# /(Re+|DoS)*/
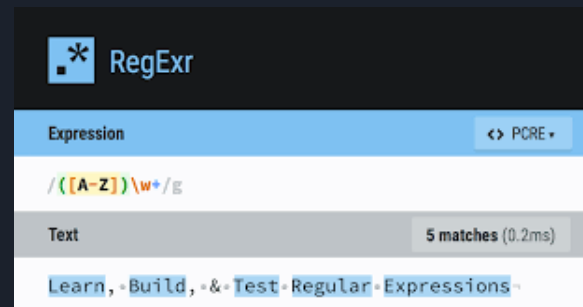
For fun and stuck!

# O que é uma expressão regular (regex)?

- Uma expressão regular (regex) é uma sequência de caracteres que especificam um padrão de busca.

- Surgiu a partir dos estudos do matemático Stephen Cole Kleene no ano de 1950.

- Definiu as bases para os primeiros algoritmos de busca.

# Pra que serve uma (expressão regular) regex?

- São utilizadas para busca, extração, substituição e validação de padrões de caracteres.

- Utilizada por softwares que realizam busca, SGDBs, IDEs, processadores de texto, WAFs, IDSs, motores de busca, parsers, utilitários de linha de comando e etc.

- MySQL, MongoDB, VS Code, grep, findstr e etc.

# Pra que serve uma (expressão regular) regex?

- Validação de IPv4

```
/\d{1,3}\.\d{1,3}\.\d{1,3}\.\d{1,3}/

\d     - Representa dígitos
{1,3}  - O caractere anterior pode aparecer 1 a 3 vezes
\.     - Equivale a um ponto final

//192.168.0.1
//127.0.0.1
//8.8.8.8
```

# Pra que serve uma (expressão regular) regex?

- Comentários HTML / XML

```
/<!--([\s\S]*?)-->/

<!-- - String literal
\s   - Espaço em branco
\S   - Não espaço em branco
*?   - Busca uma ou mais vezes de forma não gananciosa
-->  - String literal


// <!-- -->
// <!-- #HackThePlanet -->
// <!-- Hello, world! -->
```
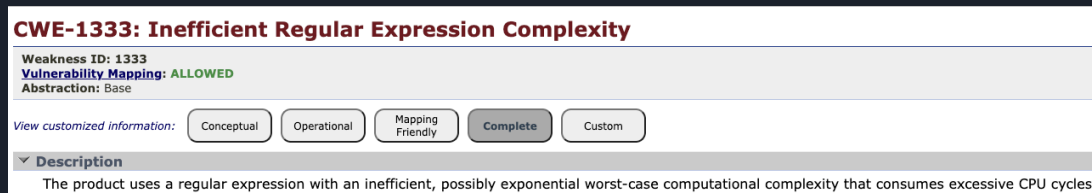
# O que é um ataque ReDoS?

- É uma abreviação de Regular Expression Denial of Service.

- Pode causar indisponibilidade de um sistema por lentidão ou por queda.

- É um ataque que ocorre quando um engine de expressão regular tentar validar uma entrada (texto) e consome muito recurso computacional para realizar a análise.



**CWE-1333: Inefficient Regular Expression Complexity**

Weakness ID: 1333
**Vulnerability Mapping: ALLOWED**
Abstraction: Base

*View customized information:*  | Conceptual | Operational | Mapping Friendly | Complete | Custom |

▽ **Description**

The product uses a regular expression with an inefficient, possibly exponential worst-case computational complexity that consumes excessive CPU cycles.

# Porque um ReDoS ocorre?

- Alguns algoritmos de análise de expressão regular implementam backtracking*.

- Dependendo do tamanho do input, esse algoritmo com  backtracking é sobrecarregado e consome muito processamento para analisar um regex.

- Expressões regulares ambíguas "triggam" esse cenário.

- Linguagens que possuem engine que  implementam backtracking são vulneráveis. (C#, JavaScript, Java, Python, PHP, Ruby, etc.)

*Algoritmo de busca profunda (Depth-First Search)

# Show me the code



Talk is cheap. Show me the code.

