



Reproducible Research RegEx

Niklaus Zemp

June 2020

Genetic Diversity Centre (GDC)

Bioinformatics

ETH Zurich



Replace pattern

>E1_L96

AATTACTTTATGACT

>E2_L119

CGAATTCGTCATTTGAAACCGATTCTGG
CTAGAATT

>E4_L96

TTTTACTTACATGGTGAAAAAATAGAAT
ACGTATTCTCTGCCAAGATTCATTA
CTAACT
CAAAGAGAAATTTTTTGAGTTAATGCA
GAGGATACGAATT

>E1

AATTACTTTATGACT

>E2

CGAATTCGTCATTTGAAACCGATTCTGG
CTAGAATT

>E4

TTTTACTTACATGGTGAAAAAATAGAAT
ACGTATTCTCTGCCAAGATTCATTA
CTAACT
CAAAGAGAAATTTTTTGAGTTAATGCA
GAGGATACGAATT

Regular expression (RegEx)

Groups and Ranges

`.` Any character except new line (`\n`)

`(a|b)` a or b

`(...)` Group

`(?:...)` Passive (non-capturing) group

`[abc]` Range (a or b or c)

`[^abc]` Not (a or b or c)

`[a-q]` Lower case letter from a to q

`[A-Q]` Upper case letter from A to Q

`[0-7]` Digit from 0 to 7

`\x` Group/subpattern number "x"

Ranges are inclusive.

Quantifiers

<code>*</code>	0 or more	<code>{3}</code>	Exactly 3
----------------	-----------	------------------	-----------

<code>+</code>	1 or more	<code>{3,}</code>	3 or more
----------------	-----------	-------------------	-----------

<code>?</code>	0 or 1	<code>{3,5}</code>	3, 4 or 5
----------------	--------	--------------------	-----------

Add a `?` to a quantifier to make it ungreedy.

Anchors

`^` Start of string, or start of line in multi-line pattern

`\A` Start of string

`$` End of string, or end of line in multi-line pattern

`\Z` End of string

`\b` Word boundary

`\B` Not word boundary

`\<` Start of word

`\>` End of word

Tutorial:

RegEx with Atom

>E1_L96

AATTATTACTTTATGACACTGACACTGACACTGACACTGACATAACAGAAATGAATTAAGTCAAGAACCAAAGCGGAGGAAGCGCTTCTAGAGAATT

>E2_L119

CGAATTCGTCATTTGAAACCGATTCAATGAGTTTTAGACTTGAGTTCACGAAGAAGTTTAATGAACTTAAAAACACCCTAGTTCTACTCTTCAAAT

>E4_L96

TTTTACTTACATGGTGAAAAAATAGAATACGTATTCTCTGCCAAGATTCATTAECTCAAAGAGAATTTTTTGAGTTAATGCAGAGGATACGAATT

>E14_L135

CGAATGTCTCCTGGGACTTCTTGGTAGCTTGACCTTCATTCCACATCGTCTGCATTAECTCAAAGAGTATCTTTTGAGTTAATACCTATCGGCTG

4 results found for '(>E[0-9]+)_L[0-9]+' Finding with Options: Regex, Case Insensitive .* Aa [List Icon] [Close Icon]

(>E[0-9]+)_L[0-9]+ 4 found Find Find All

\$1 Replace Replace All

(>E[0-9]+)_L[0-9]+