` against the string `x=` followed by repeated `x`'s



HAPPY SOUND

22

If we add a semi-colon, such as trying to find the string `foo=bar;` our test string will always fail. But instead of taking 23 steps it takes 90 steps for `x=x`

For failing to match `x= 20 x's` it takes 5353 steps.

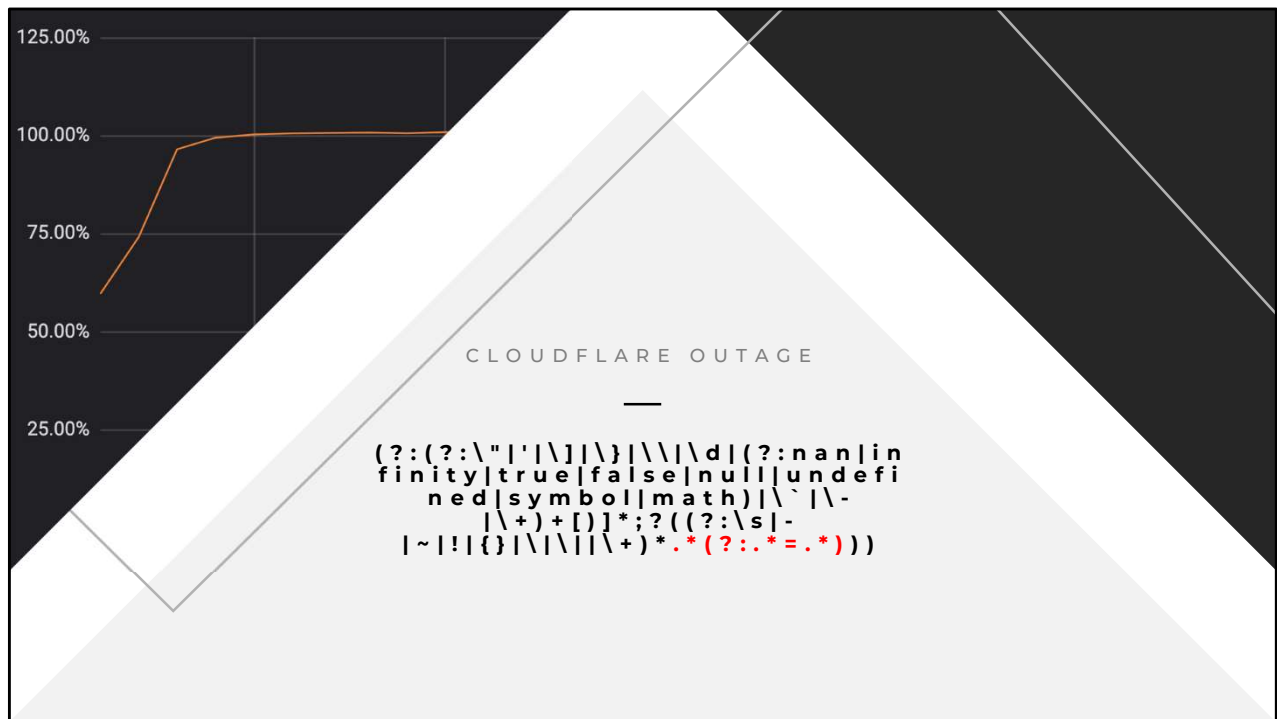
```
1. perl5.18

Starting regex match
|
V
/.*.*=.*;/

|
V
'X=XXXXXXXXXXXXXXXXXXXXX'

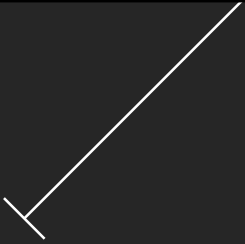
[Visual of regex at 're.pl' line 6]      [step: 0]
```

HAPPY SOUND 23




We can see that the regex they deployed had the hallmarks of a regex which allows catastrophic backtracking.

[illegible]



SOLUTION TO BACK- TRACKING

HAPPY SOUND



GET GOOD


Now knowing what commonly causes problems with RegEx – avoid them.

- Make regex as specific as you can.
- Avoid repeating groups.
- Avoid . as much as possible.


NON-BACKTRACKING REGEX ENGINE

google/re2

RE2 is a fast, safe, thread-friendly alternative to backtracking regular expression engines like those used in PCRE, Perl, and Python....



44	0	6k	860
Contributors	Issues	Stars	Forks



27

There are two simple solutions to avoiding catastrophic regex backtracking.

The first one is to simply git gud. I've explained the problem. Avoid repetitive matches. Try to avoid . and especially .* as much as possible.

Second: check out the google re2 engine. This engine doesn't backtrack. What this means, however, is you cannot do back-references and can be slower for capturing operations. However, because RE2 does not backtrack it guarantees linear run-time, not exponential as we've seen.

Step 0. a b a b
a b b b

Step 1. a b a b
a b b b

Step 2. a b a b
a b b b

Step 3. a b a b
a b b b

Step 4. a b a b
a b b b

Step 5. a b a b
a b b b

Step 6. a b a b
a b b b

Step 7. a b a b
a b b b

Step 8. a b a b
a b b b

Fails, backtracks

**tl;ncs:
RE2**

HAPPY SOUND

Step 0. a b a b
a b b b

Step 1. a b a b
a b b b

Step 2. a b a b
a b b b

Step 3. a b a b
a b b b

Step 4. a b a b
a b b b

Instead of proceeding through the regex until it fails,
it assesses all possible states at once.

This is where we start talking about NFA and DFA – Non-Deterministic, and Deterministic Finite Automaton.

The “Too Long; Not a Computer Scientist” of RE2 is that instead of proceeding through regex until it fails, then backtracking, is that it will assess all possible states at once.

Why isn’t this used more often? Well – most of the time, matching in a traditional NFA engine (such as PCRE) is simply ‘fast-enough’, and in those cases it’s usually faster than DFA.

HAPPY
SOUND

FURTHER
READING

Learn Regex: <https://regexone.com/>

Russ Cox – Regular Expression Matching can be Simple and Fast
<https://swtch.com/~rsc/regex/regexpl.html>

REDoS: <https://blog.caller.xyz/user-agent-parsing-redos-cve-2020-5243/>

DoyenSec Regexploit: <https://blog.doyensec.com/2021/03/11/regexploit.html>

Stack Overflow outage: <https://stackstatus.net/post/147710624694/outage-postmortem-july-20-2016>

Stack Overflow outage causing post:
<https://stackoverflow.com/questions/38484433/join-tiles-in-corona-sdk-into-one-word-for-a-breakout-game-grid>

30