(Linus Torvalds)

izquotes.com

# Incidentes e vulnerabilidades de ReDoS

**CloudFlare (**02/07/19**)** - DoS de domínios devido a uma regex em uma regra de WAF.

**StackOverflow (**20/07/2016**)** - DoS devido a uma regex que removia espaços no início e ao final do texto.
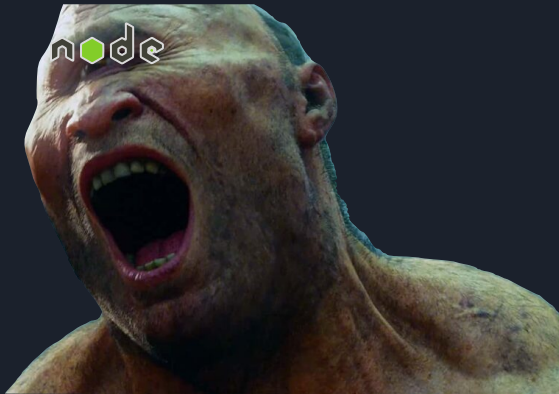
**validator.js (**CVE-2021-3765**)** - Já teve várias vulnerabilidades de ReDoS e é utilizada por mais de 5600 bibliotecas.

**jsPDF** (CVE-2021-23353) - Vulnerabilidade de ReDoS ao validar uma imagem.

# Falhas ReDoS não param de crescer

**ReDoS vulnerabilities in npm spikes by 143% and XSS continues to grow**
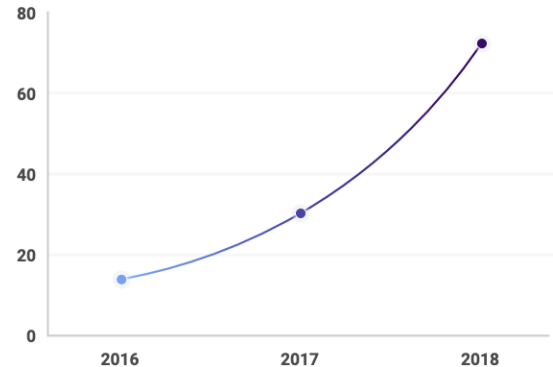
Written by: Liran Tal

## Regular expression denial of service

The Node.js runtime is known to have many strengths, but one of them, the single threaded Event Loop, can also be its weakest link if not used correctly. This happens more regularly than one might think.

Regular expression denial of service (ReDoS) attacks exploit the non-linear worst-case complexity vulnerabilities that some regex patterns can lead to. For a single-threaded runtime this could be devastating, and this is why Node.js is significantly affected by this type of vulnerability.

We found that there were a growing number of ReDoS vulnerabilities disclosed over the last three years, with a spike of 143% in 2018 alone.

**Regular expression denial of service (ReDoS) disclosures on the rise**

snyk

# O formato CSV

- É um formato de texto que utiliza vírgulas para separar valores

- É um formato de arquivo level e simples para representar dados em forma de linhas e colunas

- Muito usado em softwares como Microsoft Excel e Libre Office

# O formato CSV

```
Nome,Idade,Cidade,Hobbies
João,25,São Paulo,"Leitura,Esportes,Jogos de Tabuleiro"
Maria,30,Rio de Janeiro,"Música,Artes,Culinária"
Carlos,22,Belo Horizonte,"Viajar,Fotografia,Esportes Aquáticos"
```

| Nome | Idade | Cidade | Hobbies |
|------|-------|--------|---------|
| João | 25 | São Paulo | Leitura,Esportes,Jogos de Tabuleiro |
| Maria | 30 | Rio de Janeiro | Música,Artes,Culinária |
| Carlos | 22 | Belo Horizonte | Viajar,Fotografia,Esportes Aquáticos |

# Google e sites estilo Stack Overflow

# Google e sites estilo Stack Overflow

| Regex vulneráveis | Regex não vulneráveis |
|---|---|
| 5 | 3 |

# Chat GPT



**You**

I need a regex to get every comma that's outside a "", look the example

1,0,3,"Braund, Mr. Owen Harris", male,22,1,0,A/5 21171,7.25,,S

**ChatGPT**

You can use the following regular expression to match commas outside double quotes:

```regex
,(?=(?:[^"]*"[^"]*")*[^"]*$)
```

Here's how this regex works:

- `,` matches the comma.
- `(?= ...)` is a positive lookahead assertion, which checks if the following pattern is present without including it in the match.
- `(?:[^"]*"[^"]*")*` matches pairs of double quotes and anything in between them. This allows us to skip commas inside double quotes.
- `[^"]*$` matches anything that is not a double quote until the end of the string.

