



# PyRXP 0.5

## User Documentation

---

Lombard Business Park  
8 Lombard Road  
Wimbledon  
London, ENGLAND SW19 3TZ

103 Bayard Street  
New Brunswick  
New Jersey, 08901  
USA

# PyRXP Documentation

## Contents

1. Introduction	3
1.1 Who is this document aimed at?	3
1.2 What is PxRXP?	3
1.3 License terms	3
1.4 Why another XML toolkit?	3
1.5 Design Goals	4
1.6 Design non-goals	4
1.7 How fast is it?	5
1.8 The Tuple Tree structure	5
1.9 Can I get involved?	6
2. Installation and Setup	7
2.1 Windows binary - pyRXP.pyd	7
2.2 Source Code installation	7
2.3 Examples	7
3. Using pyRXP	8
3.1 Simple use without validation	8
3.2 Validating against a DTD	10
3.3 Interface Summary	12
3.4 Parser Object Attributes and Methods	12
3.5 List of Flags	14
3.6 Flag explanations and examples	15
4. The examples and utilities	21
4.1 Benchmarking	21
4.2 xmlutils and the TagWrapper	21
5. Future Directions	23
5.1 Test Suite	23
5.2 Standardize the wrapper	23
5.3 Other parsers	23
5.4 Better benchmark suite	23
5.5 Type Conversion Utility	23
5.6 Source file references	23
5.7 (longer term and debatable) Richer tuple tree structure	24

# 1. Introduction

## 1.1 Who is this document aimed at?

Anyone who wants to know how to use the pyRXP parser extension from Python. It's assumed that you know how to use the Python programming language and understand its terminology. We make no attempt to teach you XML, so you should already know the basics (what a DTD is, some of the syntax etc.)

## 1.2 What is PxRXP?

PyRXP is a Python language wrapper around the excellent RXP parser. RXP is a validating namespace-aware XML parser written in C. Together, they provide the fastest XML-parsing framework available to Python programmers today.

RXP was written by Richard Tobin at the Language Technology Group, Human Communication Research Centre, University of Edinburgh. PyRXP was written by Robin Becker at ReportLab.

This documentation describes pyRXP-0.5 being used with RXP 1.2.5, as well as ReportLab's emerging XML toolkit which uses it.

## 1.3 License terms

Edinburgh University have released RXP under the GPL. This is generally fine for in-house or open-source use. But if you want to use it in a closed-source commercial product, you may need to negotiate a separate license with them. By contrast, most Python software uses a less restrictive license; Python has its own license, and ReportLab uses the FreeBSD license for our PDF Toolkit, which means you CAN use it in commercial products.

We have therefore no choice but to do the same. We licensed RXP for our commercial products, but are releasing pyRXP under the GPL. If you want to use pyRXP for a commercial product, you need to negotiate with ReportLab and with Edinburgh University. We fully encourage commercial enquiries and will act as a single point of contact. Any commercial licensing deals will help to support both development teams involved.

However, the XML framework ReportLab is using and building will be under our own license. It predates pyRXP and can be made to work off any XML parser (such as expat), and we hope to produce something which can go into the Python distribution one day.

## 1.4 Why another XML toolkit?

This grew out of real world needs which others in the Python community may share. ReportLab make tools which read in some kind of data and make PDF reports. One common input format these days is XML. It's very convenient to express the interface to a system as an XML file. Some other system might send us some XML with tags like <invoice> and <customer>, and we turn these into nice looking invoices.

Also, we have a commercial product called Report Markup Language – we sell a converter to turn RML files into PDF. This has to parse XML, and do it fast and accurately.

Typically we want to get this XML into memory as fast as possible. And, if the performance penalties are not too great, we'd like the option to validate it as well. Validation is useful because we can stop bad data at the point of input; if someone else sends our system an XML 'invoice packet' which is not valid according to the agreed DTD, and gets a validation error, they will know what's going on. This is a lot more helpful than getting a strange and unrelated-sounding error during the formatting stage.

We tried to use all the parsers we could find. We found that almost all of them were constructing large object models in Python code, which took a long time and a lot of memory. Even the fastest C-based parser, expat (which was not yet a standard part of Python at the time), calls back into Python code on every start and end tag, which defeats most of the benefit. Aaron Watters sat down for a couple of days and produced his own parser, rparsexml, which uses string.find and got pretty much the same speed as pyexpat. We also evolved our own representation of a tree in memory, which became the cornerstone of our approach; and when we found RXP we found it easy to make a wrapper around it to produce the "tuple tree".

We have now reached the point in our internal bag-of-tools where XML parsing is a standard component, running entirely at C-like speeds, and we don't even think much about it any more. Which means we must be doing something right and it's time to share it :-)

## 1.5 Design Goals

This is part of an XML framework which we will polish up and release as soon as it can all be documented. The general components are:

- A standard in-memory representation of an XML document (the *tuple tree* below)
- Various parsers which can produce this – principally pyRXP
- A 'lazy wrapper' around this which gives a very friendly Pythonic interface for navigating the tree
- A lightweight transformation tool which does a lot of what XSLT can do, but again with Pythonic syntax

In general we want to get the whole structure of an XML document into memory as soon as possible. Having done so, we're going to traverse through it and move the data into our own object model anyway; so we don't really care what kind of "node objects" we're dealing with and whether they are DOM-compliant. Our goals for the whole framework are:

- Fast - XML parsing should not be an overhead for a program
- Validate when needed, with little or no performance penalty
- Construct a complete tree in memory which is easy and natural to access

Note that pyRXP is just the main parsing component and not the framework itself.

