# Introduction to pattern matching

Using Regexes

**CIRCL**
Computer Incident
Response Center
Luxembourg

Alexandre Dulaunoy
alexandre.dulaunoy@circl.lu

Jean-Louis Huynen
jean-louis.huynen@circl.lu

Sami Mokaddem
sami.mokaddem@circl.lu

info@circl.lu

October 7, 2022

## Typical Linux problem

```
1  $ cat files.txt
2  readme.md
3  document.pdf
4  image.png
5  music.mp3
6  video.mp4
7  manual.pdf
```

**Objectives**: List only PDF files

## $ man grep

```
1  GREP(1)                 User Commands                 GREP(1)
2
3  NAME
4       grep, egrep, fgrep, rgrep - print lines that match
       patterns
5
6  SYNOPSIS
7       grep [OPTION...] PATTERNS [FILE...]
8       grep [OPTION...] -e PATTERNS ... [FILE...]
9       grep [OPTION...] -f PATTERN_FILE ... [FILE...]
10
11 DESCRIPTION
12      grep searches for PATTERNS in each FILE.  PATTERNS
       is one or more patterns separated  by  newline
       characters,  and  grep prints each line that matches
        a pattern.  Typically PATTERNS should be quoted
       when grep is used in a shell command.
```

## Using grep

```
1 $ cat files.txt | grep 'pdf'
2 document.pdf
3 manual.pdf
```

Easy! However...

## Using grep

```
1 $ cat files.txt | grep 'pdf'
2 document.pdf
3 manual.pdf
```
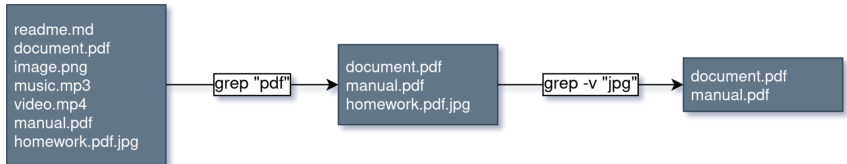
Easy! However...

```
1 $ cat files-2.txt | grep 'pdf'
2 document.pdf
3 manual.pdf
4 homework.pdf.jpg
```

How can we filter out homework.pdf.jpg?

## Using grep

```
1 $ cat files -2. txt | grep 'pdf' | grep -v 'jpg'
2 document . pdf
3 manual . pdf
```

-v allows us to perform an in**V**ert match

## Using grep

```
1 $ cat files -3.txt | grep 'pdf' | grep -v 'jpg'
2 document.pdf
3 manual.pdf
4 adobe_pdf_reader.exe
5 i_hate_pdf.mp3
6 this.is.a.weird.pdf.filename.zip
7 filename with spaces are evil.pdf
```

Using invert match is not going to scale...

## Other commonly encountered problems

- Matching valid email addresses
- Matching valid IBAN number
- Matching valid IPv4 addresses

$\rightarrow$ Regular Expressions (**Regex**) to the rescue

## Regular Expression

A **regular expression** (shortened as **regex** or **regexp**) is a sequence of characters that specifies a search pattern in text.

**Regexes** are extremely useful in extracting information from text.

## Regular Expression Basics

- What ?
  - Literal characters: `abc`
  - Quantifiers: `ab+c`
  - Operator OR: `(abc|cba)`
  - Bracket expressions: `[a-z]`
  - Meta sequences: `\S`
  - Capture group: `(abc)`
  - Anchors: `^abc$`
- New skill ?

```
/ <([a-z]+)(>(.*)<\/\1>|\s+\/>) /
```

