

Semantic Web — Spring 2009
Homework #1

Simple XML Parser using Python Lex-Yacc

Behnam Esfahbod
behnam@sharif.edu

April 16, 2009

1 Introduction

This project consists of an XML parser, *xml-ply*, written in Python using *python Lex-Yacc* library, a simple document object model (DOM), and a tree-like view to the document.

Python Lex-Yacc (PLY) [1] is a pure-Python implementation of the popular compiler construction tools *lex* and *yacc*, with the goal to stay fairly faithful to the way in which traditional lex/yacc tools work. PLY consists of two separate modules; `lex.py` and `yacc.py`, both of which are found in a Python package called `ply`.

2 Lexer

The `lex.py` module is used to break input text into a collection of tokens specified by a collection of regular expression rules. Each token is specified by writing a regular expression rule. Each of these rules are defined by making function declarations with a special prefix `t_` to indicate that it defines a token. For simple tokens, the regular expression can be specified as strings. (Python raw strings are used since they are the most convenient way to write regular expression strings)

The lexer has four *states*. The default state, `INITIAL`, is used for all the text outside the XML tags. The `tag` state is used for the tag opening and closing, the tag name, and the attributes names and assignments. The `attrvalue1` and `attrvalue2` states are used for single-quoted and double-quoted strings of attribute values respectively.

For now, we assume that the XML file doesn't contain XML special tags, DTD tags, nor comments.

3 Parser

The `yacc.py` module is used to recognize language syntax that has been specified in the form of a context free grammar. `yacc.py` uses LR parsing and generates its parsing tables using either the *LALR(1)* or *SLR* table generation algorithms. Here we use SLR table generation algorithm.

Each grammar rule is defined by a Python function where the doc-string to that function contains the appropriate context-free grammar specification. Each function accepts a single argument `p` that is a sequence containing the values of each grammar symbol in the corresponding rule.

The non-terminal tokens in the grammar are:

root The root element of the XML document

element The document nodes

opentag An opening tag of a node

closetag A closing tag of a node

lonetag A tag which doesn't have any child and ends right in the opening tag

attributes The sequence of attributes of a tag

attribute The pair of attribute name and value

attrvalue Value of an attribute

children The sequence of children of the tag

child A CDATA or Element node

4 Document Object Model

The document object model (DOM) has two models, **Element** and **Cdata**.

The **Element** model represents a node and contains the node name, the dictionary of attribute name-values, and the list of child nodes.

The **Cdata** model represents the CDATA nodes of the document, which only contains the text.

5 Source Code

Listing 1: The parser (parser.py)

```

1  #!/usr/bin/env python

import sys
from UserString import UserString

6  from ply import lex, yacc

#####
# LEXER
11

class XmlLexer:
    '''The XML lexer'''

16     # states:
    #   default:    The default context, non-tag
    #   tag:        The document tag context
    #   attrvalue1: Single-quoted tag attribute value
    #   attrvalue2: Double-quoted tag attribute value

21     states = (
        ('tag', 'exclusive'),
        ('attrvalue1', 'exclusive'),
        ('attrvalue2', 'exclusive'),
26     )

    tokens = [

31         # state: INITIAL
        'CDATA',
        'OPENTAGOPEN',
        'CLOSETAGOPEN',

        # state: tag
36         'TAGATTRNAME',