## 1.6 Design non-goals

It's often much more helpful to spell out what a system or component will NOT do. Most of all we are NOT trying to produce a standards-compliant parser.

- Not a SAX parser
- Not a DOM parser
- Does not capture full XML structure

Why not? Aren't standards good?

It's great that Python has support for SAX and DOM, but these are basically Java (or at least cross-platform) APIs. If you're doing Python, it's possible to make things simpler, and we've tried. Let's imagine you have some XML containing an *invoice* tag, that this in turn contains *lineItems* tags, and each of these has some text content and an *amount* attribute. Wouldn't it be nice if you could write some Python code this simple?

```
invoice = pyRXP.Parser().parse(myInvoiceText)
for lineItem in invoice:
    print invoice.amount
```

Likewise, if a node is known to contain text, it would be really handy to just treat it as a string. We have a preprocessor we use to insert data into HTML and RML files which lets us put Python expressions in curly braces, and we often do stuff like

```

<html><head><title>My web page</title></head>
<body>
<h1>Statement for {{xml.customer.DisplayName}}</h1>
<!-- etc etc -->
</body>
</html>
<h1></h1>

```

Try to write the equivalent in Java and you'll have loads of method calls to `getFirstElement()`, `getNextElement()` and so on. Python has beatifully compact and readable syntax, and we'd rather use it. So we're not bothering with SAX and DOM support ourselves. (Although if other people want to contribute full DOM and SAX wrappers for pyRXP, we'll accept the patches).

## 1.7 How fast is it?

The examples file includes a crude benchmarking script. It measures speed and memory allocation of a number of different parsers and frameworks. This is documented later on. Suffice to say that we can parse hamlet in 0.15 seconds with full validation on a P500 laptop. Doing the same with the *minidom* in the Python distribution takes 33 times as long and allocates 8 times as much memory, and does not validate. It also appears to have a significant edge on Microsoft's XML parser and on FourThought's cDomlette. Using pyRXP means that XML parsing will typically take a tiny amount of time compared to whatever your Python program will do with the data later.

## 1.8 The Tuple Tree structure

Most 'tree parsers' such as DOM create 'node objects' of some sort. The DOM gives one consensus of what such an object should look like. The problem is that "objects" means "class instances in Python", and the moment you start to use such beasts, you move away from fast C code to slower interpreted code. Furthermore, the nodes tend to have magic attribute names like "parent" or "children", which one day will collide with structural names.

So, we defined the simplest structure we could which captured the structure of an XML document. Each tag is represented as a tuple of

```
(tagName, dict_of_attributes, list_of_children, spare)
```

The `dict_of_attributes` can be `None` (meaning no attributes) or a dictionary mapping attribute names to values. The `list_of_children` may either be `None` (meaning a singleton tag) or a list with elements that are 4-tuples or plain strings.

A great advantage of this representation - which only uses built-in types in Python - is that you can marshal it (and then zip or encrypt the results) with one line of Python code. Another is that one can write fast C code to do things with the structure. And it does not require any classes installed on the client machine, which is very useful when moving xml-derived data around a network.

This does not capture the full structure of XML. We make decisions before parsing about whether to expand entities and CDATA nodes, and the parser deals with it; after parsing we have most of the XML file's content, but we can't get back to the original in 100% of cases. For example both of the following will (with default settings) return the string "Smith & Jones", and you can't tell from the tuple tree which one was in the file:

```
<provider>Smith & Jones</provider>
```

Alternatively one can use

```
<provider><[CDATA[Smith & Jones]]></provider>
```

So if you want a tool to edit and rewrite XML files with perfect fidelity, our model is not rich enough. However, note that RXP itself DOES provide all the hooks and could be the basis for such a parser.

## 1.9 Can I get involved?

Sure! Join us on the Reportlab-users mailing list ([www.egroups.com/group/reportlab-users](http://www.egroups.com/group/reportlab-users)), and feel free to contribute patches. The final section of this manual has a brief "wish list".

Because the Reportlab Toolkit is used in many mission critical applications and because tiny changes in parsers can have unintended consequences, we will keep checkin rights on sourceforge to a trusted few developers; but we will do our best to consider and process patches.

## 2. Installation and Setup

We make available pre-built Windows binaries. On other platforms you can build it from source using distutils. pyRXP is a single extension module with no other dependencies outside Python itself.

### 2.1 Windows binary - pyRXP.pyd

ReportLab's FTP server has a win32-dlls directory, which is sub-divided into Python versions. Each of these has the version of the pyd file suitable for use with that version of Python. So, the version we use with Python 2.1 is at

<http://www.reportlab.com/ftp/win32-dlls/2.1/pyRXP.pyd>

Download the pyRXP DLL from the ReportLab FTP site. Save the pyRXP.pyd in the DLLs directory under your Python installation (eg this is the C:\Python21\DLLs directory for a standard Windows installation of Python 2.1).

### 2.2 Source Code installation

The source code is open source under the GPL. This is available on SourceForge.