# Regular Expression Basics: Literal characters

Letters and digits from the ASCII character set match their respective value

# Regular Expression Basics: The Dot

. is a *joker* or *wildcard* that can match any single character

## Regular Expression Basics: The Period

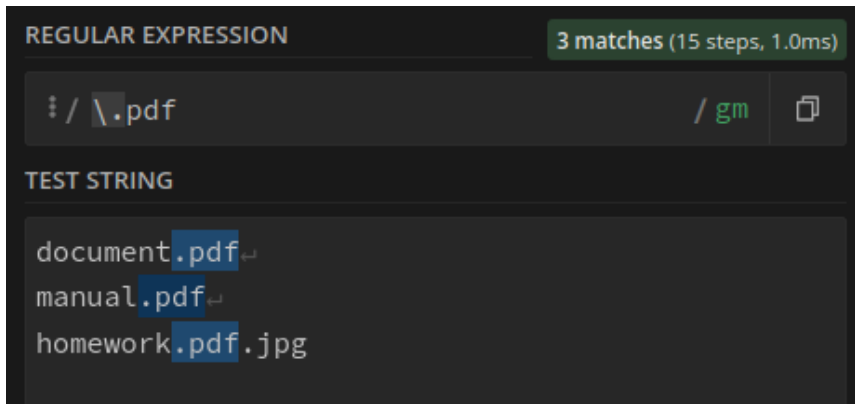The period character `.` can be matched using the escape character `\`.

# Regular Expression Basics: OR Operator

The $(\;|\;)$ structure can be used as a logical operator to match one sequence or the other

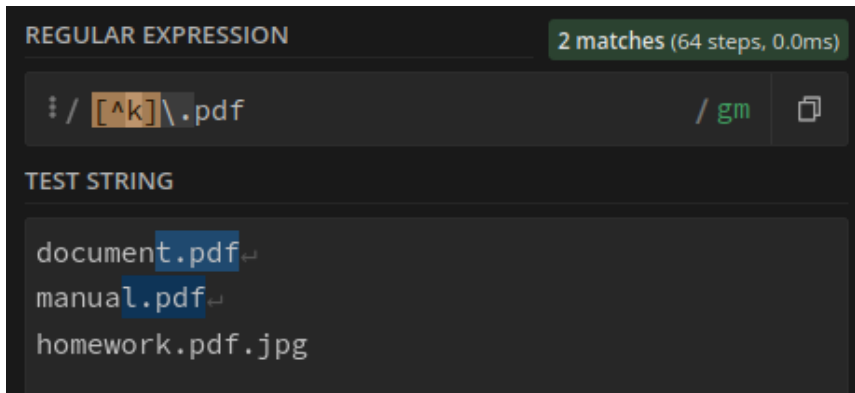# Regular Expression Basics: Bracket Expression (1)

The [ ] structure can be used to specifiy a set of characters that can match

# Regular Expression Basics: Bracket Expression (2)

The [^ ] structure can be used to exclude a specific set of characters

# Regular Expression Basics: Bracket Expression (3)

The [ - ] structure can be used to specifiy a range of sequential characters

# Regular Expression Basics: Meta Sequences (1)

- / . /
  - Any single character
- / \w /
  - Any word character
  - / [a-zA-Z0-9_] /
  - **Match**: any non-whitespace character $!-:;
- / \W /
  - Any non-word character
  - / [^a-zA-Z0-9_] /
  - **Match**: any whitespace character $!-:;

# Regular Expression Basics: Meta Sequences (2)

- `/ \d /`
  - Any digit
  - **Match**: `one:` `1` `, two:` `2` `;`
- `/ \s /`
  - Any whitespace character
  - **Match**: `any` `whitespace` `character` `$!-:;`
- `/ \S /`
  - Any non-whitespace character
  - **Match**: `any` `non-whitespace` `character` `$!-:;`

1. / facebo.k /

1. / facebo.k /
   ○ **Match**: facebook, faceboak, facebo&k
2. / 4\.2 /

1. / facebo.k /
   ○ **Match**: facebook, faceboak, facebo&k
2. / 4\.2 /
   ○ **Match**: 4.2
   ○ **Match**: Nice number:  4.2
3. / drink (beer|wine) ! /

1. / facebo.k /
   - **Match**: facebook, faceboak, facebo&k
2. / 4\.2 /
   - **Match**: 4.2
   - **Match**: Nice number:  4.2
3. / drink (beer|wine) ! /
   - **Match**: I drink beer !
   - **Match**: I drink wine !
4. / [e-h] /

# Regular Expression Basics: Reading Exercises (1)

1. `/ facebo.k /`
   ○ **Match**: `facebook`, `faceboak`, `facebo&k`

2. `/ 4\.2 /`
   ○ **Match**: `4.2`
   ○ **Match**: `Nice number:  4.2`

3. `/ drink (beer|wine) ! /`
   ○ **Match**: `I drink beer !`
   ○ **Match**: `I drink wine !`

