

## Lexer.php

```
1. <?php
2. /**
3. * Lexer adapted from Simple Test:
4. * http://sourceforge.net/projects/simpletest/
5. * For an intro to the Lexer see:
6. *
7. * @author Marcus Baker http://www.lastcraft.com
8. */
9.
10. namespace dokuwiki\Parsing\Lexer;
11.
12. /**
13. * Accepts text and breaks it into tokens.
14. *
15. * Some optimisation to make the sure the content is only scanned
16. * by the PHP regex
17. * parser once. Lexer modes must not start with leading
18. * underscores.
19. */
20. class Lexer
21. {
22.     /** @var ParallelRegex[] */
23.     protected $regexes;
24.     /** @var \Doku_Handler */
25.     protected $handler;
26.     /** @var StateStack */
27.     protected $modeStack;
28.     /** @var array mode "rewrites" */
29.     protected $mode_handlers;
30.     /** @var bool case sensitive? */
31.     protected $case;
32.     /**
33.      * Sets up the lexer in case insensitive matching by default.
34.      * @param \Doku_Handler $handler Handling strategy by
35.      * reference.
36.      * @param string $start           Starting handler.
37.      * @param boolean $case           True for case sensitive.
38.     */
39.     public function __construct($handler, $start = "accept", $case
40. = false)
41.     {
42.         $this->case = $case;
43.         $this->regexes = array();
```

```
42.         $this->handler = $handler;
43.         $this->modeStack = new StateStack($start);
44.         $this->mode_handlers = array();
45.     }
46.
47.     /**
48.      * Adds a token search pattern for a particular parsing mode.
49.      *
50.      * The pattern does not change the current mode.
51.      *
52.      * @param string $pattern      Perl style regex, but ( and )  
53.      *                             lose the usual meaning.  
54.      * @param string $mode        Should only apply this  
55.      *                             pattern when dealing with  
56.      *                             this type of input.
57.     */
58.    public function addPattern($pattern, $mode = "accept")
59.    {
60.        if (! isset($this->regexes[$mode])) {
61.            $this->regexes[$mode] = new
62.                ParallelRegex($this->case);
63.            $this->regexes[$mode]->addPattern($pattern);
64.        }
65.
66.        /**
67.         * Adds a pattern that will enter a new parsing mode.
68.         *
69.         * Useful for entering parenthesis, strings, tags, etc.
70.         *
71.         * @param string $pattern      Perl style regex, but ( and )  
72.         *                             lose the usual meaning.
73.         * @param string $mode        Should only apply this pattern  
when dealing with this type of input.
74.         * @param string $new_mode    Change parsing to this new  
nested mode.
75.     */
76.    public function addEntryPattern($pattern, $mode, $new_mode)
77.    {
78.        if (! isset($this->regexes[$mode])) {
79.            $this->regexes[$mode] = new
80.                ParallelRegex($this->case);
81.            $this->regexes[$mode]->addPattern($pattern, $new_mode);
82.        }
83.        /**
84.         * Adds a pattern that will exit the current mode and re-enter  
the previous one.
85.         *
```

```
86.      * @param string $pattern      Perl style regex, but ( and )
87.      * @param string $mode         lose the usual meaning.
88.      */
89.     public function addExitPattern($pattern, $mode)
90.     {
91.         if (! isset($this->regexes[$mode])) {
92.             $this->regexes[$mode] = new
93.             ParallelRegex($this->case);
94.         }
95.         $this->regexes[$mode]->addPattern($pattern, "__exit");
96.     }
97. /**
98. * Adds a pattern that has a special mode.
99. *
100. * Acts as an entry and exit pattern in one go, effectively
101. * calling a special
102. * parser handler for this token only.
103. * @param string $pattern      Perl style regex, but ( and )
104. * @param string $mode         lose the usual meaning.
105. * @param string $special      Should only apply this pattern
106. * when dealing with this type of input.
107. * @param string $special      Use this mode for this one
108. * token.
109. */
110. public function addSpecialPattern($pattern, $mode, $special)
111. {
112.     if (! isset($this->regexes[$mode])) {
113.         $this->regexes[$mode] = new
114.             ParallelRegex($this->case);
115.     }
116.     $this->regexes[$mode]->addPattern($pattern, "__$special");
117. }
118. /**
119. * Adds a mapping from a mode to another handler.
120. *
121. * @param string $mode         Mode to be remapped.
122. * @param string $handler      New target handler.
123. */
124. public function mapHandler($mode, $handler)
125. {
126.     $this->mode_handlers[$mode] = $handler;
127. }
128. /**
129. * Splits the page text into tokens.
130. *
131. * Will fail if the handlers report an error or if no content
```

```
is consumed. If successful then each
130.      * unparsed and parsed token invokes a call to the held
       listener.
131.      *
132.      * @param string $raw           Raw HTML text.
133.      * @return boolean              True on success, else false.
134.      */
135.     public function parse($raw)
136.     {
137.         if (! isset($this->handler)) {
138.             return false;
139.         }
140.         $initialLength = strlen($raw);
141.         $length = $initialLength;
142.         $pos = 0;
143.         while (is_array($parsed = $this->reduce($raw))) {
144.             list($unmatched, $matched, $mode) = $parsed;
145.             $currentLength = strlen($raw);
146.             $matchPos = $initialLength - $currentLength -
                     strlen($matched);
147.             if (! $this->dispatchTokens($unmatched, $matched,
148.               $mode, $pos, $matchPos)) {
149.                 return false;
150.             }
151.             if ($currentLength == $length) {