The source for pyRXP and a slightly patched version of RXP is made available by anonymous CVS at

```
:pserver:anonymous@cvs.reportlab.sourceforge.net:/cvsroot/reportlab
```

To get the source use the commands

```
cvs -d :pserver:anonymous@cvs.reportlab.sourceforge.net:/cvsroot/reportlab login
cvs -d :pserver:anonymous@cvs.reportlab.sourceforge.net:/cvsroot/reportlab co rl_addons/pyRXP
```

enter a carriage return for the password.

If you have obtained the source code in the way described above, the rl\_addons/pyRXP directory should contain a distutils script, setup.py which should be run with argument install or build. If successful a shared library pyRXP.pyd or pyRXP.so should be built.

### 2.3 Examples

We have made available a small directory of example stuff to play with. This will be superseded by the release of the framework soon. As such there is no formal package location for it; unzip anywhere you want.

[http://www.reportlab.com/ftp/pyRXP\\_examples.xip](http://www.reportlab.com/ftp/pyRXP_examples.xip).

The examples directory includes a couple of substantial XML files with DTDs, a wrapper module called *xmutils* which provides easy access to the tuple tree, and the beginnings of a benchmarking script. The benchmark script tries to find lots of XML parsers on your system. Both are documented in section 4 below.

## 3. Using pyRXP

### 3.1. Simple use without validation

Firstly you have to import the pyRXP module (using Python's `import` statement). While we are here, pyRXP has a couple of attributes that are worth knowing about: `version` gives you a string with the version number of the pyRXP module itself, and `RXPVersion` gives you string with the version information for the rxp library embedded in the module.

```
C:\Python21>python
Python 2.1.1 (#20, Jul 20 2001, 01:19:29) [MSC 32 bit (Intel)] on win32
Type "copyright", "credits" or "license" for more information.
>>> import pyRXP
>>> pyRXP.version
'0.5'
>>> pyRXP.RXPVersion
'RXP 1.2.5 Copyright Richard Tobin, LTG, HCRC, University of Edinburgh'
```

Once you have imported pyRXP, you can instantiate a parser instance using the `Parser` class.

```
>>>p=pyRXP.Parser()
```

This by itself isn't very useful. But it does allow us to create a single parser which we can reuse many times. It also allows us to type a short variable name rather than 'pyRXP.Parser' every time we need to use it. `p` is now an instance of `Parser` – `Parser` is a constructor that creates an object with its own methods and attributes. When you create a parser like this you can also set multiple flags at the same time. This can save you from having to set them separately, or having to set them all repeatedly each time you need to do a parse.

To parse some XML, you use the `parse` method. The simplest way of doing this is to feed it a string. You could create the string beforehand, or read it from disk (using something like `s=open('filename', 'r').read()`). PyRXP isn't designed to allow you to read the source directly from disk without an intermediate step like this.

We'll start with some very simple examples and leave validation for later.

```
>>> p.parse('<tag>content</tag>')
('tag', None, ['content'], None)
```

This could also more long-windedly be expressed as `pyRXP.Parser().parse('content')`

The parser has taken this single line of XML and returned a tuple-tree. This tuple is made up of the following elements: 'tag': the tag name (aka element name). None: a dictionary of the tag's attributes (null here since it doesn't have any). ['content']: a list of included textual results. This is the contents of the tag. None: the fourth element is unused. It's provided for debugging info, line numbers or other possible future uses.

This tree structure is equivalent to the input XML, at least in information content. It is theoretically possible to recreate the original XML from this tree since no information has been lost.

A tuple tree for more complex XML snippets will contain more of these tuples, but they will all use the same structure as this one.

```
>>> p.parse('<tag1><tag2>content</tag2></tag1>')
('tag1', None, [('tag2', None, ['content'], None)], None)
```



This may be easier to understand if we lay it out differently:

```
>>> p.parse('<tag1><tag2>content</tag2></tag1>')
('tag1',
 None,
 [('tag2',
  None,
  ['content'],
  None)
 ],
 None)
```

Tag1 is the name of the outer tag, which has no attributes. Its contents is a list. This contents contains Tag2, which has its own attribute dictionary (which is also empty since it has no attributes) and its content, which is the string 'content'. It has the closing null element, then the list for Tag2 is closed, Tag1 has its own final null element and it too is closed.

The XML that is passed to the parser must be balanced. Any opening and closing tags must match. They wouldn't be valid XML otherwise.

Look at the following three examples. The first is one a fairly ordinary tag with contents. The second and third can both be considered as empty tags – one is a tag with no content between its opening and closing tag, and the other is the singleton form which by definition has no content.

```
>>> p.parse('<tag>my contents</tag>')
('tag', None, ['my contents'], None)

>>> p.parse('<tag></tag>')
('tag', None, [], None)

>>> p.parse('<tag/>')
('tag', None, None, None)
```

Notice how the contents list is handled differently for the last two examples. This is how we can tell the difference between an empty tag and its singleton version. If the content list is empty then the tag doesn't have any content, but if the list is None, then it can not have any content since it's the singleton form which can't actually have any.

