

Parsing XML With SimpleXML

(<http://www.sitepoint.com/author/spanda/>)

Parsing XML essentially means navigating through an XML document and returning the relevant data. An increasing number of web services return data in JSON format, but a large number still return XML, so you need to master parsing XML if you really want to consume the full breadth of APIs available.

Using PHP's SimpleXML extension that was introduced back in PHP 5.0, working with XML is very easy to do. In this article I'll show you how.

Basic Usage

Let's start with the following sample as `languages.xml` :

```
1  <?xml version="1.0" encoding="utf-8"?>
2  <languages>
3    <lang name="C">
4      <appeared>1972</appeared>
5      <creator>Dennis Ritchie</creator>
6    </lang>
7    <lang name="PHP">
8      <appeared>1995</appeared>
9      <creator>Rasmus Lerdorf</creator>
10   </lang>
11   <lang name="Java">
12     <appeared>1995</appeared>
13     <creator>James Gosling</creator>
14   </lang>
15 </languages>
```

The above XML document encodes a list of programming languages, giving two details about each language: its year of implementation and the name of its creator.

The first step is to loading the XML using either `simplexml_load_file()` or `simplexml_load_string()` . As you might expect, the former will load the XML file a file and the later will load the XML from a given string.

```
1 <?php
2 $languages = simplexml_load_file("languages.xml");
```

Both functions read the entire DOM tree into memory and returns a `SimpleXMLElement` object representation of it. In the above example, the object is stored into the `$ languages` variable. You can then use `var_dump()` or `print_r()` to get the details of the returned object if you like.

SimpleXMLElement Object

```
(
  [lang] => Array
    (
      [0] => SimpleXMLElement Object
        (
          [@attributes] => Array
            (
              [name] => C
            )
          [appeared] => 1972
          [creator] => Dennis Ritchie
        )
      [1] => SimpleXMLElement Object
        (
          [@attributes] => Array
            (
              [name] => PHP
            )
          [appeared] => 1995
          [creator] => Rasmus Lerdorf
        )
      [2] => SimpleXMLElement Object
        (
          [@attributes] => Array
            (
              [name] => Java
            )
          [appeared] => 1995
          [creator] => James Gosling
        )
    )
)
```

The XML contained a root `language` element which wrapped three `lang` elements, which is why the `SimpleXMLElement` has the public property `lang` which is an array of three `SimpleXMLElements`. Each element of the array corresponds to a `lang` element in the XML document.

You can access the properties of the object in the usual way with the `->` operator. For example, `$languages->lang[0]` will give you a `SimpleXMLElement` object which corresponds to the first `lang` element. This object then has two public properties: `appeared` and `creator`.

```
1 <?php
2 $languages->lang[0]->appeared;
3 $languages->lang[0]->creator;
```

Iterating through the list of languages and showing their details can be done very easily with standard looping methods, such as `foreach`.

```
1 <?php
2 foreach ($languages->lang as $lang) {
3     printf(
4         "<p>%s appeared in %d and was created by %s.</p>",
5         $lang["name"],
6         $lang->appeared,
7         $lang->creator
8     );
9 }
```

Notice that I accessed the `lang` element's `name` attribute to retrieve the name of the language. You can access any attribute of an element represented as a `SimpleXMLElement` object using array notation like this.

Dealing With Namespaces

Many times you'll encounter namespaced elements while working with XML from different web services. Let's modify our `languages.xml` example to reflect the usage of namespaces:

```
1 <?xml version="1.0" encoding="utf-8"?>
2 <languages
3     xmlns:dc="http://purl.org/dc/elements/1.1/" (http://purl.org/dc/elements/1.1/)">
4     <lang name="C">
5         <appeared>1972</appeared>
6         <dc:creator>Dennis Ritchie</dc:creator>
```

```
7     </lang>
8     <lang name="PHP">
9         <appeared>1995</appeared>
10        <dc:creator>Rasmus Lerdorf</dc:creator>
11    </lang>
12    <lang name="Java">
13        <appeared>1995</appeared>
14        <dc:creator>James Gosling</dc:creator>
15    </lang>
16 </languages>
```

Now the `creator` element is placed under the namespace `dc` which points to <http://purl.org/dc/elements/1.1/>. If you try to print the creator of a language using our previous technique, it won't work. In order to read namespaced elements like this you need to use one of the following approaches.