**You**

This regex is vulnerable to redos

```
,(?=(?:[^"]*"[^"]*")*[^"]*$)
```

**ChatGPT**

You're correct; the provided regex can potentially be vulnerable to catastrophic backtracking, which is a security concern. To address this, you can make a slight modification to make it more efficient and avoid potential performance issues.

Here's an updated regex that aims to be more efficient:

```regex
,(?=(?:(?:(?:"[^"]*")*[^"]*)*"[^"]*"[^"]*)*[^"]*$)
```

This modification introduces a non-capturing group `(?:[^"]*"[^"]*"[^"]*)*` to match sequences of characters within double quotes more efficiently, reducing the risk of catastrophic backtracking. The overall structure of the regex remains the same, ensuring commas are matched outside these sequences of double quotes.

Always be cautious when dealing with regular expressions, especially in scenarios where user input is involved, to prevent potential security vulnerabilities like Regular Expression Denial of Service (ReDoS).

# Chat GPT



**You**
This regex is vulnerable to redos

,(?=(?:[^"]*"[^"]*")*[^"]*$)

**ChatGPT**
I apologize for the oversight. You are correct, and I appreciate your diligence in pointing that out. To avoid potential ReDoS vulnerabilities, we can use a different approach:

```regex
regex                                    Copy code

"[^"]*"|[^,]+
```

This regex will match either a quoted string or a sequence of characters that does not include a comma. You can use this regex to extract values from your CSV data while avoiding the ReDoS issue.

# Chat GPT

| Regex vulneráveis | Regex não vulneráveis |
|---|---|
| 2 | 1 |

# Github Copilot

I need a regex that parse this CSV 1,0,3,"Braund, Mr. Owen Harris", male,22,1,0,A/5 21171,7.25,,S

Accept  Discard  ⟳  Changed 4 lines 👍 👎

```
// BEGIN: regex
const regex = /(?:^|,)(?:(?:"([^"]*(?:""[^"]*)*)")|([^",]*))(?=,|$)/g;
// END: regex
```

I need a regex that parse this CSV 1,0,3,"Braund, Mr. Owen Harris", male,22,1,0,A/5 21171,7.25,,S

Accept  Discard  ⟳  Changed 7 lines 👍 👎

```
const csvString = '1,0,3,"Braund, Mr. Owen Harris", male,22,1,0,A/5 21171,7.25,,S';
const regexPattern = /,(?=(?:(?:[^"]*")*[^"]*$)/;

const csvArray = csvString.split(regexPattern);
console.log(csvArray);
```

I need a regex that parse this CSV 1,0,3,"Braund, Mr. Owen Harris", male,22,1,0,A/5 21171,7.25,,S

Accept  Discard  ⟳  Changed 7 lines 👍 👎

```
// BEGIN: Regex pattern
const csvString = '1,0,3,"Braund, Mr. Owen Harris", male,22,1,0,A/5 21171,7.25,,S';
const regex = /(?:^|,)(\"(?:[^\"]+|\"\")*\"|[^,]*)/g;
const result = csvString.match(regex).map((value) => value.replace(/^,/, ''));
console.log(result);
// END: Regex pattern
```

I need a regex that parse this CSV 1,0,3,"Braund, Mr. Owen Harris", male,22,1,0,A/5 21171,7.25,,S

Accept  Discard  ⟳  Changed 7 lines 👍 👎

```
// BEGIN: Regex to split CSV string
const csvString = '1,0,3,"Braund, Mr. Owen Harris", male,22,1,0,A/5 21171,7.25,,S';
const regex = /,(?=(?:(?:[^"]*")*[^"]*$)/;
const result = csvString.split(regex);
console.log(result);
// END: Regex to split CSV string
```

# Github Copilot



I need a regex that parse this CSV 1,0,3,"Braund, Mr. Owen Harris", male,22,1,0,A/5 21171,7.25,,S

**Accept**  Discard ⌄  ↻                                    Changed 7 lines 👍 👎

```
const csvString = '1,0,3,"Braund, Mr. Owen Harris", male,22,1,0,A/5 21171,7.25,,S';
const regex = /("[^"]+"|[^,]+)/g;
const result = csvString.match(regex);

console.log(result);
```

I need a regex that parse this CSV 1,0,3,"Braund, Mr. Owen Harris", male,22,1,0,A/5 21171,7.25,,S

