# Parsing with PHP

Kore Nordmann <kore@php.net>

August 22, 2009

- ▶ Kore Nordmann, <kore@php.net>, <kn@ez.no>
- ▶ Long time PHP developer
- ▶ Regular speaker, author, etc.
- ▶ Studies computer science in Dortmund
- ▶ Active open source developer:
    - ▶ eZ Components (Graph, WebDav, *Document*), Arbit, PHPUnit, Torii, PHPillow, KaForkL, Image 3D, WCV, ...

Introduction

Examples

The document component

- ▶ Parsers in PHP? Why the hell?

- Parsers in PHP? Why the hell?
- Applications for parsers
    - Markup languages
    - Domain specific languages (DSL)
    - Language interpreters (template languages)

- ▶ Parsers in PHP? Why the hell?
- ▶ Applications for parsers
    - ▶ Markup languages
    - ▶ Domain specific languages (DSL)
    - ▶ Language interpreters (template languages)
- ▶ So, who already wrote a parser in PHP?

▶ Parsing = text processing = regular expressions; Right?

▶ Parsing = text processing = regular expressions; Right?

▶ *No!* – regular expressions only work for regular languages. [1]

- ▶ Parsing = text processing = regular expressions; Right?
- ▶ *No!* – regular expressions only work for regular languages. [1]
- ▶ Regular languages
    - ▶ ... cannot express recursion
    - ▶ ... grammars must be right-linear (right-regular)

- ▶ Parsing = text processing = regular expressions; Right?
- ▶ *No!* – regular expressions only work for regular languages. [1]
- ▶ Regular languages
    - ▶ ... cannot express recursion
    - ▶ ... grammars must be right-linear (right-regular)
- ▶ What does that mean?

```
1  S ::= "(" A ")"
2  A ::= "foo"
3     | S
```

- *But:* PCRE are no real regular expression any more.

- ▶ *But:* PCRE are no real regular expression any more.
- ▶ PCRE knows backreferences

```
1  S ::= "a" S "a"  |  "b"
2
3  ( (a*) b \1 )
```

- *But:* PCRE are no real regular expression any more.
- PCRE knows backreferences

```
1  S ::= "a" S "a" | "b"
2
3  ( (a*) b \1 )
```

- ... and recursion

```
1  S ::= "(" A ")"
2  A ::= "foo" | S
3
4  (
5    \(
6      ( (?>foo) | (?R) )
7    \)
8  )
```

▶ Seems to make people think: Use PCRE for parsing!

▶ You don't even...

▶ Seems to make people think: Use PCRE for parsing!
▶ A PCRE BBCode parser:

```
1  (
2    (
3      [^\[\]]*
4        (?# Match an opening BBCode tag )
5        \[([a-z]+)(?:=([^\]]+))?\]
6          (?# The actual recursion )
7          (?>[^\[\]]* | (?R) )
8        (?# Match the closing tag )
9        \[/\2\]
10     [^\[\]]*
11   )
12 ) i x
```

▶ You don't even …

- ▶ Seems to make people think: Use PCRE for parsing!
- ▶ A PCRE BBCode parser:

```
1   (
2     (
3        [^\[\]]*
4          (?# Match an opening BBCode tag )
5          \[([a-z]+)(?:=([^\]]+))?\]
6            (?# The actual recursion )
7            (?>[^\[\]]* | (?R) )
8          (?# Match the closing tag )
9          \[/\2\]
10       [^\[\]]*
11    )
12  ) ix
```

- ▶ This is useless, because:
  - ▶ It's unmaintainable
    - ▶ You don't even ...

▶ Seems to make people think: Use PCRE for parsing!
▶ A PCRE BBCode parser:

```
1  (
2    (
3      [^\[\]]*
4        (?# Match an opening BBCode tag )
5        \[([a-z]+)(?:=([^\]]+))?\]
6          (?# The actual recursion )
7          (?>[^\[\]]* | (?R) )
8        (?# Match the closing tag )
9        \[/\2\]
10     [^\[\]]*
11   )
12 ) ix
```

▶ This is useless, because:
  ▶ It's unmaintainable
  ▶ You don't get a syntax tree (AST)
  ▶ You don't even...