Another example:

```
>>>p.parse('<outerTag><innerTag>bb</innerTag>aaa<singleTag/></outerTag>')
('outerTag', None, [('innerTag', None, ['bb'], None), 'aaa', ('singleTag',
None, None, None)], None)
```

Again, this is more understandable if we show it like this:

```
('outerTag',
 None,
 [('innerTag',
  None,
  ['bb'],
  None),
  'aaa',
  ('singleTag',
   None,
   None,
   None)
 ],
 None)
```

In this example, the tuple contains the outerTag (with no attribute dictionary), whose list of contents are the innerTag, which contains the string 'bb' as its contents, and the singleton singleTag whose contents list is replaced by a null.

By default, PyRXP ignores comments whose contents are lost (this behaviour can be changed – see the section of Flags later for details).

```
>>> p.parse('<tag/><!-- this is a comment about the tag -->')
('tag', None, None, None)

>>> p.parse('<!-- this is a comment -->')
Traceback (most recent call last):
  File "<interactive input>", line 1, in ?
Error: Error: Document ends too soon
in unnamed entity at line 1 char 22 of [unknown]
```

This causes an error, since the parser sees an empty string which isn't valid XML.

It also brings us on to the `srcName` attribute. The Parser has an attribute called `srcName` which is useful when debugging. This is the name by which pyRXP refers to your code in tracebacks. This can be useful – for example, if you have read the XML in from a file, you can use the `srcName` attribute to show the filename to the user. It doesn't get used for anything other than pyRXP Errors – `SyntaxErrors` and `IOErrors` still won't refer to your XML by name.

```
>>> p.parse('<a>aaa</a>') # note the missing '>'
Traceback (most recent call last):
  File "<stdin>", line 1, in ?
pyRXP.Error: Error: Expected > after name in end tag, but got <EOE>
in unnamed entity at line 1 char 10 of yourString
```

The XML that is passed to the parser must be balanced. Not only must the opening and closing tags match (they wouldn't be valid XML otherwise), but there must also be one tag that encloses all the others. If there are otherwise valid fragments that aren't enclosed by another valid tag, they are considered 'multiple elements' and a pyRXP Error is raised.

```
>>> p.parse('<a></a><b></b>')
Traceback (most recent call last):
  File "<stdin>", line 1, in ?
pyRXP.Error: Error: Document contains multiple elements
in unnamed entity at line 1 char 9 of [unknown]

>>> p.parse('<outer><a></a><b></b></outer>')
('outer', None, [('a', None, [], None), ('b', None, [], None)], None)
```

## 3.2. Validating against a DTD

This section describes the default behaviours when validating against a DTD. Most of these can be changed – see the section on flags later in this document for details on how to do that.

For the following examples, we're going to assume that you have a single directory with the DTD and any test files in it.

```
>>> import os
>>> os.getcwd()
'C:\tmp\pyRXP_tests'

>>> os.listdir('.')
['sample1.xml', 'sample2.xml', 'sample3.xml', 'sample4.xml', 'tinydtd.dtd']

>>> dtd = open('tinydtd.dtd', 'r').read()

>>> print dtd
<!-- A tiny sample DTD for use with the PyRXP documentation -->
<!-- $Header $-->

<!ELEMENT a (b)>
<!ELEMENT b (#PCDATA)*>
```

This is just to show you how trivial the DTD is for this example. It's about as simple as you can get – two tags, both mandatory. Both a and b must appear in an xml file for it to conform to this DTD, but you can have as many b's as you want, and they can contain any content.

```
>>> fn=open('sample1.xml', 'r').read()

>>> print fn
<?xml version="1.0" encoding="iso-8859-1" standalone="no" ?>
<!DOCTYPE a SYSTEM "tinydtd.dtd">

<a>
<b>This is the contents</b>
</a>
```

This is the simple example file. The first line is the XML declaration, and the *standalone="no"* part says that there should be an external DTD. The second line says where the DTD is, and gives the name of the root element – a in this case. If you put this in your XML document, pyRXP will attempt to validate it.

```
>>p.parse(fn)
('a',
 None,
 ['\n', ('b', None, ['This tag is the contents'], None), '\n'],
 None)
>>>
```

This is a successful parse, and returns a tuple tree in the same way as we have seen where the input was a string.

If you have a reference to a non-existent DTD file in a file (or one that can't be found over a network), then any attempt to parse it will raise a pyRXP error.

```
>>> fn=open('sample2.xml', 'r').read()

>>> print fn
<?xml version="1.0" encoding="iso-8859-1" standalone="no" ?>
<!DOCTYPE a SYSTEM "nonexistent.dtd">

<a>
<b>This is the contents</b>
</a>

>>> p.parse(fn)
C:\tmp\pyRXP_tests\nonexistent.dtd: No such file or directory
Traceback (most recent call last):
  File "<stdin>", line 1, in ?
pyRXP.Error: Error: Couldn't open dtd entity file:///C:/tmp/pyRXP_tests/nonexistent.dtd
in unnamed entity at line 2 char 38 of [unknown]
```

This is a different kind of error to one where no DTD is specified:

```
>>> fn=open('sample4.xml', 'r').read()

>>> print fn
<?xml version="1.0" encoding="iso-8859-1" standalone="no" ?>
<a>
<b>This is the contents</b>
</a>

>>> p.parse(fn)
Traceback (most recent call last):
  File "<stdin>", line 1, in ?
pyRXP.Error: Error: Document has no DTD, validating abandoned
in unnamed entity at line 3 char 2 of [unknown]
```