**Accept**  Discard ⌄  ↻                                    Changed 12 lines 👍 👎

```
const csvString = '1,0,3,"Braund, Mr. Owen Harris", male,22,1,0,A/5 21171,7.25,,S';
const regexPattern = /(?:^|,)(?:"([^"]*(?:""[^"]*)*)"|([^",]*))/g;
const parsedArray = [];
let match;

while ((match = regexPattern.exec(csvString))) {
  parsedArray.push(match[1] || match[2]);
}

console.log(parsedArray);
```

I need a regex that parse this CSV 1,0,3,"Braund, Mr. Owen Harris", male,22,1,0,A/5 21171,7.25,,S

**Accept**  Discard ⌄  ↻                                    Changed 7 lines 👍 👎

```
// BEGIN: regex
const csvString = '1,0,3,"Braund, Mr. Owen Harris", male,22,1,0,A/5 21171,7.25,,S';
const regexPattern = /("[^"]*"|[^,]+)(?:,|$)/g;
const csvArray = csvString.match(regexPattern).map((entry) => entry.replace(/^"(.*)"$/, '$1'));
console.log(csvArray);
// END: regex
```
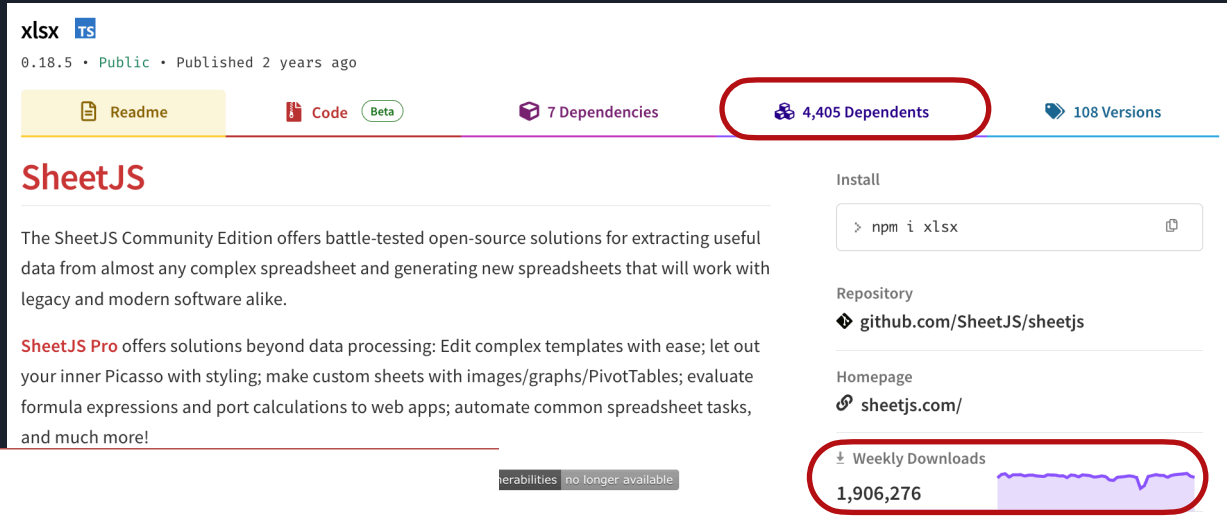
# Github Copilot

| Regex vulneráveis | Regex não vulneráveis |
|:---:|:---:|
| 4 | 3 |

# Google, Chat GPT e GitHub Copilot

| Regex vulneráveis | Regex não vulneráveis |
|:-:|:-:|
| 11 | 7 |

# Minha primeira CVE 🎉



xlsx TS

0.18.5 • Public • Published 2 years ago

📄 Readme | 🔴 Code Beta | 📦 7 Dependencies | **4,405 Dependents** | 🏷 108 Versions

## SheetJS

The SheetJS Community Edition offers battle-tested open-source solutions for extracting useful data from almost any complex spreadsheet and generating new spreadsheets that will work with legacy and modern software alike.

**SheetJS Pro** offers solutions beyond data processing: Edit complex templates with ease; let out your inner Picasso with styling; make custom sheets with images/graphs/PivotTables; evaluate formula expressions and port calculations to web apps; automate common spreadsheet tasks; and much more!

