

# Regular Expression Cheat Sheet for Paratext and RegEx Pal

Reg Ex Function	Description	Alternate expression	sample expression	matching explanation
\	Escape character— do not know what it is until the next character (Reg Ex metacharacter or actual character).		\\ \s	means the \ character means any whitespace
<b>White space</b>				
\r\n	Carriage return and linefeed (end of line)		\r\n	is both parts of a line break
\s	Any whitespace character (including nobreak, thin, en, em spaces, etc.)	[ \r\n\s\t]	\s	match: is_the_[carriage & linefeed return]
<b>Range of characters</b>				
[x-y]	Any one of the characters in the range specified within the brackets.		[a-cx-z]	match: <b>a,b,c,x,y, or z</b>
[^x-y]	Any one of the character not in the range specified within the brackets.		[^a-cx-z]	match: any thing that's NOT <b>a,b,c,x,y, and z</b>
<b>Character classes</b>				
.	Any character <b>except</b> linefeed in <i>RegExPal</i> . Any character in <i>ParaTEExt</i> . An end of line consists of two parts, the carriage return "\r" and linefeed "\n".	[^\n]		match: <u>is the</u> tiger [not the new line]
\w	Any word building character (letters & digits).			match: <u>Wá</u> sp?
\W	Any non-word building character (not a letter and not a digit).			match: Wá <u>sp?</u>
[\w-[\d]]	Any word-building character excluding digits. <b>Note:</b> "-" in front of embedded [] removes digits from the class.	\p{L}		match: <u>Wá</u> sp?
\s	Any whitespace character.	[ \r\n\s\t]		See <u>is</u> above under <i>Whitespace</i>
\S	Any non-whitespace character	[^ \r\n\s\t]		match: <u>is the tiger</u>
\d	Any digit in any script	\p{N}		match: <u>24a</u> <sub>1</sub>
\D	Any character other than a digit.	\P{N}		match: 24 <u>a</u> <sub>1</sub>
[...]	Any one character between the []		[abc]	match: <u>ab</u> acus
[^...]	Any one character not between the []		[^abc]	match: ab <u>ac</u> us
<b>Environment—Context, Anchors, Positioning (finds context but does not capture OR anchors at context)</b>				
(?=...)	Followed by ... (place expression after matched item)		a(=\s)	match: <b>a</b> when followed by a space <u>ha</u> tch, but not hat
(?!...)	Not followed by ... (place expression after matched item)		a(?!\s)	match: <b>a</b> when <b>not</b> followed by a space <u>hat</u> , but not ha t
(?<=...)	Preceded by ... (place expression before matched item)		(?<=\s)c	match: <b>c</b> preceded by a space <u>hat</u> <b>catch</b> , but not hatc
(?<!...)	Not preceded by ... (place expression before matched item)		(?<!\s)t	match: <b>t</b> not preceded by a space <u>at</u> test, testing
\b	Word boundary. Positions to but does not capture the word boundary.		\bin\b	match: word " <u>in</u> ", but not "in" as part of a word as in: bin, or cinch
\B	Not a word boundary. Positions to but does not capture other word building characters.		\Bin\b	match: <u>bin</u> ary, <u>fin</u> e, but not bin, inch, or in
<b>Anchors</b>				
^	Start of record	a record is a <b>chapter</b> in <i>RegExPal</i> a record is a <b>book</b> in <i>ParaTEExt</i>	\A	
\$	End of record		\Z	

**Metacharacters** | when **finding** the actual character place a \ before the metacharacter  
 \ ( ) { } . ? \* + ^ \$ | when **replacing** the actual character only the \ needs to be preceded by a \

In *ParaTEExt* to use a regular expression in a find, **press ctrl-f**, then in the find box key in: **regex**: immediately followed by the regular expression. Regular expressions **cannot** be used in the *replace*.