```

```

        'TAGCLOSE',
        'LONETAGCLOSE',
        'ATTRASSIGN',

41     # state: attrvalue1
        'ATTRVALUE1OPEN',
        'ATTRVALUE1STRING',
        'ATTRVALUE1CLOSE',

46     # state: attrvalue2
        'ATTRVALUE2OPEN',
        'ATTRVALUE2STRING',
        'ATTRVALUE2CLOSE',

51 ]

# Complex patterns
re_digit      = r'([0-9])'
56 re_nondigit = r'([_A-Za-z])'
re_identifier = r'(' + re_nondigit + r'(' + re_digit + r'|' + re_nondigit + r')
        *)'

# ANY

61 def t_ANY_error(self, t):
    raise SyntaxError("Illegal character '%s'" % t.value[0])
    t.lexer.skip(1)
    pass

66

# INITIAL

t_ignore = ''

71

def t_CLOSETAGOPEN(self, t):
    r'<'
    t.lexer.push_state('tag')
    return t

76

def t_OPENTAGOPEN(self, t):
    r'<'
    t.lexer.push_state('tag')
    return t

81

def t_CDATA(self, t):
    r'[^<]+'
    return t

86

# tag: name

t_tag_ignore = ' \t'

91

def t_tag_TAGATTRNAME(self, t):
    return t
t_tag_TAGATTRNAME.__doc__ = re_identifier

def t_tag_TAGCLOSE(self, t):
96     r'>'
    t.lexer.pop_state()
    return t

def t_tag_LONETAGCLOSE(self, t):
101     r'>'
    t.lexer.pop_state()
    return t

```

```

106     # tag: attr

    t_tag_ATTRASSIGN      = r'='

    def t_tag_ATTRVALUE1OPEN(self, t):
111        r'\''
        t.lexer.push_state('attrvalue1')
        return t

    def t_tag_ATTRVALUE2OPEN(self, t):
116        r'"'
        t.lexer.push_state('attrvalue2')
        return t


121     # attrvalue1

    def t_attrvalue1_ATTRVALUE1STRING(self, t):
        r'[\''+']+'
        t.value = unicode(t.value)
126        return t

    def t_attrvalue1_ATTRVALUE1CLOSE(self, t):
        r'\''
        t.lexer.pop_state()
131        return t

    t_attrvalue1_ignore = ''


136     # attrvalue2

    def t_attrvalue2_ATTRVALUE2STRING(self, t):
        r'[""]+'
        t.value = unicode(t.value)
141        return t

    def t_attrvalue2_ATTRVALUE2CLOSE(self, t):
        r'"'
        t.lexer.pop_state()
146        return t

    t_attrvalue2_ignore = ''


151     # misc

    literals = '$%^'

    def t_ANY_newline(self, t):
156        r'\n'
        t.lexer.lineno += len(t.value)


    # Build the lexer
161     def build(self, **kwargs):
        self.lexer = lex.lex(object=self, **kwargs)

    # Test it output
166     def test(self, data):
        self.lexer.input(data)

        while 1:
            tok = self.lexer.token()
            if not tok: break
            _debug_print_('LEXER', '%-12s' %s' % (self.lexer.lexstate, tok))
171

```

```

# Customization
class SyntaxError(Exception):
176     pass

#####
# PARSER
181
tag_stack = []

# Grammar
186 def p_root_element(p):
    '''root : element
        | element CDATA
    '''
    _parser_trace(p)
191     p[0] = p[1]

def p_root_cdata_element(p):
    '''root : CDATA element
196         | CDATA element CDATA
    '''
    _parser_trace(p)

    p[0] = p[2]
201
def p_element(p):
    '''element : opentag children closetag
                | lonetag
    '''
206     _parser_trace(p)

    if len(p) == 4:
        p[1].children = p[2]

211     p[0] = p[1]

# tag
def p_opentag(p):
    '''opentag : OPENTAGOPEN TAGATTRNAME attributes TAGCLOSE
216     '''
    _parser_trace(p)

    tag_stack.append(p[2])
    p[0] = DOM.Element(p[2], p[3])
221
def p_closetag(p):
    '''closetag : CLOSETAGOPEN TAGATTRNAME TAGCLOSE
    '''
    _parser_trace(p)
226
    n = tag_stack.pop()
    if p[2] != n:
        raise ParserError('Close tag name ("%s") does not match the corresponding
            open tag ("%s").' % (p[2], n))

231 def p_lonetag(p):
    '''lonetag : OPENTAGOPEN TAGATTRNAME attributes LONETAGCLOSE
    '''
    _parser_trace(p)

236     p[0] = DOM.Element(p[2], p[3])

# attr