**Install**

> `npm i xlsx`

**Repository**
🔷 github.com/SheetJS/sheetjs

**Homepage**
🔗 sheetjs.com/

⬇ **Weekly Downloads**
**1,906,276**

nerabilities no longer available

---

## CVE-2024-22363

### Summary

All versions of SheetJS CE through 0.20.1 are vulnerable to "Regular Expression Denial of Service" (ReDoS). For more details, see https://regexide.com

### Categorization

CVSS Vector: CVSS:3.1/AV:N/AC:L/PR:N/UI:N/S:U/C:N/I:N/A:H **(Score 7.5 - High)**

CWE-1333 Inefficient Regular Expression Complexity [1]

# Minha primeira CVE 🎉

## Why the move away from npm registry? #2667

**Closed**   JSin opened this issue on Apr 26, 2022 · 24 comments

**JSin** commented on Apr 26, 2022

I noticed from this `ed18acd` that you moved away from publishing to the npm registry and recommend people download using tarballs on the CDN. Why did you move away? The npm registry is an extremely common way to download packages.

👍 29   👀 2

**SheetJSDev** commented on Apr 27, 2022   Contributor

Hi, sheetjs!

On May 31, 2022, your account will be enrolled in mandatory 2FA as you are a maintainer of xlsx, one of the Top-500 packages by dependents on the npm registry. Late last year, we announced our intent to enforce 2FA for all high-impact maintainers of package on the npm registry and on February 1, 2022, we began rolling out mandatory 2FA to our first cohort.

**How will this affect my account?**

The package on the npm registry is surprisingly popular (one of the top-500 by dependents). npm invalidated the old publish token and is forcing 2FA on the publishing account.

Due to ongoing legal matters between SheetJS LLC and npm, Inc. (which will not be discussed here), it did not make sense to continue using the public npm registry for distribution.

With GitHub (the owner of NPM) sunsetting the git.io domain with only 4 days notice, we are reminded of the ephemerality of the Internet and the inherent risks of relying on platforms.

---

## Remediation

Users should upgrade to version 0.20.2 or later. Official releases are available on the SheetJS CDN [2]. SheetJS CE documentation includes installation instructions for common deployments [7].

## Acknowledgements

Special thanks to Asadbek Karimov and Jardel Matias for discovering and helping to resolve the issue.

# Minha quase segunda CVE

**nodemailer** DT

6.9.13 • Public • Published 17 days ago

📄 Readme    🦀 Code (Beta)    📦 0 Dependencies    **6,911 Dependents**    🏷 269 Versions

## Nodemailer

**Install**

```
> npm i nodemailer
```

**Repository**
◈ github.com/nodemailer/nodemailer

**Homepage**
🔗 nodemailer.com/

⬇ **Weekly Downloads**
**3,150,937**

---

## nodemailer ReDoS when trying to send a specially crafted email

**Moderate severity**   GitHub Reviewed   Published on Jan 31 in **nodemailer/nodemailer** • Updated on Feb 1

**Vulnerability details**    Dependabot alerts 0

| Package | Affected versions | Patched versions |
|---|---|---|
| ⓝ **nodemailer** (npm) | <= 6.9.8 | 6.9.9 |

**Severity**
Moderate   5.3 / 10

### Description

#### Summary

A ReDoS vulnerability occurs when nodemailer tries to parse img files with the parameter `attachDataUrls` set, causing the stuck of event loop.
Another flaw was found when nodemailer tries to parse an attachments with a embedded file, causing the stuck of event loop.

#### Details

Regex: /^data:((?:[^;];)(?:[^;])),(.)$/

Path: compile -> getAttachments -> _processDataUrl

Regex: /(<img\b[^>]* src\s*=[\s"']*)(data:([^;]+);[^"'>\s]+)/

Path: _convertDataImages

#### PoC

https://gist.github.com/francoatmega/890dd5053375333e40c6fdbcc8c58df6
https://gist.github.com/francoatmega/9aab042b0b24968d7b7039818e8b2698

**CVSS base metrics**

| | |
|---|---|
| Attack vector | Network |
| Attack complexity | Low |
| Privileges required | None |
| User interaction | None |
| Scope | Unchanged |
| Confidentiality | None |
| Integrity | None |
| Availability | Low |

CVSS:3.1/AV:N/AC:L/PR:N/UI:N/S:U/C:N/I:N/A:L

**Weaknesses**
CWE-1333

**CVE ID**
No known CVE

**GHSA ID**
GHSA-9h6g-pr28-7cqp

# Fun facts

Your reported issue stems from the regular expression /<!--([\s\S]*?)-->/mg . This general pattern of a non-greedy match of HEADER + any + FOOTER is common (and is easily "pumped" by repeating HEADER). There are a number of workarounds. In this case it would suffice to ensure that FOOTER exists:

if(str.indexOf("-->") > -1) str = str.replace(/<!--([\s\S]*?)-->/mg,"")