# Regular Expression Cheat Sheet for Paratext and RegEx Pal

Reg Ex Function	Description	Alternate expression	sample expression	matching explanation
<b>Options—switches</b>				
(?i)	Ignore case—Match either upper or lower case		(?i)a Matches one <b>a</b> at a time	match: lower <u>and</u> uppercase a <u>Adams</u> <u>apple</u>
(?s)	At start of expression dot also matches linefeed.		(?s)is.*tiger" Matches every thing including a newline	match: <b>is the [carriage &amp; line return feed tiger</b>
<b>Repetition</b>				
{n,m}	Match the previous item at least <i>n</i> times but no more than <i>m</i> times.		xa{2,3}l	match: <u>xaal</u> and <u>xaaal</u> , but not xal or xaaal
{n,}	Match the previous item at least <i>n</i> times.		xa{2,}l	match: <u>xaal</u> , <u>xaaal</u> , and <u>xaaal</u> but not xal
{n}	Match exactly <i>n</i> of the previous item.		a{2}	match: only <u>aa</u>
?	Match 0 or 1 times of previous item (It does not or does exist)	{0,1}	fa?ir	match: <u>fir</u> , <u>fair</u> , and <u>afir</u> , but not faair
*	Match 0 or more occurrences of previous item until the <b>last</b> occurrence of that item. <b>GREEDY</b>	{0,}	\\f.*\\f*	match: \\f a \\fr 1.18 \\ft first footnote\\f* and more\\f b \\fr 1.18 \\ft 2nd footnote\\f*
*?	Adding ? matches all occurrences of previous item until <b>first</b> occurrence of the next item. <b>NOT GREEDY</b>	{0,}?	\\f.*?\\f* matches footnote a followed by footnote b	match: \\f a \\fr 1.18 \\ft first footnote\\f* and more\\f b \\fr 1.18 \\ft 2nd footnote\\f*
+	Match 1 or more occurrences of previous item until the <b>last</b> occurrence of that item.	{1,} <b>GREEDY</b>	b(an)+a	match: <u>urbanana</u> and <u>bananana</u>
+?	Match 1 or more occurrences of previous item until <b>first</b> occurrence of that item.	{1,}? <b>NOT GREEDY</b>	b(an)+?a	match: <u>urbanana</u> and <u>bananana</u>
<b>Consider the following scripture text with 2 footnotes and with the start and ending footnote markers <u>underlined</u>:</b>				
\\v 18 This is some scripture text\\f a \\fr 1.18 \\ft first footnote\\f* and more\\f b \\fr 1.18 \\ft second footnote\\f* until the end.				
With 2 footnotes in the verse a greedy match for footnotes \\f.*?\\f* would match the start of the 1st all the way thru the end of the 2nd footnote:				
\\v 18 This is some scripture text\\f a \\fr 1.18 \\ft first footnote\\f* and more\\f b \\fr 1.18 \\ft second footnote\\f* until the end.				
With 2 footnotes in the verse a non greedy match for footnotes \\f.*?\\f*? would match first on footnote a and then on footnote b				
\\v 18 This is some scripture text\\f a \\fr 1.18 \\ft first footnote\\f* and more\\f b \\fr 1.18 \\ft second footnote\\f* until the end.				
<b>Groups—groups are numbered in order of "(" starting from the left. Don't include environment "(" as in "(?"</b>				
(...)	Match and capture what's in parenthesis ( ), store in a group for later reference. Groups are numbered 1–9 based on sequence from left to right of open ( .	EXAMPLE	GROUP #	1 2 3
		find:	(?s)(?<=\\c \\d\\s+)(\\s.*?)(\\s+)(\\r.*)	
		replace:	\\3\\2\\1	
		Swap order of: \\c, \\s, \\r to \\r, \\s, \\c.		
		<b>NOTE:</b> Parenthesis ( followed by a ? as in (?s) are a function and are not assigned a group #.		
	Alternation. Match either side of the		cat dog	match: <u>cat</u> nip or <u>dog</u> ma.
\\1	Match text captured in group 1— first set of ( ). You can reference up to 9 groups.		c.(r e)\\1_	match: <u>carrion</u> and <u>cheech</u> , but not caret or cherish.

