# Shortcodes

# Common Shortcodes

[0-9]

[A-Za-z0-9_]

[\t\f\r\n ]

\d

\w

\s

# Negated Shortcodes

[**^**0-9]

[**^**A-Za-z0-9_]

[**^**\t\f\r\n ]

\D

\W

\S

# PCRE 7.2+ Shortcodes

[\t\f ]  [\r\n]  |  [^\t\f ]  [^\r\n]

\h  \v  |  \H  \V

*PCRE 7.2 was released in June 2007*

# PCRE vs. POSIX Syntax

**PCRE (common)**

\s

\s+

[\s\d]+

**POSIX**

[:space:]

[[:space:]]+

[[:space:][:digit:]]+

| Character classes | | PCRE | | POSIX | |
|---|---|---|---|---|---|
| [0-9] | [^0-9] | \d | \D | [[:digit:]] | [^[:digit:]] |
| [A-Za-z0-9_] | [^A-Za-z0-9_] | \w | \W | [[:word:]] | [^[:word:]] |
| [\t\f\r\n \v] | [^\t\f\r\n \v] | \s | \S | [[:space:]] | [^[:space:]] |
| [\t\f ] | [^\t\f ] | \h | \H | [[:blank:]] | [^[:blank:]] |
| [\r\n] | [^\r\n] | \v | \V | - | - |

PCRE vs. POSIX

# Locale

**English (en)**

| déjà vu | [\w ]+ | ❌ |

**French (fr)**

| déjà vu | [\w ]+ | ✅ |

\b | Uses \w and \W to find boundaries.

# Locale

**English (en)**

déjà vu          [\w ]+     ✕

**French (fr)**

déjà vu          [\w ]+     ✓

# Pitfalls

| [^\D\S] | PCRE vs. POSIX | Locale | Engines & Implementations |

# Inconsistent Implementations

[ \f\n\r\t]

[\f\n\r\t\v\x85\p{Z}]

[ \f\n\r\t\x0B]

[ \f\n\r\t\v\u1680\u180e
\u2000\u2001\u2002
\u2003\u2004\u2005
\u2006\u2007\u2008
\u2009\u200a\u2028
\u2029\u202f\u205f\u3000]

[ \f\n\r\t\v]

[\f\n\r\t\p{Z}]

# To Use or Not to Use ?

## Advantages

- Adjust to locale

## Disadvantages

- Adjust to locale
- Inconsistent implementations
- Low portability
- Difficult to unit test

# Unicode Shortcodes

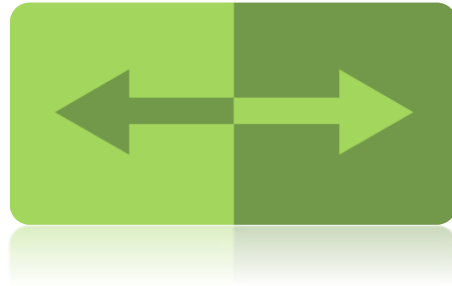# Is Unicode Supported ?