The first approach is to use the namespace URI directly in your code when accessing namespaced elements. The following example demonstrates how:

```
1  <?php
2  $dc = $languages->lang[1]->children("http://purl.org/dc/elements/1.1/ (http://purl
3  echo $dc->creator;
```



The `children()` method takes a namespace and returns the children of the element that are prefixed with it. It accepts two arguments; the first one is the XML namespace and the latter is an optional Boolean which defaults to false. If you pass true, the namespace will be treated as a prefix rather than the actual namespace URI.

The second approach is to read the namespace URI from the document and use it while accessing namespaced elements. This is actually a cleaner way of accessing elements because you don't have to hardcode the URI.

```
1  <?php
2  $namespaces = $languages->getNamespaces(true);
3  $dc = $languages->lang[1]->children($namespaces["dc"]);
4
5  echo $dc->creator;
```

The `getNamespaces()` method returns an array of namespace prefixes with their associated URIs. It accepts an optional parameter which defaults to false. If you set it true then the method will return the namespaces used in parent and child nodes. Otherwise, it finds namespaces used within the parent node only.

Now you can iterate through the list of languages like so:

```
1  <?php
2  $languages = simplexml_load_file("languages.xml");
3  $ns = $languages->getNamespaces(true);
4
5  foreach($languages->lang as $lang) {
6      $dc = $lang->children($ns["dc"]);
7      printf(
8          "<p>%s appeared in %d and was created by %s.</p>",
9          $lang["name"],
10         $lang->appeared,
11         $dc->creator
12     );
13 }
```

A Practical Example – Parsing YouTube Video Feed

Let's walk through an example that retrieves the RSS feed from a YouTube channel displays links to all of the videos from it. For this we need to make a call to the following URL:

<http://gdata.youtube.com/feeds/api/users//uploads>

The URL returns a list of the latest videos from the given channel in XML format. We'll parse the XML and get the following pieces of information for each video:

- Video URL
- Thumbnail
- Title

We'll start out by retrieving and loading the XML:

```
1  <?php
```

```
2 $channel = "channelName";
3 $url = "http://gdata.youtube.com/feeds/api/users/_ (http://gdata.youtube.com/feeds/a
4 $xml = file_get_contents($url);
5
6 $feed = simplexml_load_string($xml);
7 $ns=$feed->getNamespaces(true);
```

If you take a look at the XML feed you can see there are several `entry` elements each of which stores the details of a specific video from the channel. But we are concerned with only thumbnail image, video URL, and title. The three elements are children of `group`, which is a child of `entry`:

```
1 <entry>
2   ...
3   <media:group>
4     ...
5     <media:player url="video url"/>
6     <media:thumbnail url="video url" height="height" width="width"/>
7     <media:title type="plain">Title...</media:title>
8     ...
9   </media:group>
10  ...
11 </entry>
```

We simply loop through all the `entry` elements, and for each one we can extract the relevant information. Note that `player`, `thumbnail`, and `title` are all under the `media` namespace. So, we need to proceed like the earlier example. We get the namespaces from the document and use the namespace while accessing the elements.

```
1 <?php
2 foreach ($feed->entry as $entry) {
3     $group=$entry->children($ns["media"]);
4     $group=$group->group;
5     $thumbnail_attrs=$group->thumbnail[1]->attributes();
6     $image=$thumbnail_attrs["url"];
7     $player=$group->player->attributes();
8     $link=$player["url"];
9     $title=$group->title;
10    printf('<p><a href="%s"></a></p>',
```

```
11     $player, $image, $title);  
12 }
```

Conclusion

Now that you know how to use SimpleXML to parse XML data, you can improve your skills by parsing different XML feeds from various APIs. But an important point to consider is that SimpleXML reads the entire DOM into memory, so if you are parsing large data sets then you may face memory issues. In those cases it's advisable to use something other than SimpleXML, preferably an event-based parser such as XML Parser. To learn more about SimpleXML, check out its [documentation \(http://php.net/manual/en/book.simplexml.php\)](http://php.net/manual/en/book.simplexml.php).