4. `/ [e-h] /`
   ○ **Match**: `fefe`, `hehe`
   ○ **No match**: ~~haha~~

# Regular Expression Basics: Writing Exercises (1)

1. **Match**: *red_light*, *green_light* and *!=_light*

# Regular Expression Basics: Writing Exercises (1)

1. **Match**: *red_light*, *green_light* and *!=_light*

   / _light /

2. **Match**: *red_light* and *green_light* **but not** *white_light*

## Regular Expression Basics: Writing Exercises (1)

1. **Match**: *red_light*, *green_light* and *!=_light*
   / _light /
2. **Match**: *red_light* and *green_light* **but not** *white_light*
   / (red|green)_light /
3. **Match**: *\*_light* where * is any digit

# Regular Expression Basics: Writing Exercises (1)

1. **Match**: *red_light*, *green_light* and *!=_light*
   / _light /
2. **Match**: *red_light* and *green_light* **but not** *white_light*
   / (red|green)_light /
3. **Match**: *\*_light* where \* is any digit
   / [0-9]_light /
4. **Match**: *?_light* where ? is 4-letters color name

## Regular Expression Basics: Writing Exercises (1)

1. **Match**: *red_light*, *green_light* and *!=_light*
   / _light /
2. **Match**: *red_light* and *green_light* **but not** *white_light*
   / (red|green)_light /
3. **Match**: *\*_light* where * is any digit
   / [0-9]_light /
4. **Match**: *?_light* where ? is 4-letters color name
   / [a-z][a-z][a-z][a-z]_light /

# Regular Expression Basics: Writing Exercises (1)

1. **Match**: *red_light*, *green_light* and *!=_light*
   / _light /
2. **Match**: *red_light* and *green_light* **but not** *white_light*
   / (red|green)_light /
3. **Match**: *\*_light* where \* is any digit
   / [0-9]_light /
4. **Match**: *?_light* where ? is 4-letters color name
   / [a-z][a-z][a-z][a-z]_light /

**Question**: *?_light* where ? is any color between 3 and 6 letters

   We need a way to express occurences... Introducing **quantifiers**

# Regular Expression Basics: Quantifiers (1)

The * control character can be used to describe **zero or more** occurences

# Regular Expression Basics: Quantifiers (2)

The + control character can be used to describe **one or more** occurences

# Regular Expression Basics: Quantifiers (3)

- / a? /
  - Match 0 or one a character
- / a{3} /
  - Match exactly 3 a character
- / a{3,} /
  - Match 3 or more a character
- / a{3,6} /
  - Match between 3 and 6 a character

1. / colou?r /

1. `/ colou?r /`
   - **Match**: `colour` and `color`
2. `/ go*gle /`

# Regular Expression Basics: Reading Exercises (2)

1. `/ colou?r /`
   - **Match**: colour and color
2. `/ go*gle /`
   - **Match**: gogle, gooooogle, ggle, ...
3. `/ waz+up /`

1. / colou?r /
   - **Match**: colour and color
2. / go*gle /
   - **Match**: gogle, gooooogle, ggle, ...
3. / waz+up /
   - **Match**: wazup, wazzzzzup, ...
4. / +352[0-9]{6,8} /

1. / colou?r /
   - **Match**: colour and color
2. / go*gle /
   - **Match**: gogle, gooooogle, ggle, ...
3. / waz+up /
   - **Match**: wazup, wazzzzzup, ...
4. / +352[0-9]{6,8} /
   - **Match**: +352791648, +35226791349

# Regular Expression Basics: Writing Exercises (2)

1. **Match**: The time (16:42, 03:59)

# Regular Expression Basics: Writing Exercises (2)

1. **Match**: The time (16:42, 03:59)
   / [0-9][0-9]:[0-9][0-9] /
   (not perfect but good enough for the exercise)
2. **Match**: Luxembourg postal code (L-4253, L-1110)

# Regular Expression Basics: Writing Exercises (2)

1. **Match**: The time (16:42, 03:59)
   / [0-9][0-9]:[0-9][0-9] /
   (not perfect but good enough for the exercise)
2. **Match**: Luxembourg postal code (L-4253, L-1110)
   / L-[0-9]{4} /
3. **Match**: *_light where * is any color?