A good Regular Expression Web Site — <http://www.regular-expressions.info/unicode.html>  
 Try out regular expressions Web Site — <https://regex101.com/>

**NOTE:** In *ParaText*  
 // is used to denote a line break often used in section heads  
 ~ is used to denote a non-breaking space.

# Regular Expression Cheat Sheet for Paratext and RegEx Pal

Reg Ex Function	Description	samples
<b>Unicode — \p and \P for matching and nonmatching Unicode expressions</b>		
\uFFFF	specific Unicode code point	\u0301 combining acute \u2013 en dash \u201C left double quote
\p{L}	any letter (does not include numbers)	alternate expression [\w-[\d]] not the same as \w, since \w includes numbers
\p{Ll}	any lowercase letter	a-z, à, á, â, ã, è, é, ê, ì, í, î, ï, ð, ñ, ò, ó, ô, õ, etc.
\p{Lu}	any uppercase letter	A-Z, À, Á, Â, Ã, È, É, Ê, Ì, Í, Î, Ï, Ð, Ñ, Ò, Ó, Ô, Õ, etc.
<b>White space</b>		
\p{Z}	any white space character	tab (\u, space, carriage return (\r), newline (\n), enspace(\u2002), etc.
\p{Zs}	any white space character that does not take up space	Zero-width space (\u200b), etc.
<b>Numbers</b>		
\p{N}	any number in any script	1 ١ 2 ٢ includes roman, Arabic-Indic, ideographic, etc.
\p{Nd}	any non-ideographic digit	includes roman, Arabic-Indic (١, ٢, ٣, ...), etc.
\p{No}	superscript or subscript digit, or any digit not 0-9 (excluding ideographic digits)	
<b>Combining characters</b>		
\p{M}	combining characters	includes both \p{Mc} and \p{Mn}
\p{Mn}	zero width combining	combining accents, circumflex, etc
\p{Mc}	combining characters that	middle eastern vowels
<b>Punctuation</b>		
\p{P}	any punctuation characters	
\p{Pd}	any kind of hyphen or dash	includes hyphen, nobreak hyphen, en-dash, em-dash, figure-dash
\p{Ps}	any kind of open/left bracket	includes braces {}, square brackets [], parenthesis ()
\p{Pe}	any kind of close/right bracket	
\p{Pi}	any kind of opening quote	Includes following open quotes: « < ‘ ‚ “ ” „ • “
\p{Pf}	any kind of closing quote	Includes following close quotes: » > ’ ” • ”
\p{Pc}	a punctuation character such as an underscore (low line) that connects words.	_ ~ abc_def
\p{Po}	any punctuation character that is not a dash, bracket, quote or connector.	?¿¡!,:; (to name a few)
<b>::: — Search within a search ONLY works in RegExPal</b>		
:::	match on the expression to the left of the :::, then within that match, match on the expression to the right of the :::	Find cross references that contain a book/chapter separator. <ul style="list-style-type: none"> <li>• first match on \xt and its contents</li> <li>• then on semicolon ; "within \xt match"</li> </ul>

Regular Expression	in RegEx Pal, select	sample output	BACKGROUND—What are you doing? INTERPRET EXPRESSION—What does it mean? ANALYSIS—Interpret results
--------------------	----------------------	---------------	---