152.                 return false;
153.             }
154.             $length = $currentLength;
155.             $pos = $initialLength - $currentLength;
156.         }
157.         if (!$parsed) {
158.             return false;
159.         }
160.         return $this->invokeHandler($raw, DOKU_LEXER_UNMATCHED,
161.           $pos);
162.     }
163.     /**
164.      * Gives plugins access to the mode stack
165.      *
166.      * @return StateStack
167.      */
168.     public function getModeStack()
169.     {
170.         return $this->modeStack;
171.     }
172.     /**
173.      * Sends the matched token and any leading unmatched
174.      * text to the parser changing the lexer to a new
```

```
175.     * mode if one is listed.
176.     *
177.     * @param string $unmatched Unmatched leading portion.
178.     * @param string $matched Actual token match.
179.     * @param bool|string $mode Mode after match. A boolean false
180.       mode causes no change.
181.     * @param int $initialPos
182.     * @param int $matchPos Current byte index location in raw doc
183.       thats being parsed
184.     * @return boolean           False if there was any error
185.       from the parser.
186.     */
187.     protected function dispatchTokens($unmatched, $matched, $mode,
188.     $initialPos, $matchPos)
189.     {
190.         if (! $this->invokeHandler($unmatched,
191.           DOKU_LEXER_UNMATCHED, $initialPos)) {
192.             return false;
193.         }
194.         if ($this->isModeEnd($mode)) {
195.             if (! $this->invokeHandler($matched, DOKU_LEXER_EXIT,
196.               $matchPos)) {
197.                 return false;
198.             }
199.             return $this->modeStack->leave();
200.         }
201.         if ($this->isSpecialMode($mode)) {
202.             $this->modeStack->enter($this->decodeSpecial($mode));
203.             if (! $this->invokeHandler($matched,
204.               DOKU_LEXER_SPECIAL, $matchPos)) {
205.                 return false;
206.             }
207.             return $this->invokeHandler($matched, DOKU_LEXER_MATCHED,
208.               $matchPos);
209.         /**
210.          * Tests to see if the new mode is actually to leave the
211.          * current mode and pop an item from the matching
212.          * mode stack.
213.          * @param string $mode    Mode to test.
214.          * @return boolean        True if this is the exit mode.
215.         */
```

```
216.     protected function isModeEnd($mode)
217.     {
218.         return ($mode === "__exit");
219.     }
220.
221.     /**
222.      * Test to see if the mode is one where this mode is entered
223.      * for this token only and automatically
224.      * leaves immediately afterwards.
225.      *
226.      * @param string $mode Mode to test.
227.      * @return boolean True if this is the exit mode.
228.      */
229.     protected function isSpecialMode($mode)
230.     {
231.         return (strncmp($mode, "_", 1) == 0);
232.
233.     /**
234.      * Strips the magic underscore marking single token modes.
235.      *
236.      * @param string $mode Mode to decode.
237.      * @return string Underlying mode name.
238.      */
239.     protected function decodeSpecial($mode)
240.     {
241.         return substr($mode, 1);
242.     }
243.
244.     /**
245.      * Calls the parser method named after the current mode.
246.      *
247.      * Empty content will be ignored. The lexer has a parser
248.      * handler for each mode in the lexer.
249.      *
250.      * @param string $content Text parsed.
251.      * @param boolean $is_match Token is recognised rather
252.      * than unparsed data.
253.      * @param int $pos Current byte index location in raw doc
254.      * thats being parsed
255.      * @return bool
256.      */
257.     protected function invokeHandler($content, $is_match, $pos)
258.     {
259.         if (($content === "") || ($content === false)) {
260.             return true;
261.         }
262.         $handler = $this->modeStack->getCurrent();
263.         if (isset($this->mode_handlers[$handler])) {
264.             $handler = $this->mode_handlers[$handler];
```

```

264.         }
265.
266.         // modes starting with plugin_ are all handled by the same
267.         // handler but with an additional parameter
268.         if (substr($handler, 0, 7)=='plugin_') {
269.             list($handler,$plugin) = sexplode('_', $handler, 2,
270.                  '');
271.             return $this->handler->$handler($content, $is_match,
272.                  $pos, $plugin);
273.         }
274.     }
275.
276.     /**
277.      * Tries to match a chunk of text and if successful removes
278.      * the recognised chunk and any leading
279.      * unparsed data. Empty strings will not be matched.
280.      *
281.      * @param string $raw           The subject to parse. This is
282.      * the content that will be eaten.
283.      * @return array|bool          Three item list of unparsed
284.      * content followed by the    recognised token and finally the
285.      * action the parser is to take.
286.      *
287.      * is a parsing error.
288.      */
289.     protected function reduce(&$raw)
290.     {
291.         if (!
292.             iset($this->regexes[$this->modeStack->getCurrent()])) {
293.             if ($raw === "") {
294.                 return true;
295.             }
296.             if ($action =
297.                 $this->regexes[$this->modeStack->getCurrent()]->split($raw,
298.                     $split)) {
299.                     list($unparsed, $match, $raw) = $split;
300.                     return array($unparsed, $match, $action);
301.                 }
302.                 return true;
303.             }
304.         /**
305.          * Escapes regex characters other than (, ) and /
306.          *
307.          * @param string $str