- ▶ Seems to make people think: Use PCRE for parsing!
- ▶ A PCRE BBCode parser:

```
1  (
2    (
3      [^\[\]]*
4        (?# Match an opening BBCode tag )
5        \[([a-z]+)(?:=([^\]]+))?\]
6          (?# The actual recursion )
7          (?>[^\[\]]* | (?R) )
8        (?# Match the closing tag )
9        \[/\2\]
10     [^\[\]]*
11   )
12 ) ix
```

- ▶ This is useless, because:
  - ▶ It's unmaintainable
  - ▶ You don't get a syntax tree (AST)
  - ▶ You don't even get proper match arrays

- ▶ Regular expressions are *not* for:

- ▶ Regular expressions are *not* for:
  - ▶ Building the full parser for your markup language

- ▶ Regular expressions are *not* for:
    - ▶ Building the full parser for your markup language
    - ▶ Finding contents in HTML documents

- ▶ Regular expressions are *not* for:
  - ▶ Building the full parser for your markup language
  - ▶ Finding contents in HTML documents
- ▶ Regular expressions are perfect for:

- ▶ Regular expressions are *not* for:
  - ▶ Building the full parser for your markup language
  - ▶ Finding contents in HTML documents
- ▶ Regular expressions are perfect for:
  - ▶ Tokenizing (example follows)

- ▶ Regular expressions are *not* for:
    - ▶ Building the full parser for your markup language
    - ▶ Finding contents in HTML documents
- ▶ Regular expressions are perfect for:
    - ▶ Tokenizing (example follows)
    - ▶ Parse regular DSLs (seldom)

Introduction

## Examples

The document component

# Find links

▶ http://www.the-art-of-web.com/php/parse-links/

```
1  $url = "http://www.example.net/somepage.html";
2  $input = @file_get_contents($url) or die('Could not
       access file: $url');
3  $regexp = "<a\s[^>]*href=(\"??)([^\" >]*?)\\1[^>]*>(.*)
       <\/a>";
4  if ( preg_match_all( "/$regexp/siU", $input, $matches )
       )
5  {
6    // $matches[2] = array of link addresses
7    // $matches[3] = array of link text - including HTML
          code
8  }
```

▶ The correct way: [2]

```
1  $oldSetting = libxml_use_internal_errors( true );
2
3  $html = new DOMDocument();
4  $html->loadHtmlFile( $url );
5  $xpath = new DOMXPath( $html );
6  $links = $xpath->query( '//a' );
7  foreach ( $links as $link ) {
8    echo $link->getAttribute( 'href' ), "\n";
9  }
```

▶ The correct way: [2]

```php
1  $oldSetting = libxml_use_internal_errors( true );
2
3  $html = new DOMDocument();
4  $html->loadHtmlFile( $url );
5  $xpath = new DOMXPath( $html );
6  $links = $xpath->query( '//a' );
7  foreach ( $links as $link ) {
8    echo $link->getAttribute( 'href' ), "\n";
9  }
```

▶ Maintainable

▶ The correct way: [2]

```
1  $oldSetting = libxml_use_internal_errors( true );
2
3  $html = new DOMDocument();
4  $html->loadHtmlFile( $url );
5  $xpath = new DOMXPath( $html );
6  $links = $xpath->query( '//a' );
7  foreach ( $links as $link ) {
8    echo $link->getAttribute( 'href' ), "\n";
9  }
```

▶ Maintainable

▶ Correct

▶ The correct way: [2]

```
1  $oldSetting = libxml_use_internal_errors( true );
2
3  $html = new DOMDocument();
4  $html->loadHtmlFile( $url );
5  $xpath = new DOMXPath( $html );
6  $links = $xpath->query( '//a' );
7  foreach ( $links as $link ) {
8    echo $link->getAttribute( 'href' ), "\n";
9  }
```

▶ Maintainable

▶ Correct

▶ Handles recursion properly

▶ http://www.tutorials.de/forum/php-tutorials/279124-bbcode-mit-php-parsen.html

```php
1  function parseBBCode2HTML( $bb )
2  {
3    $bb = preg_replace(
4      '(\[b\](.*?)\[/b\])', '<b>$1</b>', $bb );
5    $bb = preg_replace(
6      '(\[i\](.*?)\[/i\])', '<i>$1</i>', $bb );
7    $bb = preg_replace(
8      '(\[color=([a-f\d]{6}?).*\](.*?)\[/color\])',
9      '<font color="#$1">$2</font>', $bb );
10   $bb = preg_replace(
11     '(\[url=([^ ]+).*\](.*)\[/url\])',
12     '<a href="$1">$2</a>', $bb );
13   $bb = preg_replace(
14     '(\n)', "<br/>\n", $bb );
15   return $bb;
16 }
```