For good measure, we will go back and review some of the regular expressions with an eye towards different options. That replace literally looks for a closing `-->` after an opening `<!--`, which is probably better implemented as a loop that scans for HEADER and FOOTER using `indexOf`

Jardel Matias
Para: SheetJS Security Team

☼   ← Responder     ← Responder a todos     → Encaminhar     [

Seg, 08/01/2

evilFile.xlsx
1    <!--><!--<!--<!--<!--<!
     !--<!--<!--<!--<!--<!--<
     <!--<!--<!--<!--<!--<!--
     !--<!--<!--<!--<!--<!--<
     <!--<!--<!--<!--<!--<!--
     -<!--<!--<!--<!--<!--<!-
     <!--<!--<!--<!--<!--<!--
     -<!--<!--<!--<!--<!--<!-
     !--<!--<!--<!--<!--<!--<
     -<!--<!--<!--<!--<!--<!-

Hi again, just thinking a little bit more but I think that this line does not solve the problem

if(str.indexOf("-->") > -1) str = str.replace(/<!--([\s\S]*?)-->/mg,"")

The conditional there can be easily bypassed via a crafted file that start like this *<!--><!--<!--* and still trigger the ReDos.

Best regards.
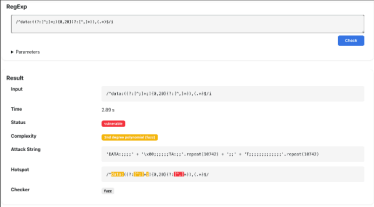
# Fun facts



Jardel Matias
Para: andris@reinman.eu

Ter, 20/02/2024 02:24

I'm sorry about that, I don't like being an annoying guy, but the fix for the ReDoS flaw that I mentioned still persists, the main point is:

I recommend you this link to check vulnerable regex, it is not a 100% guarantee, but during my tests it almost never failed.

Best regards,

Jardel Matias

## v6.9.10

### 6.9.10 (2024-02-22)

#### Bug Fixes

- **data-uri:** Do not use regular expressions for parsing data URI schemes (12e65e9)
- **data-uri:** Moved all data-uri regexes to use the non-regex parseDataUri method (edd5dfe)

▸ **Assets** 2

🙂 👍 5 You and 4 others reacted

# Fun facts

```
473    –                let parts = (content.path || content.href).match(/^data:((?:[^;]*;)*)(?:[^,]*)),(.*)$/i);
474    –                if (!parts) {
       522  +           let parsedDataUri = module.exports.parseDataURI(content.path || content.href);
       523  +
       524  +           if (!parsedDataUri || !parsedDataUri.data) {
```

```
        parseDataURI(uri) {
            let input = uri;
565  +      let commaPos = input.indexOf(',');
566  +      if (!commaPos) {
567  +          return uri;
568  +      }
569  +
570  +      let data = input.substring(commaPos + 1);
571  +      let metaStr = input.substring('data:'.length, commaPos);
572  +
573  +      let encoding;
574  +
575  +      let metaEntries = metaStr.split(';');
576  +      let lastMetaEntry = metaEntries.length > 1 ? metaEntries[metaEntries.length - 1] : false;
577  +      if (lastMetaEntry && lastMetaEntry.indexOf('=') < 0) {
578  +          encoding = lastMetaEntry.toLowerCase();
579  +          metaEntries.pop();
580  +      }
581  +
582  +      let contentType = metaEntries.shift() || 'application/octet-stream';
583  +      let params = {};
584  +      for (let entry of metaEntries) {
585  +          let sep = entry.indexOf('=');
586  +          if (sep >= 0) {
587  +              let key = entry.substring(0, sep);
588  +              let value = entry.substring(sep + 1);
589  +              params[key] = value;
```