```

```

def p_attributes(p):
    '''attributes : attribute attributes
    | empty
    '''
    _parser_trace(p)

    if len(p) == 3:
        if p[2]:
            p[1].update(p[2])
            p[0] = p[1]
        else:
            p[0] = p[1]
    else:
        p[0] = {}

def p_attribute(p):
    '''attribute : TAGATTRNAME ATTRASSIGN attrvalue
    '''
    _parser_trace(p)

    p[0] = {p[1]: p[3]}

def p_attrvalue(p):
    '''attrvalue : ATTRVALUE1OPEN ATTRVALUE1STRING ATTRVALUE1CLOSE
    | ATTRVALUE2OPEN ATTRVALUE2STRING ATTRVALUE2CLOSE
    '''
    _parser_trace(p)

    p[0] = _xml_unescape(p[2])

# child
def p_children(p):
    '''children : child children
    | empty
    '''
    _parser_trace(p)

    if len(p) > 2:
        if p[2]:
            p[0] = [p[1]] + p[2]
        else:
            p[0] = [p[1]]
    else:
        p[0] = []

def p_child_element(p):
    '''child : element'''
    _parser_trace(p)

    p[0] = p[1]

def p_child_cdata(p):
    '''child : CDATA'''
    _parser_trace(p)

    p[0] = DOM.Cdata(p[1])

# empty
def p_empty(p):
    '''empty :'''
    pass

# Error rule for syntax errors
class ParserError(Exception):
    pass

def p_error(p):
    raise ParserError("Parse error: %s" % (p,))

```

```

pass

# Customization
def _parser_trace(x):
311     _debug_print_('PARSER', '%-16s' % (sys._getframe(1).f_code.co_name, x))

def _yacc_production__str__(p):
    #return "YaccProduction(%s, %s)" % (str(p.slice), str(p.stack))
    return "YaccP%s" % (str([i.value for i in p.slice]))
316 yacc.YaccProduction.__str__ = _yacc_production__str__

#####
# DOM
321
class DOM:
    class Element:
        # Document object model
        #
326        # Parser returns the root element of the XML document

        def __init__(self, name, attributes={}, children=[]):
            self.name = name
            self.attributes = attributes
331            self.children = children

        def __str__(self):
            attributes_str = ''
            for attr in self.attributes:
336                attributes_str += ' %s="%s"' % (attr, _xml_escape(self.attributes[
                    attr]))

            children_str = ''
            for child in self.children:
341                if isinstance(child, self.__class__):
                    children_str += str(child)
                else:
                    children_str += child

            return '<%s%s>%s</%s>' % (self.name, attributes_str, children_str, self.
                name)

        def __repr__(self):
            return str(self)

    class Cdata(UserString):
351        pass

#####
# ESCAPE
356
_xml_escape_table = {
    "&": "&amp;",
    "'": "&quot;",
    "'": "&apos;",
361    ">": "&gt;",
    "<": "&lt;",
}

def _xml_escape(text):
366    L=[]
    for c in text:
        L.append(_xml_escape_table.get(c,c))
    return "".join(L)

371 def _xml_unescape(s):
    rules = _xml_escape_table.items()

```

```

rules.reverse()

for x, y in rules:
    s = s.replace(y, x)

return s

#####
# INTERFACE

def xml_parse(data):

    _debug_header('INPUT')
    _debug_print_('INPUT', data)
    _debug_footer('INPUT')

    # Tokenizer
    xml_lexer = XmlLexer()
    xml_lexer.build()

    _debug_header('LEXER')
    xml_lexer.test(data)
    _debug_footer('LEXER')

    # Parser
    global tokens
    tokens = XmlLexer.tokens

    yacc.yacc(method="SLR")

    _debug_header('PARSER')
    root = yacc.parse(data, lexer=xml_lexer.lexer, debug=False)
    _debug_footer('PARSER')

    _debug_header('OUTPUT')
    _debug_print_('OUTPUT', root)
    _debug_footer('OUTPUT')

    return root

def tree(node, level=0, init_prefix=''):
    'Returns a tree view of the XML data'

    prefix = '    '
    attr_prefix = '@'
    tag_postfix = ':\\t'
    attr_postfix = ':\\t'

    s_node = init_prefix + node.name + tag_postfix
    s_attributes = ''
    s_children = ''

    for attr in node.attributes:
        s_attributes += init_prefix + prefix + attr_prefix + attr + attr_postfix +
            node.attributes[attr] + '\\n'

    if len(node.children) == 1 and not isinstance(node.children[0], DOM.Element):
        s_node += node.children[0] + '\\n'

    else:
        for child in node.children:
            if isinstance(child, DOM.Element):
                s_children += tree(child, level+1, init_prefix + prefix)

        s_node += '\\n'

