

# Regex Programming Porn

Les regex comme vous ne avez jamais vus

Human Talks Lyon, 08/10/13 / @gautier\_difolco

*regular expressions is #programming #porn, you  
always feel guilty and dirty after that :'(*

Deux solutions :

- Plus de pr0n
- Plus de Regex

# Notes

- Présentation des regex PCRE (et non POSIX)
- Nous parlerons de regex pour éviter le débat expressions régulières vs expressions rationnelles
- Présentation et environnement de benchmarks disponibles sur [https://github.com/blackheaven/regex\\_programming\\_porn](https://github.com/blackheaven/regex_programming_porn)

# Plan

- Constitutifs de base
- Changements du comportement d'évaluation
- Un exemple qui fait saigner du nez

# Constitutifs de base

# /PATTERN/OPTIONS

/at/

*A **cat** eats a black **fat** rat*

# Alternative monospace

/[cft]/

*A cat eats a black fat rat*



# Exclusion monocaractère

`/[^a]/`

*A cat eats a black fat rat*

`/[^rt]/`

*A cat eats a black fat rat*

# Multiplicité

0 ou 1  $/at?/$  -  $/at\{0,1\}/$  -  $/at\{,1\}/$

*A cat eats a black fat rat*

0 à inf.  $/at^*/$  -  $/at\{0,\}/$

*A cat eats a black fat rat*

1 à inf.  $/at+ /$  -  $/at\{1,\}/$

*A cat eats a black fat rat*

$n$  à  $m$   $/at\{n,m\}/$

*A cat eats a black fat rat*

$n$   $/at\{n\}/$

*A cat eats a black fat rat*

# Classes de caractères

Intervals / `[a-z]` /

*Hier, j'ai terminé à 6h00.*

Classes prédéfinies / `[\w]` / vs / `[:alnum]` /

*Hier, j'ai terminé à 6h00.*

Jocker / `.\w/`

*Hier, j'ai\_terminé\_à\_6h00.*

## Non-équivalences à cause des accents dans certains langages

Accents / `[=e=]` /

*Hier, j'ai terminé à 6h00.*

# Capture / Alternatives multi-caractères

`/(ca|ts)/`

*A **cat** eats a black fat rat*

`/(.)(?:la|at)/`

*A **cat** eats a black **fat** rat*

En PCRE `/(bla|blac)/`

*A cat eats a **black** fat rat*

En POSIX `/(bla|blac)/`

*A cat eats a **black** fat rat*

# Changements du comportement d'évaluation

# Ancres

Début / ^ . /

***A** cat eats a black fat rat*

Fin / . \$ /

*A cat eats a black fat **rat***

Délimiteur de mots / . at \b /

*A **cat** eats a black **fat rat***

Qui est différent de l'espace / . at \s /

*A **cat**\_eats a black **fat\_rat**\$*

# Options

/PATTERN/OPTIONS

Insensibilité à la casse /a/i

***A cat eats a black fat rat***

Prise en compte des retours à la ligne des ancres /at\$/m

***A cat\n***

***eats a black fat rat\$***

Prise en compte des retours pour le point /at./s voire /at./s

***A cat\n***

***eats a black fat\_rat***

# Commentaires

/at# Cherche at/x

*A cat eats a black fat rat*



# Comportements glouton

Chercher la plus longue correspondance /a.\*t/

*A cat eats a black fat rat*

Chercher la plus courte correspondance /a.\*?t/

*A cat eats a black fat rat*

# Références aux correspondances

Les correspondances capturées sont accessible via \1 à \9

Chercher les balise à contenu /<(\w+)>.\*<\/\1>/

```
<p></br></p>
```

# Motifs prospectifs

Positif / ( ? = . \* [ A - Z ] ) ( ? = . \* \ d ) . { 5 , } /

*MyPassword*

Négatif / ( ? = . \* [ A - Z ] ) ( ? ! . \* \ d ) . { 5 , } /

*MyPassword*

# Motifs retrospectifs

Positif / (?<=if) (true|false) /

*if **true** while false*

Négatif / (?<!if) (true|false) /

*if true while **false***

# Motifs conditionnels

```
/if (true)?(? (1) (?! else)|false else true)/
```

*if true then false*

*if true*

*if false then true*

```
/if (true)?(? (1) (?! else)) /
```

*if true then false*

*if true*

*if false*

# Et plus encore

- Motifs retardés (récursifs) en Perl
- Motifs autonomes ("super-glouton") en JAVA

# Un exemple

```
/\b(?:[\w_][\w_\d]*)\b(?  
<!if|else|while|for|return)\s*\b([\w_  
[\w_\d]*)\b(?:\s*\(\s*(?!.*\b\1\b)/
```



```
/
\b
(?:[\w_][\w_\d]*) # Groupe 1
\b
(?<!if|else|while|for|return) # Rétrospectif négatif 1
\s*
\b
([\w_][\w_\d]*) # Groupe 2
\b
(?:\s*\() # Prospectif négatif 1
(?:!\s*\b1\b) # Prospectif négatif 2
/xm
```

```
/
\b
(?:[\w_][\w_\d]*) # Identificateur de type non capturé
\b
(?<!if|else|while|for|return) # qui n'est pas un mot clef
\s*
\b
([\w_][\w_\d]*) # Identificateur capturé
\b
(?:!\s*\() # Prospectif négatif 1
(?:!\.*\b\1\b) # Prospectif négatif 2
/xm
```

```
/
\b
(?:[\w_][\w_\d]*) # Identificateur de type non capturé
\b
(?<!if|else|while|for|return) # qui n'est pas un mot clef
\s*
\b
([\w_][\w_\d]*) # Identificateur capturé
\b
(?:!\s*\() # qui n'est pas le nom d'une fonction
(?:!\.*\b1\b) # Prospectif négatif 2
/xm
```

```
/
\b
(?:[\w_][\w_\d]*) # Identificateur de type non capturé
\b
(?<!if|else|while|for|return) # qui n'est pas un mot clef
\s*
\b
([\w_][\w_\d]*) # Identificateur capturé
\b
(?:!\s*\() # qui n'est pas le nom d'une fonction
(?:!\.*\b\1\b) # et qui ne réapparaît plus
/xm
```

```
/
\b
(?:[\w_][\w_\d]*) # Identificateur de type non capturé
\b
(?<!if|else|while|for|return) # qui n'est pas un mot clef
\s*
\b
([\w_][\w_\d]*) # Identificateur capturé
\b
(?:!\s*\() # qui n'est pas le nom d'une fonction
(?:!.*\b\1\b) # et qui ne réapparaît plus
/xm
```

**Liste les variables déclarées mais non utilisées**

# Conclusion

- Utile pour extraire/valider des informations basiques (adresses email, dates, unités lexicales)
- Nombreuses variations en fonction du langage ou de la bibliothèque
- Un "vrai" analyseur lexical/syntaxique devient rapidement vital

# Questions ?

Pour continuer l'entrainement : <http://regexcrossword.com/>