RegEx Pal—Insert Regular Expression via: File, USFM

<p>COUNT FOOTNOTE MARKUP</p> <p><code>\\f.*?\\f*</code></p>	<p>Tools Count/Extract</p> <ul style="list-style-type: none"> <li><input checked="" type="radio"/> count</li> <li><input type="radio"/> sort</li> <li><input type="radio"/> combine nonmarker text</li> </ul> <p>Count marker patterns in footnotes (displays "x" for text. NOTE: Be consistent with what precedes \f*. No white space.)</p>	<p>7: \f x \fr x \ft x \fq xlf*</p> <p>1: \f x \fr x \ft x \fq xlf xlf*</p> <p>1: \fr x \ft xlf*</p> <p>9: TOTAL</p>	<p><b>Objective:</b> Count and list footnotes and show the marker patterns collapsing all data in between markers into the letter x.</p> <p><u><code>\\f</code></u> finds the start of a footnote.</p> <p><u><code>.*?</code></u> is a non-greedy match of any character until you find the first occurrence of what follows the "?".</p> <p><u><code>\\f*</code></u> matches closing footnote marker (because it follows "**?" it's the first one following the open footnote marker).</p> <hr/> <p><b>Analysis:</b> 1 footnote starts with \fr and is missing the \f caller id.</p>
---	--	--	---

<p>EXTRACT SECTION HEADS</p> <p><code>\\sld?.*</code></p>	<p>Tools Count/Extract</p> <ul style="list-style-type: none"> <li><input checked="" type="radio"/> extract</li> </ul>	<p><b>will find</b></p> <p>\s The Arrival of the Lover</p> <p>\s2 The Adjuration Refrain</p> <p><b>will not find</b></p> <p>\sp The Beloved to Her...</p> <p>\sc mss\sc* read...</p>	<p><b>Objective:</b> List all section head markers. Include level number when it exists</p> <p><u><code>\\sld?</code></u> finds \s marker followed by an optional number.</p> <p><u><code>.*</code></u> matches everything up to but not including the line break character.</p> <hr/> <p><b>Analysis:</b> Without a space following the optional digit <code>ld?</code> and the anything character <code>.</code> will match <code>\sp</code> and <code>\sc</code>.</p> <p><b>NOTE:</b> Insert space after <code>?</code> to match only section head levels.</p>
---	---	--	---

<p>COUNT CROSS REFERENCE MARKUP</p> <p><code>\\x.*?\\x *</code></p>	<p>Tools Count/Extract</p> <ul style="list-style-type: none"> <li><input checked="" type="radio"/> count</li> <li><input type="radio"/> sort</li> <li><input type="radio"/> combine non-marker text</li> </ul>	<p>254: \x x \xo x \xt x x*</p> <p>2: \x x x*</p> <p>1: \xo x \xt x x*</p> <p>257: TOTAL</p>	<p><b>Objective:</b> Count and list cross references and show the marker patterns collapsing all data in between markers into the letter x.</p> <p><u><code>\\x</code></u> finds a cross reference marker.</p> <p><u><code>.*?</code></u> is non-greedy match of any character until first occurrence of “\”.</p> <p><u><code>\\x *</code></u> matches first closing cross ref.</p> <hr/> <p><b>Analysis:</b> 2 cross refs are missing <code>\xo</code> and <code>\xt</code>. 1 cross ref is missing the opening cross ref marker <code>\x</code>.</p>
---	--	--	--

<p>COUNT ALL USFM</p> <p><code>\\ w+ *?</code></p>	<p>Tools Count/Extract</p> <ul style="list-style-type: none"> <li><input checked="" type="radio"/> count</li> <li><input type="radio"/> sort</li> <li><input type="radio"/> combine non-marker text</li> </ul>	<p>1: \c</p> <p>1: \h</p> <p>1: \id</p> <p>1: \mt1</p> <p>8: \p</p> <p>14: \v</p> <p>29: TOTAL</p>	<p><b>Objective:</b> List all markers</p> <p><u><code>\\</code></u> find start of a marker \.</p> <p><u><code>\\w+</code></u> find 1 or more letters/numbers for marker name.</p> <p><u><code>\\*?</code></u> find optional end marker indicator.</p> <hr/> <p><b>Analysis:</b> It's a one chapter book with header, main title, 14 verses, and 8 paragraphs.</p>
--	--	--	---