▶ Does not handle invalid markup:

```
1  Input: " Hello ␣[b] world !":
2  ⇒ Hello [b] world !
3
4  Input: " Hello ␣world [/b]!":
5  ⇒ Hello world [/b]!
6
7  Input: " [i] Hello ␣[b] world [/i][/b]!":
8  ⇒ <i>Hello <b>world </i></b>!
```

▶ Fails:
  ▶ Does not report formatting errors
  ▶ Creates invalid markup

- ▶ It is trivial to do it correct.
- ▶ Define tokens first:

```
1  protected $tokens = array(
2    'open'  => '(\\A\\[(?P<value>[a−z]+)\\])',
3    'close' => '(\\A\\[/(?P<value>[a−z]+)\\])',
4    'text'  => '(\\A(?P<value>[^\\[]+|\\[))',
5  );
```

▶ Tokenize input string:

```php
1    public function tokenize( $string ) {
2      $scanned = array();
3      while (strlen($string)) {
4        foreach ($this->tokens as $type => $expression) {
5          if (preg_match($expression, $string, $match)) {
6            $scanned[] = array(
7              'type'    => $type,
8              'content' => $match['value'],
9            );
10           $string = substr($string, strlen($match[0]));
11           continue 2;
12         }
13       }
14       throw new Exception("Could_not_process:_'$string'");
15     }
16     return $scanned;
17   }
```

▶ Build AST from token stream

```php
1  public function parse( array &$tokens, $tag = null ) {
2    $ast = array();
3    while ( $token = array_shift( $tokens ) ) {
4      switch ( $token['type'] ) {
5        case 'text':
6          $ast[] = $token['content'];
7          break;
8        case 'open':
9          $ast[] = array(
10           'tag'     => $token['content'],
11           'content' => $this->parse( $tokens, $token['content'] ),
12          );
13          break;
14        case 'close':
15          if ( $token['content'] !== $tag ) throw new Exception( "Unexpected
                  closing tag: {$token['content']}." );
16          return $ast;
17        }
18      }
19    if ( $tag !== null ) throw new Exception( "Missing closing tag for $tag." );
20    return $ast;
21  }
```

▶ Example result

```
1   Input : " [ i ] Hello ␣ [ b ] world [ / b ] [ / i ] ! " :
2   AST :
3   array ( 2 ) {
4     [ 0 ] ⇒ array ( 2 ) {
5       [ " tag " ]     ⇒ string ( 1 ) " i "
6       [ " content " ] ⇒ array ( 2 ) {
7         [ 0 ] ⇒ string ( 6 ) " Hello ␣
8         [ 1 ] ⇒ array ( 2 ) {
9           [ " tag " ]     ⇒ string ( 1 ) " b "
10          [ " content " ] ⇒ array ( 1 ) {
11            [ 0 ] ⇒ string ( 5 ) " world "
12          }
13        }
14      }
15    }
16    [ 1 ] ⇒ string ( 1 ) " ! "
17  }
```

► Proper error messages:

```
1  Input : " Hello ⌴ [ b ] world ! " :
2  ⟹ Exception : Missing closing tag for b .
3
4  Input : " Hello ⌴ world [ / b ] ! " :
5  ⟹ Exception : Unexpected closing tag : b .
6
7  Input : " [ i ] Hello ⌴ [ b ] world [ / i ] [ / b ] ! " :
8  ⟹ Exception : Unexpected closing tag : i .
```

► TODO (trivial):
  ► Add context information to errors (line, position)
  ► Parse attributes in tags

▶ Parse simplified CSS specifications:

```
page {
    page-size: "A4";
    page-orientation: "portrait";
    padding: "22mm 16mm";

    // Margin for pages specifies an additional
        outer border, which can be used
    // to cut if off later, f.e. in printing
    margin: "0mm";
}

para {
    margin: "3mm 0mm 1mm 0mm";
}
```