Engine support | Compilation | Input encoding

# Graphmeme Clusters

à
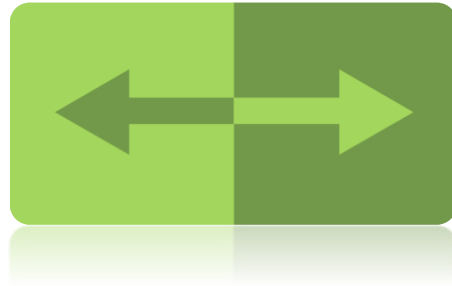


à

U+00E0

U+0061 + U+0300

# Codepoints vs. Graphmeme Clusters

à

à

U+00E0

U+0061 + U+0300

# Codepoints vs. Graphmeme Clusters

à                à

U+00E0           U+0061 + U+0300

**Matches:**

à

# Unicode Wildcard *

## \X

- Matches graphmeme clusters
  - ~ equivalent to \P{M}\p{M}*+
- Matches new line

# Unicode Range Identifiers

Categories

Scripts

Blocks

Binary properties

# Unicode Shortcode Syntax

## Positive

`\p{Identifier}`

## Negative

`\P{Identifier}`

# Unicode Shortcodes

## \p{..}
## \P{..}

- {Identifier}
- Mind: Capitalization
- Can be used anywhere

# Unicode Range Identifiers

| Categories | Scripts |
|---|---|
| Blocks | Binary properties |

# Category-based Unicode Shortcodes

Letters            \p{L}

\p{Ll}    Letter: Lowercase
\p{Lm}    Letter: Mark/modifier
\p{Lo}    Letter: Other
\p{Lt}    Letter: Titlecase
\p{Lu}    Letter: Uppercase

Marks              \p{M}

Numbers            \p{N}

Punctuation        \p{P}

Symbols            \p{S}                    Alternative Syntaxes:

Separators         \p{Z}                    \pX

Other              \p{C}                    \p{Category}

\w ≠ \p{L}

Close approximation:

[\p{L}\p{M}\p{Nd}\p{Nl}\p{Pc}\u200c\u200d]

# Unicode Range Identifiers

Categories

Scripts

Blocks

Binary properties

# Blocks vs. Scripts



Codepoint Block

# Blocks vs. Scripts



Script

# Script

- \p{Scriptname}
- \p{IsScriptName}
- \p{script=ScriptName}
- \p{sc=ScriptName}

# Block

- \p{Blockname}
- \p{InBlockName}
- \p{IsBlockName}
- \p{block=BlockName}
- \p{blk=BlockName}

---

\p{Cyrillic}

\p{InCyrillic}

\p{InCyrillic_Supplementary}

# Modifiers

/[a-z0-9]+/im Modifiers

# Applying Modifiers

/regex/m

m/regex/

match( 'regex', modifiers )

new Re( /regex/, flags )

preg_match()

*vs.*

preg_match_all()

(?m)

g* i m s x

# Modifiers

g*

i

m

s

x

# g

- GLOBAL
- Return all matches vs. first match
- Non-overlapping

# Modifiers

g*

**i**

m

s

x

- CASE-INSENSITIVE

- Mind locales

  - German: FUSSBALL vs. fußball

# Modifiers

g*

i

**m**

s

x

**m**

- MULTILINE
- Affects ^ and $ behaviour

# Modifiers

g*

i

m

s

x

## S

- DOTALL or SINGLELINE
- Affects . (dot) to match \n
- Slow [n]