# Regular Expression Basics: Writing Exercises (2)

1. **Match**: The time (16:42, 03:59)
   / [0-9][0-9]:[0-9][0-9] /
   (not perfect but good enough for the exercise)
2. **Match**: Luxembourg postal code (L-4253, L-1110)
   / L-[0-9]{4} /
3. **Match**: *_light where * is any color?
   / [a-z]+_light /
4. **Match**: any hexadecimnal color (#ff0000, #f7f8f9)

# Regular Expression Basics: Writing Exercises (2)

1. **Match**: The time (16:42, 03:59)
   / [0-9][0-9]:[0-9][0-9] /
   (not perfect but good enough for the exercise)
2. **Match**: Luxembourg postal code (L-4253, L-1110)
   / L-[0-9]{4} /
3. **Match**: *_light* where * is any color?
   / [a-z]+_light /
4. **Match**: any hexadecimnal color (#ff0000, #f7f8f9)
   / #[a-f0-9]{6} /

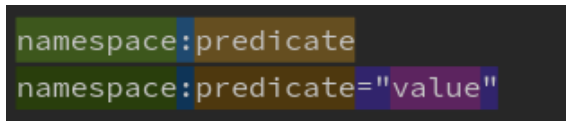Create a Regex that validates the following tags:

```
1 namespace : predicate
2 namespace : predicate = " value "
3 name _ space : pred _ icate = " value "
4 namespace : predicate = " qwert  _ + $ -  yuiop "
5 namespace : predicate = " qwert = : yuiop "
```

But not these:

```
1 tag
2 name  space : pred  icate = " value "
3 name - space : predicate = " value "
4 namespace : predicate = " qwert " yuiop "
```

# Regular Expressions: Tag matching (2)

A valid tag is composed of 2 or 3 parts



1. The namespace
2. The predicate
3. The optional value

Tag validator

```
/ ^([\w]+):([\w]+)(="([^\n"]+)")?$ /
```

# Regular Expressions: Tag matching (4)

▼ **3rd Capturing Group** `(="([^\n"]+)")?`

? matches the previous token between zero and one times, as many times as possible, giving back as needed (greedy)

  ▶ =" matches the characters =" literally (case sensitive)
  ▶ **4th Capturing Group** `([^\n"]+)`
  " matches the character " with index $34_{10}$ ($22_{16}$ or $42_8$) literally (case sensitive)

▼ **4th Capturing Group** `([^\n"]+)`
  ▼ **Match a single character not present in the list below** `[^\n"]`
    + matches the previous token between one and unlimited times, as many times as possible, giving back as needed (greedy)
    \n matches a line-feed (newline) character (ASCII 10)
    " matches the character " with index $34_{10}$ ($22_{16}$ or $42_8$) literally (case sensitive)

# Regular Expressions: Final question

What does these regexes do?

1. / ([12]\d{3}-(0[1-9]|1[0-2])-(0[1-9]|[12]\d|3[01])) /

2. / <([a-z]+)(>(.*)<\/\1>|\s+\/>) /
   - \1 is used to reference the first capturing group
   - First capturing group is ([a-z]+)

# Regexes: Going further

- ˆ and **$** anchors
- Capture **groups**
- **Greedy** and **Lazy** quantifiers
- **Possessive** quantifier

# Fun with Regular Expressions

https://regexcrossword.com/

# Fun with Regular Expressions

`https://regexcrossword.com/`

# Fun with Regular Expressions

https://regexcrossword.com/



|  | (FI\|A)+ | (YE\|OT)K | (.)[IF]+ | [NODE]+ | (FY\|F\|RG)+ |
|---|---|---|---|---|---|
| (Y\|F)(.)\2[DAF]\1 |  |  |  |  |  |
| (U\|O\|I)*T[FRO]+ |  |  |  |  |  |
| [KANE]*[GIN]* |  |  |  |  |  |

# Fun with Regular Expressions

https://regexcrossword.com/



|  | (FI\|A)+ | (YE\|OT)K | (.)[IF]+ | [NODE]+ | (FY\|F\|RG)+ |
|---|---|---|---|---|---|
| (Y\|F)(.)\2[DAF]\1 | F | O | O | D | F |
| (U\|O\|I)*T[FRO]+ | I | T | F | O | R |
| [KANE]*[GIN]* | A | K | I | N | G |