```



```

    return s_node + s_attributes + s_children
441

#####
# DEBUG
446 _DEBUG = {
    'INPUT': False,
    'LEXER': False,
    'PARSER': False,
    'OUTPUT': False,
451 }

def _debug_header(part):
    if _DEBUG[part]:
        print '-----'
456     print '%s:' % part

def _debug_footer(part):
    if _DEBUG[part]:
        pass
461

def _debug_print_(part, s):
    if _DEBUG[part]:
        print s
466

#####
# MAIN

def main():
471     data = open(sys.argv[1]).read()
    root = xml_parse(data)
    print tree(root)

if __name__ == '__main__':
476     main()

```

6 Generated Grammar

Listing 2: The generated grammar (parser.out)

```

Created by PLY version 3.2 (http://www.dabeaz.com/ply)
3 Grammar

Rule 0      S' -> root
Rule 1      root -> element
Rule 2      root -> element CDATA
8 Rule 3      root -> CDATA element
Rule 4      root -> CDATA element CDATA
Rule 5      element -> opentag children closetag
Rule 6      element -> lonetag
Rule 7      opentag -> OPENTAGOPEN TAGATTRNAME attributes TAGCLOSE
13 Rule 8      closetag -> CLOSETAGOPEN TAGATTRNAME TAGCLOSE
Rule 9      lonetag -> OPENTAGOPEN TAGATTRNAME attributes LONETAGCLOSE
Rule 10     attributes -> attribute attributes
Rule 11     attributes -> empty
Rule 12     attribute -> TAGATTRNAME ATTRASSIGN attrvalue
18 Rule 13     attrvalue -> ATTRVALUE1OPEN ATTRVALUE1STRING ATTRVALUE1CLOSE
Rule 14     attrvalue -> ATTRVALUE2OPEN ATTRVALUE2STRING ATTRVALUE2CLOSE
Rule 15     children -> child children
Rule 16     children -> empty
Rule 17     child -> element
23 Rule 18     child -> CDATA

```

```

Rule 19      empty -> <empty>

Terminals, with rules where they appear

28 ATTRASSIGN      : 12
   ATTRVALUE1CLOSE : 13
   ATTRVALUE1OPEN  : 13
   ATTRVALUE1STRING : 13
   ATTRVALUE2CLOSE : 14
33 ATTRVALUE2OPEN  : 14
   ATTRVALUE2STRING : 14
   CDATA           : 2 3 4 4 18
   CLOSETAGOPEN    : 8
   LONETAGCLOSE    : 9
38 OPENTAGOPEN     : 7 9
   TAGATTRNAME     : 7 8 9 12
   TAGCLOSE        : 7 8
   error           :

43 Nonterminals, with rules where they appear

   attribute       : 10
   attributes      : 7 9 10
   attrvalue       : 12
48 child           : 15
   children        : 5 15
   closetag        : 5
   element         : 1 2 3 4 17
   empty           : 11 16
53 lonetag        : 6
   opentag         : 5
   root            : 0

Parsing method: SLR

58 state 0

   (0) S' -> . root
   (1) root -> . element
63   (2) root -> . element CDATA
   (3) root -> . CDATA element
   (4) root -> . CDATA element CDATA
   (5) element -> . opentag children closetag
   (6) element -> . lonetag
68   (7) opentag -> . OPENTAGOPEN TAGATTRNAME attributes TAGCLOSE
   (9) lonetag -> . OPENTAGOPEN TAGATTRNAME attributes LONETAGCLOSE

   CDATA           shift and go to state 3
   OPENTAGOPEN      shift and go to state 1

73   lonetag                shift and go to state 4
   element                shift and go to state 2
   root                   shift and go to state 5
   opentag                shift and go to state 6

78 state 1

   (7) opentag -> OPENTAGOPEN . TAGATTRNAME attributes TAGCLOSE
   (9) lonetag -> OPENTAGOPEN . TAGATTRNAME attributes LONETAGCLOSE

83   TAGATTRNAME           shift and go to state 7

state 2

88   (1) root -> element .
   (2) root -> element . CDATA