# Modifiers

g*

i

m     X          - EXTENDED

s

x

```
/^((
  25[0-5]|                    # Match 250-255 range
  2[0-4][0-9]|                # Match 200-249 range
  [01]?[0-9]{1,2}             # Match   0-199 range
)\.){3}                       # Repeat 3 times with period
(25[0-5]|2[0-4][0-9]|[01]?[0-9]{1,2}) # and once without
$/x
```
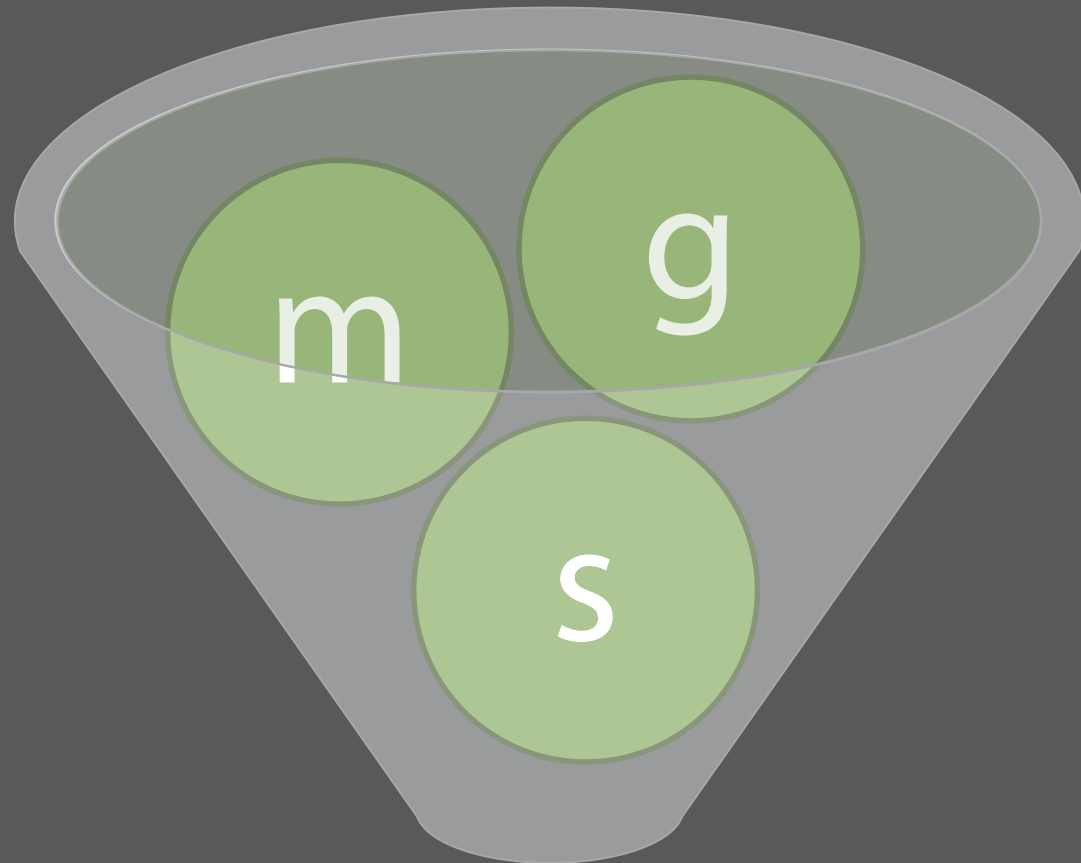
# Modifiers

**g***

**i**

**m**

**s**

**x**

**X**

- EXTENDED

- Ignore whitespace & # to end of line

- Mind: escaping

  - \x20

  - \#

# Inline Modifiers

Setting:

(?i)

(?i)caseless(?-i)cased(?i)caseless

cased(case(?i)insen|sitive)cased

# Inline Modifiers

**Setting:**

**(?i)**

**Unsetting:**

**(?-i)**

**Combined:**

**(?im-sx)**

**Apply to subpattern (non-capturing):**

**(?i:subp)**

# Explore

| S Pattern Analysis | U Ungreedy | l Locale | p Preserve |
| y Sticky | D Dollar end-only | a Ascii | r Return replacement, don't modify |
| u UTF-8 | e EUC-JP | o Interpolation only | c No position reset |

# Delimiters

/[a-z0-9]+/im

Delimiters

# Delimiters

- Enclose the pattern

- Not always needed

- Alternative delimiters

# Alternative Delimiter Requirements

Non-alphanumeric | Non-backslash | Non-whitespace

# Alternative Delimiters

![0-9]+!

#[0-9]+#

@[0-9]+@

`[0-9]+`

~[0-9]+~

%[0-9]+%

# Did you know ?

You can use brackets as delimiters:

**(**p[at]{2}te(rn)**)**    **{**p[at]{2}te(rn)**}**

**[**p[at]{2}te(rn)**]**    **<**p[at]{2}te(rn)**>**

/http:\/\/\/\p{L}+\.[a-z]+\/\//

**vs.**

`http://\p{L}+\.[a-z]+/`

# Choose Wisely

# Next up:



Working with Matches

Juliette Reinders Folmer

@jrf_nl | regexcheatsheets.com