# Fun facts

```
421  + module.exports.parseDataURI = uri => {
422  +     let input = uri;
423  +     let commaPos = input.indexOf(',');
424  +     if (!commaPos) {
425  +         return uri;
426  +     }
427  +
428  +     let data = input.substring(commaPos + 1);
429  +     let metaStr = input.substring('data:'.length, commaPos);
430  +
431  +     let encoding;
432  +
433  +     let metaEntries = metaStr.split(';');
434  +     let lastMetaEntry = metaEntries.length > 1 ? metaEntries[metaEntries.length - 1] : false;
435  +     if (lastMetaEntry && lastMetaEntry.indexOf('=') < 0) {
436  +         encoding = lastMetaEntry.toLowerCase();
437  +         metaEntries.pop();
438  +     }
439  +
440  +     let contentType = metaEntries.shift() || 'application/octet-stream';
441  +     let params = {};
442  +     for (let entry of metaEntries) {
443  +         let sep = entry.indexOf('=');
444  +         if (sep >= 0) {
445  +             let key = entry.substring(0, sep);
446  +             let value = entry.substring(sep + 1);
447  +             params[key] = value;
448  +         }
```

```
537  537          * @return {Object} Parsed element
538  538          */
539  539         _processDataUrl(element) {
540      -            let parts = (element.path || element.href).match(/^data:((?:[^;]*;){0,20}(?:[^,]*)),(.*)$/i);
541      -            if (!parts) {
     540  +            let parsedDataUri;
     541  +            if ((element.path || element.href).match(/^data:/)) {
     542  +                parsedDataUri = this.parseDataURI(element.path || element.href);
     543  +            }
     544  +
     545  +            if (!parsedDataUri) {
```

# Como prevenir ataques ReDoS?

- Utilize bibliotecas de terceiros.

- Evite criar regex que possam levar a resultados ambíguos, evite agrupamentos com repetições ou alternância. Como por exemplo:

  `^data:(\w*\/\w*);*(charset=[\w=-]*)*;*$`

  `^([a-zA-Z0-9])(([\-.]|[_]+)?([a-zA-Z0-9]+))*(@){1}[a-z0-9]+[.]{1}(([a-z]{2,3})|([a-z]{2,3}[.]{1}[a-z]{2,3}))$`

  `^-?\D*[\d.]+,\d{2}$`

  `^(a|a)*$`

- Utilize engines de regex que não sejam vulneráveis a backtracing, como por exemplo a re2 da Google.*

- Mantenham suas dependências atualizadas.

*Linguagens como Go e Rust não são vulneráveis.

# Como prevenir ataques ReDoS?

- Valide as suas regex

- https://devina.io/redos-checker

- Projeto **recheck**

- Biblioteca JS, biblioteca Scala e Plugin ESLint

- https://makenowjust-labs.github.io/recheck/





recheck

The trustworthy ReDoS checker

# Como prevenir ataques ReDoS?

- Utilize as funcionalidades da linguagem que você utiliza

- C# (Backtrack abortion)

## Regex.MatchTimeout Property

Reference                                                    👍 Feedback

## Definition

Namespace:  System.Text.RegularExpressions
Assembly:  System.Text.RegularExpressions.dll
Source:  Regex.Timeout.cs ⧉

Gets the time-out interval of the current instance.

```C#
public TimeSpan MatchTimeout { get; }
```

# Como prevenir ataques ReDoS?

## Linear regex algorithm #51659

✓ Closed  francoatmega opened this issue on Feb 4 · 4 comments

francoatmega commented on Feb 4 · edited ▾

### What is the problem this feature will solve?

Recently, we have witnessed a surge in ReDoS (Regular Expression Denial of Service) issues, and the root cause behind these occurrences lies in the regex engine we employ for matching patterns. There is a library called re2, which serves as a wrapper for Google's re2 engine. This library resolves all regex patterns in linear time, albeit lacking support for backreferences and look-arounds.

Upon researching this matter, it becomes evident that prominent projects are grappling with this problem. Languages such as Rust and Go have successfully mitigated these issues long ago. Considering the impact of ReDoS on Node.js, particularly its tendency to stall the event loop, it seems prudent to integrate this functionality into the Node.js codebase.

For further insights into the severity of ReDoS vulnerabilities, you may refer to the following resources:

https://snyk.io/blog/redos-vulnerabilities-in-npm-spikes-by-143-and-xss-continues-to-grow/
https://swtch.com/~rsc/regexp/regexp1.html

rluvaton commented on Feb 5                                                          Member  ⋯

Since RegExp is part of ECMAScript and therefore the JS engine should implement this and not NodeJS, you can open issue in the V8 issue tracking

😊  👍 1

# Como prevenir ataques ReDoS?

**Aurele-Barriere** commented on Feb 6

Have you looked into using the "Experimental" regex engine in V8, available in Node?