If you have errors in your XML and it does not validate against the DTD, you will get a different kind of `pyRXPError`.

```
>>> fn=open('sample3.xml', 'r').read()

>>> print fn
<?xml version="1.0" encoding="iso-8859-1" standalone="no" ?>
<!DOCTYPE a SYSTEM "tinydtd.dtd">

<x>
<b>This is the contents</b>
</x>

>>> p.parse(fn)
Traceback (most recent call last):
  File "<stdin>", line 1, in ?
pyRXP.Error: Error: Start tag for undeclared element x
in unnamed entity at line 4 char 3 of [unknown]
>>>
```

Whether PyRXP validates against a DTD, together with a number of other behaviours is decided by how the various flags are set.

By default, `ErrorOnValidityErrors` is set to 1, as is `NoNoDTDWarning`. If you want the XML you are parsing actually validate against your DTD, you should have both of these set to 1 (which is the default value), otherwise instead of raising a `pyRXP` error XML that doesn't conform to the DTD (which may or may not exist) will be silently ignored. You should also have `Validate` set to 1, otherwise validation won't even be attempted.

Note that the first examples in this chapter - the ones without a DTD - only worked because we had carefully chosen what seem like the sensible defaults. It is set to validate, but not to complain if the DTD is missing. So when you feed it something without a DTD declaration, it notices the DTD is missing but continues in non-validating mode. There are numerous flags set out below which affect the behaviour.

### 3.3 Interface Summary

The python module exports the following:

<code>Error</code>	a python exception
<code>Version</code>	the string version of the module
<code>RXPVersion</code>	the version string of the rxp library embedded in the module
<code>parser_flags</code>	a dictionary of parser flags the values are the defaults for parsers
<code>Parser(*kw)</code>	Create a parser

### 3.4 Parser Object Attributes and Methods

`parse(src)`

We have already seen that this is the main interface to the parser. It returns ReportLab's standard tuple tree representation of the xml source. The string `src` contains the xml.

`srcName`

A name used to refer to the source text in error and warning messages. It is initially set as '<unknown>'. If you know that the data came from "spam.xml" and you want error messages to say so, you can set this to the filename.

`warnCB 0,`

Warning callback. Should either be None, 0, or a callable object (e.g. a function) with a single argument which will receive warning messages. If None is used then warnings are thrown away. If the default 0 value is used then warnings are written to the internal error message buffer and will only be seen if an error occurs.

eoCB

Entity-opening callback. The argument should be None or a callable method with a single argument. This method will be called when external entities are opened. The method should return a (possibly modified) URI. So, you could intercept a declaration referring to *http://some.slow.box/somefile.dtd* and point at the local copy you know you have handy, or implement a DTD-caching scheme.

### 3.5 List of Flags

Flag attributes corresponding to the rxp flags; the values are the module standard defaults. `pyRXP.parser_flags` returns these as a dictionary if you need to refer to these inline.

Flag (1=on, 0=off)	Default
AllowMultipleElements	0
AllowUndeclaredNSAttributes	0
CaseInsensitive	0
ErrorOnBadCharacterEntities	1
ErrorOnUndefinedAttributes	0
ErrorOnUndefinedElements	0
ErrorOnUndefinedEntities	1
ErrorOnUnquotedAttributeValue	1
ErrorOnValidityErrors	1
ExpandCharacterEntities	1
ExpandGeneralEntities	1
IgnoreEntities	0
IgnorePlacementErrors	0
MaintainElementStack	1
MergePCData	1
NoNoDTDWarning	1
NormaliseAttributeValue	1
ProcessDTD	0
RelaxedAny	0
ReturnComments	0
ReturnDefaultedAttributes	1
ReturnNamespaceAttributes	0
SimpleErrorFormat	0
TrustSDD	1
Validate	1
WarnOnRedefinitions	0
XMLExternalIDs	1
XMLLessThan	0
XMLMiscWFErrors	1
XMLNamespaces	0
XMLPredefinedEntities	1
XMLSpace	0
XMLStrictWFErrors	1
XMLSyntax	1

## 3.6 Flag explanations and examples

With so many flags, there is a lot of scope for interaction between them. These interactions are not documented yet, but you should be aware that they exist.

### AllowMultipleElements

Default: 0

Description:

A piece of XML that does not have a single root-tag enclosing all the other tags is described as having multiple elements. By default, this will raise a pyRXP error. Turning this flag on will ignore this and not raise those errors.

Example:

```
>>> p.AllowMultipleElements = 0
>>> p.parse('<a></a><b></b>')
Traceback (most recent call last):
  File "<stdin>", line 1, in ?
pyRXP.Error: Error: Document contains multiple elements
in unnamed entity at line 1 char 9 of [unknown]

>>> p.AllowMultipleElements = 1
>>> p.parse('<a></a><b></b>')
('a', None, [], None)
>>>
```

### AllowUndeclaredNSAttributes

Default: 0

Description:

*[to be added]*

Example:

*[to be added]*

### CaseInsensitive

Default: 0

Description:

This flag controls whether the parse is case sensitive or not.