```

```

93      $end          reduce using rule 1 (root -> element .)
      CDATA          shift and go to state 8

state 3

98      (3) root -> CDATA . element
      (4) root -> CDATA . element CDATA
      (5) element -> . opentag children closetag
      (6) element -> . lonetag
      (7) opentag -> . OPENTAGOPEN TAGATTRNAME attributes TAGCLOSE
103     (9) lonetag -> . OPENTAGOPEN TAGATTRNAME attributes LONETAGCLOSE

      OPENTAGOPEN    shift and go to state 1

      lonetag        shift and go to state 4
108     element      shift and go to state 9
      opentag        shift and go to state 6

state 4

113     (6) element -> lonetag .

      $end          reduce using rule 6 (element -> lonetag .)
      CDATA          reduce using rule 6 (element -> lonetag .)
      OPENTAGOPEN    reduce using rule 6 (element -> lonetag .)
118     CLOSETAGOPEN reduce using rule 6 (element -> lonetag .)

state 5

123     (0) S' -> root .

state 6

128     (5) element -> opentag . children closetag
      (15) children -> . child children
      (16) children -> . empty
      (17) child -> . element
133     (18) child -> . CDATA
      (19) empty -> .
      (5) element -> . opentag children closetag
      (6) element -> . lonetag
      (7) opentag -> . OPENTAGOPEN TAGATTRNAME attributes TAGCLOSE
138     (9) lonetag -> . OPENTAGOPEN TAGATTRNAME attributes LONETAGCLOSE

      CDATA          shift and go to state 11
      TAGCLOSE       reduce using rule 19 (empty -> .)
      LONETAGCLOSE   reduce using rule 19 (empty -> .)
143     CLOSETAGOPEN reduce using rule 19 (empty -> .)
      OPENTAGOPEN    shift and go to state 1

      element        shift and go to state 10
      child          shift and go to state 12
148     lonetag      shift and go to state 4
      children       shift and go to state 13
      empty          shift and go to state 14
      opentag        shift and go to state 6

153 state 7

      (7) opentag -> OPENTAGOPEN TAGATTRNAME . attributes TAGCLOSE
      (9) lonetag -> OPENTAGOPEN TAGATTRNAME . attributes LONETAGCLOSE
      (10) attributes -> . attribute attributes
158     (11) attributes -> . empty
      (12) attribute -> . TAGATTRNAME ATTRASSIGN attrvalue

```

```

(19) empty -> .

TAGATTRNAME      shift and go to state 15
TAGCLOSE         reduce using rule 19 (empty -> .)
LONETAGCLOSE     reduce using rule 19 (empty -> .)
CLOSETAGOPEN     reduce using rule 19 (empty -> .)

attributes       shift and go to state 17
empty            shift and go to state 18
attribute        shift and go to state 16

state 8

(2) root -> element CDATA .

$end             reduce using rule 2 (root -> element CDATA .)

state 9

(3) root -> CDATA element .
(4) root -> CDATA element . CDATA

$end             reduce using rule 3 (root -> CDATA element .)
CDATA            shift and go to state 19

state 10

(17) child -> element .

CDATA            reduce using rule 17 (child -> element .)
OPENTAGOPEN      reduce using rule 17 (child -> element .)
CLOSETAGOPEN     reduce using rule 17 (child -> element .)

state 11

(18) child -> CDATA .

CDATA            reduce using rule 18 (child -> CDATA .)
OPENTAGOPEN      reduce using rule 18 (child -> CDATA .)
CLOSETAGOPEN     reduce using rule 18 (child -> CDATA .)

state 12

(15) children -> child . children
(15) children -> . child children
(16) children -> . empty
(17) child -> . element
(18) child -> . CDATA
(19) empty -> .
(5) element -> . opentag children closetag
(6) element -> . lonetag
(7) opentag -> . OPENTAGOPEN TAGATTRNAME attributes TAGCLOSE
(9) lonetag -> . OPENTAGOPEN TAGATTRNAME attributes LONETAGCLOSE

CDATA            shift and go to state 11
TAGCLOSE         reduce using rule 19 (empty -> .)
LONETAGCLOSE     reduce using rule 19 (empty -> .)
CLOSETAGOPEN     reduce using rule 19 (empty -> .)
OPENTAGOPEN      shift and go to state 1

element          shift and go to state 10
child            shift and go to state 12
lonetag          shift and go to state 4
children         shift and go to state 20