It has been described here:
https://v8.dev/blog/non-backtracking-regexp
It uses exactly the Thompson NFA simulation algorithm described in the blog post you linked.

You can start node with `node --enable-experimental-regexp-engine`.
Then, a new Regex flag is available, `l`.
Any regex with the `l` flag will be executed using the linear Experimental engine.

For instance, if I define `const re = /(a*)*b/`,
Using the usual backtracking algorithm on a string of "a"s (for instance `"a".repeat(1000).match(re)`), it takes forever.
But if I define instead `const re = /(a*)*b/l` with the `l` flag, then `"a".repeat(1000).match(re)` finishes almost immediately.

As you said, this engine lacks support for backreferences and lookarounds.
However, a recent commit adds support for lookbehinds without capture groups inside:
https://chromium-review.googlesource.com/c/v8/v8/+/5093860

Is this Experimental engine what you are looking for?

👍 1    🎉 1

# Como prevenir ataques ReDoS?



**Claudia** <impinball@gmail.com>                                    qua., 7 de fev., 13:32
para v8-dev ▾

You may be interested in the following links:

- https://v8.dev/blog/non-backtracking-regexp
- https://es.discourse.group/t/safe-regex-engine-to-prevent-redos-attack/450/6
- https://es.discourse.group/t/re2-consider-having-it-as-alternative-engine-choice/1388/15

V8 does impose a bound currently on some regexps, but not all.

On Tuesday, February 6, 2024 at 4:59:04 AM UTC-8 jardelm...@gmail.com wrote:

Recently, we have witnessed a surge in ReDoS (Regular Expression Denial of Service) issues, and the root cause behind these occurrences lies in the regex engine we employ for matching patterns. There is a library called re2, which serves as a wrapper for Google's re2 engine. This library resolves all regex patterns in linear time, albeit lacking support for backreferences and look-arounds.

Upon researching this matter, it becomes evident that prominent projects are grappling with this problem. Languages such as Rust and Go have successfully mitigated these issues long ago. Considering the impact of ReDoS on Node.js, particularly its tendency to stall the event loop, it seems prudent to integrate this functionality into the Node.js codebase.

For further insights into the severity of ReDoS vulnerabilities, you may refer to the following resources:

https://snyk.io/blog/redos-vulnerabilities-in-npm-spikes-by-143-and-xss-continues-to-grow/
https://swtch.com/~rsc/regexp/regexp1.html

To solve this problem maybe we can implementing a regex algorithm that operates in linear time, such as the Thompson implementation detailed in this link, is a worthwhile endeavor.

A had open the issues in Node Js issues page, take a look https://github.com/nodejs/node/issues/51659

Best regards,

Jardel Matias

# Como prevenir ataques ReDoS?

# Como prevenir ataques ReDoS?

- Habilitar a feature flag **--enable-experimental-regexp-engine**

- Adicionar a flag *l* nas regex

- Be happy ;)

# Como prevenir ataques ReDoS?



Talk is cheap. Show me the code.

(Linus Torvalds)

izquotes.com

# Referências

https://owasp.org/www-community/attacks/Regular_expression_Denial_of_Service_-_ReDoS

https://en.wikipedia.org/wiki/Regular_expression

https://www.regular-expressions.info/catastrophic.html

https://lirantal.medium.com/node-js-pitfalls-how-a-regex-can-bring-your-system-down-cbf1dc6c4e02

https://levelup.gitconnected.com/the-regular-expression-denial-of-service-redos-cheat-sheet-a78d0ed7d865

# Dúvidas

# Post sobre ReDos

## ReDoS for fun and stuck

francoatmega
5 min read · Oct 22, 2023

👏 52

### What is a regular expression?

A regular expression (regex) is a sequence of characters that specify a search pattern. It is defined as a formal language and originated from the studies of the mathematician Stephen Cole Kleene in the 50s.

### Where are regular expressions used?

Regular expressions specify patterns, so they are widely used for data validation and extraction. Below is an example of a regex that identifies whether a given text is a IP address:

\d{1,3}.\d{1,3}.\d{1,3}.\d{1,3}

# Contato

# Contato



Jardel Matias ✓ (He/Him)
Software Engineer | InfoSec Enthusiast | App Sec
Juazeiro do Norte, Ceará, Brazil · Contact info
802 followers · 500+ connections