▶ The grammar:

```
1  File          ::=  Directive+
2  Directive     ::=  Address '{' Formatting* '}'
3  Formatting    ::=  Name ':' '"' Value '"' ';'
4  Name          ::=  [A–Za–z–]+
5  Value         ::=  [^"]+
6
7  Address       ::=  Element ( Rule )*
8  Rule          ::=  '>'? Element
9  Element       ::=  ElementName ( '.' ClassName | '#'
       ElementId )
10
11 ClassName     ::=  [A–Za–z–]+
12 ElementName   ::=  XMLName* | '*'
13 ElementId     ::=  XMLName
14
15 * XMLName references to http://www.w3.org/TR/REC–
       xml/#NT–Name
```

▶ Tokens

```
1   T_WHITESPACE      => '(\\A\\s+)S',
2   T_COMMENT         => '(\\A/\\*.*\\*/)SUs',
3   T_COMMENT         => '(\\A//.*$)Sm',
4   T_START           => '(\\A\\{)S',
5   T_END             => '(\\A\\})S',
6   T_FORMATTING      => '(\\A(?P<name>[A-Za-z-]+)\\s*:\\
        s*"(?P<value>[^"]+)"\\s*;)S',
7   T_ADDRESS         => '(\\A' . $xmlName . ')S',
8   T_DESC_ADDRESS    => '(\\A>[\\t\\x20]+' . $xmlName .
        ')S',
9   T_ADDRESS_CLASS   => '(\\A\\.[A-Za-z_-]+)S',
10  T_ADDRESS_ID      => '(\\A#' . $xmlName . ')S',
```

▶ Common read() method

```php
1    private function read( array $types, array &$tokens ) {
2        $token = array_shift( $tokens );
3
4        if ( !in_array( $token['type'], $types, true ) ) {
5            $names = array();
6            foreach ( $types as $type )
7            {
8                $names[] = $this->tokenNames[$type];
9            }
10
11           $this->triggerError( E_PARSE,
12               "Expected one of " . implode( ',', $names ) . ", found " .
                       $this->tokenNames[$token['type']] . '.',
13               $this->file, $token['line'], $token['position']
14           );
15       }
16
17       return $token;
18   }
```

► Simple domain specific LL(1) parser

```php
1   $directives     = array();
2   $addressTokens = array( self::T_ADDRESS, self::T_DESC_ADDRESS, self::
        T_ADDRESS_ID, self::T_ADDRESS_CLASS );
3
4   while ( count( $tokens ) > 1 ) {
5       $formats = array(); $address = array();
6
7       do {
8           $addressToken = $this->read( $addressTokens, $tokens );
9           $address[] = $addressToken['match'][0];
10      } while ( $tokens[0]['type'] !== self::T_START );
11
12      $this->read( array( self::T_START ), $tokens );
13
14      while ( $tokens[0]['type'] !== self::T_END ) {
15          $format = $this->read( array( self::T_FORMATTING ), $tokens );
16          $formats[$format['match']['name']] = $format['match']['value'];
17      }
18
19      $this->read( array( self::T_END ), $tokens );
20
21      $directives[] = new ezcDocumentPdfCssDirective(
22          $address,
23          $formats,
24          $this->file, $addressToken['line'], $addressToken['position']
25      );
26  }
```

- Some domain specific languages are regular languages

- ▶ Some domain specific languages are regular languages
- ▶ Can be parsed using regular expressions

- ▶ Some domain specific languages are regular languages
- ▶ Can be parsed using regular expressions
- ▶ Like CSS border specifications

```
1  border: 1px solid #f00 2px dotted black;
```

- ▶ Size definition

```
1  (?:[+-]?\s*(?:\d*\.)?\d+)(?:mm|px|pt|in)?
```

- ▶ Size definition

```
1  (?:[+-]?\s*(?:\d*\.)?\d+)(?:mm|px|pt|in)?
```

- ▶ Border definition

```
1  (?:none|dotted|dashed|solid|double|groove|ridge|
     inset|outset|inherit)
```

- ▶ Size definition

```
1  (?:[+-]?\s*(?:\d*\.)?\d+)(?:mm|px|pt|in)?
```

- ▶ Border definition

```
1  (?: none | dotted | dashed | solid | double | groove | ridge |
       inset | outset | inherit )
```