```

```

228      empty                shift and go to state 14
      opentag                shift and go to state 6

state 13

233      (5) element -> opentag children . closetag
      (8) closetag -> . CLOSETAGOPEN TAGATTRNAME TAGCLOSE

      CLOSETAGOPEN          shift and go to state 21

238      closetag            shift and go to state 22

state 14

      (16) children -> empty .

243      CLOSETAGOPEN        reduce using rule 16 (children -> empty .)

state 15

248      (12) attribute -> TAGATTRNAME . ATTRASSIGN attrvalue

      ATTRASSIGN            shift and go to state 23

253      state 16

      (10) attributes -> attribute . attributes
      (10) attributes -> . attribute attributes
258      (11) attributes -> . empty
      (12) attribute -> . TAGATTRNAME ATTRASSIGN attrvalue
      (19) empty -> .

      TAGATTRNAME           shift and go to state 15
263      TAGCLOSE            reduce using rule 19 (empty -> .)
      LONETAGCLOSE          reduce using rule 19 (empty -> .)
      CLOSETAGOPEN          reduce using rule 19 (empty -> .)

      attribute             shift and go to state 16
268      empty               shift and go to state 18
      attributes            shift and go to state 24

state 17

273      (7) opentag -> OPENTAGOPEN TAGATTRNAME attributes . TAGCLOSE
      (9) lonetag -> OPENTAGOPEN TAGATTRNAME attributes . LONETAGCLOSE

      TAGCLOSE              shift and go to state 26
278      LONETAGCLOSE        shift and go to state 25

state 18

      (11) attributes -> empty .

283      TAGCLOSE            reduce using rule 11 (attributes -> empty .)
      LONETAGCLOSE          reduce using rule 11 (attributes -> empty .)

288      state 19

      (4) root -> CDATA element CDATA .

      $end                  reduce using rule 4 (root -> CDATA element CDATA .)

293      state 20

```

```

(15) children -> child children .
298
CLOSETAGOPEN      reduce using rule 15 (children -> child children .)

state 21
303
(8) closetag -> CLOSETAGOPEN . TAGATTRNAME TAGCLOSE
TAGATTRNAME      shift and go to state 27

308
state 22

(5) element -> opentag children closetag .

313
$end              reduce using rule 5 (element -> opentag children closetag .)
CDATA            reduce using rule 5 (element -> opentag children closetag .)
OPENTAGOPEN      reduce using rule 5 (element -> opentag children closetag .)
CLOSETAGOPEN      reduce using rule 5 (element -> opentag children closetag .)

318
state 23

(12) attribute -> TAGATTRNAME ATTRASSIGN . attrvalue
(13) attrvalue -> . ATTRVALUE1OPEN ATTRVALUE1STRING ATTRVALUE1CLOSE
323 (14) attrvalue -> . ATTRVALUE2OPEN ATTRVALUE2STRING ATTRVALUE2CLOSE

ATTRVALUE1OPEN   shift and go to state 29
ATTRVALUE2OPEN   shift and go to state 30

328 attrvalue                      shift and go to state 28

state 24

(10) attributes -> attribute attributes .

333 TAGCLOSE      reduce using rule 10 (attributes -> attribute attributes .)
LONETAGCLOSE     reduce using rule 10 (attributes -> attribute attributes .)

338 state 25

(9) lonetag -> OPENTAGOPEN TAGATTRNAME attributes LONETAGCLOSE .

$end              reduce using rule 9 (lonetag -> OPENTAGOPEN TAGATTRNAME
343 attributes LONETAGCLOSE .)
CDATA            reduce using rule 9 (lonetag -> OPENTAGOPEN TAGATTRNAME
attributes LONETAGCLOSE .)
OPENTAGOPEN      reduce using rule 9 (lonetag -> OPENTAGOPEN TAGATTRNAME
attributes LONETAGCLOSE .)
CLOSETAGOPEN     reduce using rule 9 (lonetag -> OPENTAGOPEN TAGATTRNAME
attributes LONETAGCLOSE .)

348 state 26

(7) opentag -> OPENTAGOPEN TAGATTRNAME attributes TAGCLOSE .

CDATA            reduce using rule 7 (opentag -> OPENTAGOPEN TAGATTRNAME
353 attributes TAGCLOSE .)
OPENTAGOPEN      reduce using rule 7 (opentag -> OPENTAGOPEN TAGATTRNAME
attributes TAGCLOSE .)
CLOSETAGOPEN     reduce using rule 7 (opentag -> OPENTAGOPEN TAGATTRNAME
attributes TAGCLOSE .)