Example:

```
>>> p.CaseInsensitive=1
>>> p.parse('<a></A>')
('A', None, [], None)

>>> p.CaseInsensitive=0
>>> p.parse('<a></A>')
Traceback (most recent call last):
  File "<stdin>", line 1, in ?
pyRXP.Error: Error: Mismatched end tag: expected </a>, got </A>
in unnamed entity at line 1 char 7 of [unknown]
>>>
```

### ErrorOnBadCharacterEntities

Default: 1

Description:

If this is set, character entities which expand to illegal values are an error, otherwise they are ignored with a warning.

Example:

```
>>> p.ErrorOnBadCharacterEntities=0
>>> p.parse('<a>&#999;</a>')
('a', None, [''], None)

>>> p.parse('<a>&#999;</a>')
```

```
Traceback (most recent call last):
  File "<stdin>", line 1, in ?
pyRXP.Error: Error: 0x3e7 is not a valid 8-bit XML character
in unnamed entity at line 1 char 10 of [unknown]
```

### ErrorOnUndefinedAttributes

Default: 0

Description:

If this is set and there is a DTD, references to undeclared attributes are an error.

See also: ErrorOnUndefinedElements

### ErrorOnUndefinedElements

Default: 0

Description:

If this is set and there is a DTD, references to undeclared elements are an error.

See also: ErrorOnUndefinedAttributes

### ErrorOnUndefinedEntities

Default: 1

Description:

If this is set, undefined general entity references are an error, otherwise a warning is given and a fake entity constructed whose value looks the same as the entity reference.

Example:

```
>>> p.ErrorOnUndefinedEntities=0
>>> p.parse('<a>&dud;</a>')
('a', None, ['&dud;'], None)

>>> p.ErrorOnUndefinedEntities=1
>>> p.parse('<a>&dud;</a>')
Traceback (most recent call last):
  File "<stdin>", line 1, in ?
pyRXP.Error: Error: Undefined entity dud
in unnamed entity at line 1 char 9 of [unknown]
```

### ErrorOnUnquotedAttributeValues

Default: 1

Description:

*[to be added]*

### ErrorOnValidityErrors

Default: 1

Description:

If this is on, validity errors will be reported as errors rather than warnings. This is useful if your program wants to rely on the validity of its input.

### ExpandCharacterEntities

Default: 1

Description:

If this is set, entity references are expanded. If not, the references are treated as text, in which case any text returned that starts with an ampersand must be an entity reference (and provided MergePCData is off, all entity references will be returned as separate pieces).

See also: ExpandGeneralEntities, ErrorOnBadCharacterEntities

Example:

```
>>> p.ExpandCharacterEntities=1
>>> p.parse('<a>&#109;</a>')
('a', None, ['m'], None)
```



```
>>> p.ExpandCharacterEntities=0
>>> p.parse('<a>#109;</a>')
('a', None, ['#109;'], None)
```

### ExpandGeneralEntities

Default: 1

Description:

If this is set, entity references are expanded. If not, the references are treated as text, in which case any text returned that starts with an ampersand must be an entity reference (and provided MergePCData is off, all entity references will be returned as separate pieces).

See also: ExpandCharacterEntities

Example:

```
>>> p.ExpandGeneralEntities=0
>>> p.parse('<a>&lt;</a>')
('a', None, ['&lt;'], None)

>>> p.ExpandGeneralEntities=1
>>> p.parse('<a>&lt;</a>')
('a', None, ['&'], None)
```

### IgnoreEntities

Default: 0

Description:

If this flag is on, normal entity substitution takes place. If it is off, entities are passed through unaltered.

Example:

```
>>> p.IgnoreEntities=0
>>> p.parse('<a>&lt;</a>')
('a', None, ['&'], None)

>>> p.IgnoreEntities=1
>>> p.parse('<a>&lt;</a>')
('a', None, ['&lt;'], None)
```

### IgnorePlacementErrors

Default: 0

Description:

*[to be added]*

### MaintainElementStack

Default: 1

Description:

*[to be added]*

### MergePCData

Default: 1

Description:

If this is set, text data will be merged across comments and entity references.

### NoNoDTDWarning

Default: 1

Description:

Usually, if `validate` is set, the parser will produce a warning if the document has no DTD. This flag suppresses the warning (useful if you want to validate if possible, but not complain if not).

### NormaliseAttributeValues

Default: 1

**Description:**

If this is set, attributes are normalised according to the standard. You might want to not normalise if you are writing something like an editor.

**ProcessDTD**

Default: 0

**Description:**

If `TrustSDD` is set and a `DOCTYPE` declaration is present, the internal part is processed and if the document was not declared standalone or if `Validate` is set the external part is processed. Otherwise, whether the `DOCTYPE` is automatically processed depends on `ProcessDTD`; if `ProcessDTD` is not set the user must call `ParseDtd()` if desired.

See also: `TrustSDD`

**RelaxedAny**

Default: 0

**Description:**

*[to be added]*

**ReturnComments**

Default: 0

**Description:**

If this is set, comments are returned, otherwise they are ignored.

Example:

```
>>> p.ReturnComments = 1
>>> p.parse('<a><!-- this is a comment --></a>')
('a', None, ['<!-- this is a comment -->'], None)

>>> p.ReturnComments = 0
>>> p.parse('<a><!-- this is a comment --></a>')
('a', None, [], None)
```

**ReturnDefaultedAttributes**

Default: 1

**Description:**

If this is set, the returned attributes will include ones defaulted as a result of `ATTLIST` declarations, otherwise missing attributes will not be returned.