- ▶ Color definitions:

```
1  (?:#?([0-9a-f])([0-9a-f])([0-9a-f])([0-9a-f])?)
2  (?:#?([0-9a-f]{2})([0-9a-f]{2})([0-9a-f]{2})([0-9a-
      f]{2})?)
3  (?:\s*rgb\s*\(\s*([0-9]+)\s*,\s*([0-9]+)\s*,\s
      *([0-9]+)\s*\)\s*)
```

► Border style definition:

```
1  (?:
2    (?:(?:[+-]?\s*(?:\d*\.)?\d+)(?:mm|px|pt|in)?\s*)?
3    (?:(?:none|dotted|dashed|solid|double|groove|
          ridge|inset|outset|inherit)\s*)?
4    (?:
5      (?: transparent | none |
6        (?:#?([0-9a-f])([0-9a-f])([0-9a-f])([0-9a-f])
            ?) |
7        (?:#?([0-9a-f]{2})([0-9a-f]{2})([0-9a-f]{2})
            ([0-9a-f]{2})?) |
8        (?:\s*rgb\s*\(\s*([0-9]+)\s*,\s*([0-9]+)\s*,\
            s*([0-9]+)\s*\)\s*) |
9      )
10    )?
11  )
```
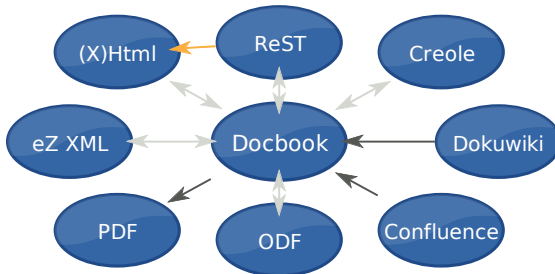
# Parsing CSS border specifications

```
1  (^\s*(?:(?P<m0>)?P<m00>)?(?:(?:(?:[+ —]?\s*(?:\d*\.)?\d+)(?:mm|px|pt|in)?\s*)?(?:(?:i:transparent|none|(?:#?([0—9a—f
   dashed|solid|double|groove|ridge|inset|outset|inherit)\s*)?(?:(?:i:transparent|none|(?:#?([0—9a—f
   ])([0—9a—f])([0—9a—f])([0—9a—f])) ?)|(?:#?([0—9a—f]{2})([0—9a—f]{2})([0—9a—f]{2})([0—9a—f]{2})?)
   |(?:\s*rgb\s*\(\s*([0—9]+)\s*,\s*([0—9]+)\s*,\s*([0—9]+)\s*\)\s*)|(?:\s*rgba\s*\(\s*([0—9]+)\s
   *,\s*([0—9]+)\s*,\s*([0—9]+)\s*,\s*([0—9]+)\s*\)\s*)))?))))|((?P<m1>)?P<m10>(?:(?:(?:[+ —]?\s*(?:\d
   *\.)?\d+)(?:mm|px|pt|in)?\s*)?(?:(?:none|dotted|dashed|solid|double|groove|ridge|inset|outset|
   inherit)\s*)?(?:(?:i:transparent|none|(?:#?([0—9a—f])([0—9a—f])([0—9a—f])([0—9a—f])) ?)|(?:#?([0—9a
   —f]{2})([0—9a—f]{2})([0—9a—f]{2})([0—9a—f]{2})?)|(?:\s*rgb\s*\(\s*([0—9]+)\s*,\s*([0—9]+)\s*,\s
   *([0—9]+)\s*\)\s*)|(?:\s*rgba\s*\(\s*([0—9]+)\s*,\s*([0—9]+)\s*,\s*([0—9]+)\s*\)\s*\)
   s*)))?))\s+(?P<m11>(?:(?:(?:[+ —]?\s*(?:\d*\.)?\d+)(?:mm|px|pt|in)?\s*)?(?:(?:none|dotted|dashed|
   solid|double|groove|ridge|inset|outset|inherit)\s*)?(?:(?:i:transparent|none|(?:#?([0—9a—f])([0—9
   a—f])([0—9a—f])([0—9a—f])) ?)|(?:#?([0—9a—f]{2})([0—9a—f]{2})([0—9a—f]{2})([0—9a—f]{2})?)|(?:\s*
   rgb\s*\(\s*([0—9]+)\s*,\s*([0—9]+)\s*,\s*([0—9]+)\s*\)\s*)|(?:\s*rgba\s*\(\s*([0—9]+)\s*,\s
   *([0—9]+)\s*,\s*([0—9]+)\s*,\s*([0—9]+)\s*\)\s*)))?)))|(?P<m2>(?P<m20>(?:(?:(?:[+ —]?\s*(?:\d*\.)