```

```
304.     * @return string
305.     */
306.     public static function escape($str)
307.     {
308.         $chars = array(
309.             '/\\\\\\\\/',
310.             '/\\.\\/',
311.             '/\\+\\/',
312.             '/\\*\\/',
313.             '/\\?\\/',
314.             '/\\[\\/',
315.             '/\\^\\/',
316.             '/\\]\\/',
317.             '/\\$\\/',
318.             '/\\{\\/',
319.             '/\\}\\/',
320.             '/\\=\\/',
321.             '/\\!\\/',
322.             '/\\<\\/',
323.             '/\\>\\/',
324.             '/\\|\\/',
325.             '/\\:\\/'
326.         );
327.
328.         $escaped = array(
329.             '\\\\\\\\\\\\\\\\\\',
330.             '\\.\\',
331.             '\\+\\',
332.             '\\*\\',
333.             '\\?\\',
334.             '\\[\\',
335.             '\\^\\',
336.             '\\]\\',
337.             '\\$\\',
338.             '\\{\\',
339.             '\\}\\',
340.             '\\=\\',
341.             '\\!\\',
342.             '\\<\\',
343.             '\\>\\',
344.             '\\|\\',
345.             '\\:\\'
346.         );
347.         return preg_replace($chars, $escaped, $str);
348.     }
349. }
```

From:  
<https://wwoss.ru/> - **worldwide open-source software**



Permanent link:  
<https://wwoss.ru/doku.php?id=wiki:xref:dokuwiki:inc:parsing:lexer:lexer.php>

Last update: **2025/01/16 18:55**