**ReturnNamespaceAttributes**

Default: 0

**Description:**

*[to be added]*

**SimpleErrorFormat**

Default: 0

**Description:**

This causes the output on errors to get shorter and more compact.

Example:

```
>>> p.SimpleErrorFormat=0
>>> p.parse('<a>causes an error</b>')
Traceback (most recent call last):
  File "<stdin>", line 1, in ?
pyRXP.Error: Error: Mismatched end tag: expected </a>, got </b>
in unnamed entity at line 1 char 22 of [unknown]

>>> p.SimpleErrorFormat=1
>>> p.parse('<a>causes an error</b>')
Traceback (most recent call last):
```

```
File "<stdin>", line 1, in ?  
pyRXP.Error: [unknown]:1:22: Mismatched end tag: expected </a>, got </b>
```

## TrustSDD

Default: 1

Description:

If `TrustSDD` is set and a `DOCTYPE` declaration is present, the internal part is processed and if the document was not declared standalone or if `Validate` is set the external part is processed. Otherwise, whether the `DOCTYPE` is automatically processed depends on `ProcessDTD`; if `ProcessDTD` is not set the user must call `ParseDtd()` if desired.

See also: `ProcessDTD`

## Validate

Default: 1

Description:

If this is on, the parser will validate the document. If it's off, it won't. It is not usually a good idea to set this to 0.

## WarnOnRedefinitions

Default: 0

Description:

If this is on, a warning is given for redeclared elements, attributes, entities and notations.

## XMLExternalIDs

Default: 1

Description:

*[to be added]*

## XMLLessThan

Default: 0

Description:

*[to be added]*

## XMLMiscWFErrors

Default: 1

Description:

To do with well-formedness errors.

See also: `XMLStrictWFErrors`

## XMLNamespaces

Default: 0

Description:

If this is on, the parser processes namespace declarations (see below). Namespace declarations are *not* returned as part of the list of attributes on an element.

See also: `XMLSpace`

## XMLPredefinedEntities

Default: 1

Description:

If this is on, pyRXP recognises the standard preset XML entities &lt; > " and ' . If this is off, all entities including the standard ones must be declared in the DTD or an error will be raised.

Example:

```
>>> p.XMLPredefinedEntities=1
>>> p.parse('<a>&amp;</a>')
('a', None, ['&'], None)

>>> p.parse('<a>&amp;</a>')
Traceback (most recent call last):
  File "<stdin>", line 1, in ?
pyRXP.Error: Error: Undefined entity amp
in unnamed entity at line 1 char 9 of [unknown]
```

### **XMLSpace**

Default: 0

Description:

If this is on, the parser will keep track of xml:space attributes

See also: XMLNamespaces

### **XMLStrictWFErrors**

Default: 1

Description:

If this is set, various well-formedness errors will be reported as errors rather than warnings.

### **XMLSyntax**

Default: 1

Description:

*[to be added]*

## 4. The examples and utilities

The zip file of examples contains a couple of validatable documents in xml, the samples used in this manual, and two utility modules: one for benchmarking and one for navigating through tuple trees.

### 4.1 Benchmarking

*benchmarks.py* is a script aiming to compare performance of various parsers. We include it to make our results reproducible. It is not a work of art and if you think you can make it fairer or better, tell us how! Here's an example run.

```
C:\code\rlextra\radxml\samples>benchmarks.py

Interactive benchmark suite for Python XML parsers.
Parsers available:

opened sample XML file 444220 bytes long
  1. pyRXP
  2. rparsexml
  3. minidom
  4. msxml30
  5. 4dom
  6. cdomlette

Shall we do memory tests? i.e. you look at Task Manager? y/n y
Test number (or x to exit)>1
testing pyRXP
Pre-parsing: please input python process memory in kb >2904
Post-parsing: please input python process memory in kb >7180
12618 tags, 8157 attributes
pyRXP: init 0.0315, parse 0.3579, traverse 0.1594, mem used 4276kb, mem factor 9.86
```

Instead of the traditional example (hamlet), we took as our example an early version of the Report Markup Language user guide, which is about half a megabyte. Hamlet's XML has almost no attributes; ours contains lots of attributes, many of which will need conversion to numbers one day, and so it provides a more rounded basis for benchmarks

We measure several factors. First there is speed. Obviously this depends on your PC. The script exits after each test so you get a clean process. We measure (a) the time to load the parser and any code it needs into memory (important if doing CGI); (b) time to produce the tree, using whatever the parser natively produces; and (c) time to traverse the tree counting the number of tags and attributes. (c) might be important with a 'very lazy' parser which searched the source text on every request. Also, later on we will be able to look at the difference between traversing a raw tuple tree and some objects with friendlier syntax.

Next is memory. Actually you have to measure that! If anyone can give us the API calls on Windows and other platforms to find out the current process size, we'd be grateful! What we are interested in is how big the structure is in memory. The above shows that the mempry allocated is 9.86 times as big as the original XML text. That sounds a lot, but it's actually much less than most DOM parsers.