   ?\d+)(?:mm|px|pt|in)?\s*)?(?:(?:none|dotted|dashed|solid|double|groove|ridge|inset|outset|
   inherit)\s*)?(?:(?:i:transparent|none|(?:#?([0—9a—f])([0—9a—f])([0—9a—f])([0—9a—f])) ?)|(?:#?([0—9a
   —f]{2})([0—9a—f]{2})([0—9a—f]{2})([0—9a—f]{2})?)|(?:\s*rgb\s*\(\s*([0—9]+)\s*,\s*([0—9]+)\s*,\s
   *([0—9]+)\s*\)\s*)|(?:\s*rgba\s*\(\s*([0—9]+)\s*,\s*([0—9]+)\s*,\s*([0—9]+)\s*\)\s
   *)))?))\s+(?P<m21>(?:(?:(?:[+ —]?\s*(?:\d*\.)?\d+)(?:mm|px|pt|in)?\s*)?(?:(?:none|dotted|dashed|
   solid|double|groove|ridge|inset|outset|inherit)\s*)?(?:(?:i:transparent|none|(?:#?([0—9a—f])([0—9
   a—f])([0—9a—f])([0—9a—f])) ?)|(?:#?([0—9a—f]{2})([0—9a—f]{2})([0—9a—f]{2})([0—9a—f]{2})?)|(?:\s*
   rgb\s*\(\s*([0—9]+)\s*,\s*([0—9]+)\s*,\s*([0—9]+)\s*\)\s*)|(?:\s*rgba\s*\(\s*([0—9]+)\s*,\s
   *([0—9]+)\s*,\s*([0—9]+)\s*,\s*([0—9]+)\s*\)\s*)))?))|(?P<m22>(?:(?:(?:[+ —]?\s*(?:\d*\.)?\d+)
   (?:mm|px|pt|in)?\s*)?(?:(?:none|dotted|dashed|solid|double|groove|ridge|inset|outset|inherit)\s
   *)?(?:(?:i:transparent|none|(?:#?([0—9a—f])([0—9a—f])([0—9a—f])([0—9a—f])) ?)|(?:#?([0—9a—f]{2})
   ([0—9a—f]{2})([0—9a—f]{2})([0—9a—f]{2})?)|(?:\s*rgb\s*\(\s*([0—9]+)\s*,\s*([0—9]+)\s*,\s
   *([0—9]+)\s*\)\s*)|(?:\s*rgba\s*\(\s*([0—9]+)\s*,\s*([0—9]+)\s*,\s*([0—9]+)\s*\)\)\s
   *)))?))))|(?P<m3>(?P<m30>(?:(?:(?:[+ —]?\s*(?:\d*\.)?\d+)(?:mm|px|pt|in)?\s*)?(?:(?:none|dotted|
   dashed|solid|double|groove|ridge|inset|outset|inherit)\s*)?(?:(?:i:transparent|none|(?:#?([0—9a—f
   ])([0—9a—f])([0—9a—f])([0—9a—f])) ?)|(?:#?([0—9a—f]{2})([0—9a—f]{2})([0—9a—f]{2})([0—9a—f]{2})?)
   |(?:\s*rgb\s*\(\s*([0—9]+)\s*,\s*([0—9]+)\s*,\s*([0—9]+)\s*\)\s*)|(?:\s*rgba\s*\(\s*([0—9]+)\s*
   *,\s*([0—9]+)\s*,\s*([0—9]+)\s*\)\s*)))?))\s+(?P<m31>(?:(?:(?:[+ —]?\s*(?:\d*\.)?\s
   d+)(?:mm|px|pt|in)?\s*)?(?:(?:none|dotted|dashed|solid|double|groove|ridge|inset|outset|inherit)
   \s*)?(?:(?:i:transparent|none|(?:#?([0—9a—f])([0—9a—f])([0—9a—f])([0—9a—f])) ?)|(?:#?([0—9a—f]{2})
   ([0—9a—f]{2})([0—9a—f]{2})([0—9a—f]{2})?)|(?:\s*rgb\s*\(\s*([0—9]+)\s*,\s*([0—9]+)\s*,\s
   *([0—9]+)\s*\)\s*)|(?:\s*rgba\s*\(\s*([0—9]+)\s*,\s*([0—9]+)\s*,\s*([0—9]+)\s*\)\)
   s*)))?))\s+(?P<m32>(?:(?:(?:[+ —]?\s*(?:\d*\.)?\d+)(?:mm|px|pt|in)?\s*)?(?:(?:none|dotted|...
```