```

```

state 27
358      (8) closetag -> CLOSETAGOPEN TAGATTRNAME . TAGCLOSE
      TAGCLOSE      shift and go to state 31

state 28
      (12) attribute -> TAGATTRNAME ATTRASSIGN attrvalue .
368      TAGATTRNAME      reduce using rule 12 (attribute -> TAGATTRNAME ATTRASSIGN
      attrvalue .)
      TAGCLOSE      reduce using rule 12 (attribute -> TAGATTRNAME ATTRASSIGN
      attrvalue .)
      LONETAGCLOSE      reduce using rule 12 (attribute -> TAGATTRNAME ATTRASSIGN
      attrvalue .)

state 29
373      (13) attrvalue -> ATTRVALUE1OPEN . ATTRVALUE1STRING ATTRVALUE1CLOSE
      ATTRVALUE1STRING shift and go to state 32

state 30
383      (14) attrvalue -> ATTRVALUE2OPEN . ATTRVALUE2STRING ATTRVALUE2CLOSE
      ATTRVALUE2STRING shift and go to state 33

state 31
388      (8) closetag -> CLOSETAGOPEN TAGATTRNAME TAGCLOSE .
      $end      reduce using rule 8 (closetag -> CLOSETAGOPEN TAGATTRNAME
      TAGCLOSE .)
      CDATA      reduce using rule 8 (closetag -> CLOSETAGOPEN TAGATTRNAME
      TAGCLOSE .)
393      OPENTAGOPEN      reduce using rule 8 (closetag -> CLOSETAGOPEN TAGATTRNAME
      TAGCLOSE .)
      CLOSETAGOPEN      reduce using rule 8 (closetag -> CLOSETAGOPEN TAGATTRNAME
      TAGCLOSE .)

state 32
398      (13) attrvalue -> ATTRVALUE1OPEN ATTRVALUE1STRING . ATTRVALUE1CLOSE
      ATTRVALUE1CLOSE shift and go to state 34

state 33
403      (14) attrvalue -> ATTRVALUE2OPEN ATTRVALUE2STRING . ATTRVALUE2CLOSE
      ATTRVALUE2CLOSE shift and go to state 35

state 34
413      (13) attrvalue -> ATTRVALUE1OPEN ATTRVALUE1STRING ATTRVALUE1CLOSE .
      TAGATTRNAME      reduce using rule 13 (attrvalue -> ATTRVALUE1OPEN
      ATTRVALUE1STRING ATTRVALUE1CLOSE .)