By contrast, here's the result for the *minidom* parser included in the official Python distro:

```
minidom: init 0.3039, parse 12.6435, traverse 0.0000, mem used 29136kb, mem factor 67.16
```

Even though minidom uses pyexpat (which is in C) to parse the XML, it's 36 times slower and uses 7 times more memory. And of course it does not validate. Phew!

### 4.2 xmlutils and the TagWrapper

Finally, we've included a 'tag wrapper' class which makes it easy to navigate around the tuple tree. This is randomly selected from many such modules we have used in various projects; the next task for us is to pick ONE and publish it! Essentially, it uses lazy evaluation to try and figure out which part of the XML you want. If you ask for 'tag.spam', it will check if (a) there is an attribute called spam, or (b) if there is a child tag whose tag name is 'spam'. And you can iterate over child nodes as a sequence. And, the str() method of a tag which just contains text is just the text. The snippets below should make it clear what we are doing.

```
>>> tree = pyRXP.Parser().parse(srcText)
>>> srcText = open('rml_a.xml').read()
>>> tree = pyRXP.Parser().parse(srcText)
>>> import xmlutils
>>> tw = xmlutils.TagWrapper(tree)
>>> tw
TagWrapper<document>
>>> tw.filename
'RML_UserGuide_1_0.pdf'
>>> len(tw.story) # how many tags in the story?
1566
>>> tw.template.pageSize
'(595, 842)'

>>> for elem in tw.story:
...     if elem.tagName == 'h1':
...         print elem
...
RML User Guide

Part I - The Basics
Part II - Advanced Features
Part III - Tables
Appendix A - Colors recognized by RML
Appendix B - Glossary of terms and abbreviations
Appendix C - Letters used by the Greek tag
Appendix D - Command reference
Generic Flowables (Story Elements)
Graphical Drawing Operations
Graphical State Change Operations
Style Elements
Page Layout Tags
Special Tags
>>>
```

We are NOT saying this is a particularly good or complete wrapper; but we do intend to standardize on one such wrapper module in the near future, because it makes access to XML information much more 'pythonic' and pleasant. It could be used with tuple trees generated by any parser. Please let us know if you have any suggestions on how it should behave.

## 5. Future Directions

### 5.1 Test Suite

We urgently need a unittest-based suite full of samples saying ‘parse this XML with these flags and assert fact X about the output’. If done right, this could be used to generate the documentation on the parser flags as well. It will be very important when allowing pluggable parsers.

In the meantime, there are some simple tests. Look at the file `test\t.py`.

### 5.2 Standardize the wrapper

A standard wrapper class to let you ‘drill down’ into the tuple tree. This should be as pythonic as possible

### 5.3 Other parsers

Include tuple tree constructors based on other parsers. One could use pyexpat (in fact a few lines could be added to pyexpat itself to produce a tuple tree in some future version of Python). This would be useful for people who cannot install extensions but have Python 2.0 or above. We also have our own parser, Aaron Watters’ `rparsexml`, which uses no C code and is thus useful in places where you cannot build extensions. The latter is not guaranteed to be 100% standards compliant, but this means we can modify it to handle bad XML.

### 5.4 Better benchmark suite

Extend this so that it knows about more parsers and (if possible) can detect the memory used by them without needing to pause and look in Task Manager. Ensure we are being fair to competitors and using their parsers optimally.

### 5.5 Type Conversion Utility

In the parsed output, everything is a string. Yet XML is full of attributes which “mean” numeric values. In particular our own Report Markup Language has numerous attributes like *x*, *y*, *width*, *height*, as well as color attributes. It would be really useful to generalize the conversion step. Let’s say you can provide a mapping like this

1. (tag, attribute) -> reader function
2. attribute -> reader function

Many of the reader functions are just *int* or *float*; others could be written in Python or C. For example we have standard length expressions like “3cm” or “8.5in” which we convert to float values in points. This could say that (a) if this tag name and attribute name has a converter function, use it in-place; (b) if the attribute name has a converter, use that; and if (c) there is nothing specified, leave it as a string.

So the tree could be converted “in place” with a simple API call, at C-like speeds. And we’d be able to remove a lot of code from our application and replace it with a very simple mapping. Expect this Real Soon Now :-)

Note that this type-conversion is not an XML standard. The One True Way is probably to use XML Schema; but for now this would be practically immensely useful, we don’t have a schema-validating parser, and we are big fans of stuff that works now.

### 5.6 Source file references

Debug/trace info: add an extra structure to show the position in the original source file where the tag starts and finished. This would be a parse-time option, as you might not want to take the time and memory. This would let an app raise an error saying not just that the color tag contained a bad color value, but also that it occurred at line 2352 of the input. Useful! This is why we reserved the final tuple element for future use.

## **5.7 (longer term and debatable) Richer tuple tree structure**

It has been suggested that we expand the structure in a couple of ways. Instead of tuples we could make a new C-based node object with a richer model.

Each node should have some pointer back to its parent. This makes navigation a lot easier, but means a little more housekeeping.

We could then also let you distinguish things like CDATA and entity nodes and make it a full rewritable DOM implementation, running at C-like speeds. We could even go further and keep references to things like comments, which are not part of the XML standard.

PyRXP meets our needs already and we won't rush into this. Still, it might be an attractive enhancement for a future version of Python; essentially one would make a lightweight XML node into a built-in type.