- ▶ Support coversions between document markup formats
- ▶ Applications in content management
    - ▶ Different input mechanisms
        - ▶ WYSISWYG editor (HTML)
        - ▶ Simple text editor (wiki markup)
        - ▶ Emails (ReST)
    - ▶ Different output formats
        - ▶ Web frontend (HTML)
        - ▶ Technical documentation management (Docbook)
        - ▶ Print (PDF)

- ▶ Currently supported formats
    - ▶ Docbook
    - ▶ (X)Html
    - ▶ eZ XML
    - ▶ ReST
    - ▶ Wiki
        - ▶ Dokuwiki, popular PHP based wiki (wiki.php.net) (read-only)
        - ▶ Creole, wiki markup standardization initiative
        - ▶ Confluence, Apache Atlassian wiki dialect (read-only)
- ▶ Currently in development
    - ▶ PDF (write only)
    - ▶ ODF

▶ Docbook as central conversion format
  ▶ Possible conversion shortcuts
  ▶ Conversions always configurable and extensible

▶ Text based markup (RST, Wiki) is easy to write

- ▶ Text based markup (RST, Wiki) is easy to write
- ▶ ... but can be horrible to parse properly.

- ▶ Text based markup (RST, Wiki) is easy to write
- ▶ ... but can be horrible to parse properly.
- ▶ RST is a context-sensitive (type 1) language.
  - ▶ Easy to prove with the pumping lemma for context-free languages [3] on the title markup.

- ▶ Text based markup (RST, Wiki) is easy to write
- ▶ . . . but can be horrible to parse properly.
- ▶ RST is a context-sensitive (type 1) language.
    - ▶ Easy to prove with the pumping lemma for context-free languages [3] on the title markup.
    - ▶ There are no general parser approaches for context-sensitive languages.

- ▶ Text based markup (RST, Wiki) is easy to write
- ▶ . . . but can be horrible to parse properly.
- ▶ RST is a context-sensitive (type 1) language.
    - ▶ Easy to prove with the pumping lemma for context-free languages [3] on the title markup.
    - ▶ There are no general parser approaches for context-sensitive languages.
    - ▶ The document component uses a manually crafted pseudo shift-reduce-parser for those languages.

▶ Design your language with care.

- ▶ Design your language with care.
- ▶ **Design your parser with language properties in mind.**

- ▶ Design your language with care.
- ▶ Design your parser with language properties in mind.
- ▶ Do *not* try to *parse* with regular expressions.

- ▶ Open questions?
- ▶ Further remarks?
- ▶ Contact
    - ▶ Mail: <kore@php.net>
    - ▶ Web: http://kore-nordmann.de/ (Slides will be available here soonish)
    - ▶ Twitter: http://twitter.com/koredn

- ▶ Some further links
  - ▶ http://kore-nordmann.de/blog/0081_parse_html_
    extract_data_from_html.html
  - ▶ http://kore-nordmann.de/blog/do_NOT_parse_using_
    regexp.html
  - ▶ http://ezcomponents.org/docs/tutorials/Document

[1] K. Nordmann.
Do not try parsing with regular expressions.
http://kore-nordmann.de/blog/do_NOT_parse_using_regexp.html,
July 2007.

[2] K. Nordmann.
Extracting data from html.
http://kore-nordmann.de/blog/0081_parse_html_extract_data_
from_html.html, Februrary 2009.

[3] Wikipedia.
Pumping lemma for context-free languages — wikipedia, the free
encyclopedia, 2009.
[Online; accessed 4-August-2009].