```

```

TAGCLOSE      reduce using rule 13 (attrvalue -> ATTRVALUE1OPEN
ATTRVALUE1STRING ATTRVALUE1CLOSE .)
LONETAGCLOSE   reduce using rule 13 (attrvalue -> ATTRVALUE1OPEN
ATTRVALUE1STRING ATTRVALUE1CLOSE .)
418
state 35
(14) attrvalue -> ATTRVALUE2OPEN ATTRVALUE2STRING ATTRVALUE2CLOSE .
423
TAGATTRNAME    reduce using rule 14 (attrvalue -> ATTRVALUE2OPEN
ATTRVALUE2STRING ATTRVALUE2CLOSE .)
TAGCLOSE       reduce using rule 14 (attrvalue -> ATTRVALUE2OPEN
ATTRVALUE2STRING ATTRVALUE2CLOSE .)
LONETAGCLOSE   reduce using rule 14 (attrvalue -> ATTRVALUE2OPEN
ATTRVALUE2STRING ATTRVALUE2CLOSE .)

```

7 Examples

7.1 Example 1

Listing 3: Sample input 1 (sample-1.xml)

```

<Products>
  <Product pid="p123">
3    <Name>gizmo</Name>
    <Price>22.99</Price>
    <Description>great</Description>
    <Store>
      <Name>wiz</Name>
8      <Phone>555-1234</Phone>
      <Markup>25</Markup>
    </Store>
    <Store>
      <Name>Econo-Wiz</Name>
13     <Phone>555-6543</Phone>
      <Markup>15</Markup>
    </Store>
  </Product>
  <Product pid="p231">
18   <Name>gizmoPlus</Name>
    <Price>99.99</Price>
    <Description>more features</Description>
    <Store>
      <Name>wiz</Name>
23     <Phone>555-1234</Phone>
      <Markup>10</Markup>
    </Store>
  </Product>
  <Product pid="p312">
28   <Name>gadget</Name>
    <Price>59.99</Price>
    <Description>good value</Description>
  </Product>
</Products>

```

Listing 4: Output 1 (sample-1.out)

```

Products:
  Product:
3    @pid: p123
      Name: gizmo
      Price: 22.99
      Description: great

```



```

8      Store:
      Name: wiz
      Phone: 555-1234
      Markup: 25
      Store:
      Name: Econo-Wiz
13     Phone: 555-6543
      Markup: 15
      Product:
      @pid: p231
      Name: gizmoPlus
18     Price: 99.99
      Description: more features
      Store:
      Name: wiz
      Phone: 555-1234
23     Markup: 10
      Product:
      @pid: p312
      Name: gadget
      Price: 59.99
28     Description: good value

```

7.2 Example 2

Listing 5: Sample input 2 (sample-2.xml)

```

1  <recipe name="bread" prep_time="5 mins" cook_time="3 hours">
    <title>Basic bread</title>
    <ingredient amount="8" unit="dL">Flour</ingredient>
    <ingredient amount="10" unit="grams">Yeast</ingredient>
6   <ingredient amount="4" unit="dL" state="warm">Water</ingredient>
    <ingredient amount="1" unit="teaspoon">Salt</ingredient>
    <instructions>
        <step>Mix all ingredients together.</step>
        <step>Knead thoroughly.</step>
11    <step>Cover with a cloth, and leave for one hour in warm room.</step>
        <step>Knead again.</step>
        <step>Place in a bread baking tin.</step>
        <step>Cover with a cloth, and leave for one hour in warm room.</step>
        <step>Bake in the oven at 180(degrees)C for 30 minutes.</step>
16    </instructions>
    </recipe>

```

Listing 6: Output 2 (sample-2.out)

```

recipe:
2   @cook_time: 3 hours
   @name: bread
   @prep_time: 5 mins
   title: Basic bread
   ingredient: Flour
7   @amount: 8
   @unit: dL
   ingredient: Yeast
   @amount: 10
   @unit: grams
12  ingredient: Water
   @amount: 4
   @state: warm
   @unit: dL
   ingredient: Salt
17  @amount: 1
   @unit: teaspoon

```

```
22  instructions:
    step: Mix all ingredients together.
    step: Knead thoroughly.
    step: Cover with a cloth, and leave for one hour in warm room.
    step: Knead again.
    step: Place in a bread baking tin.
    step: Cover with a cloth, and leave for one hour in warm room.
    step: Bake in the oven at 180(degrees)C for 30 minutes.
```

References

- [1] PLY (Python Lex-Yacc) <http://www.dabeaz.com/ply/>