_____

# XML PROCESSING

XML is an important technology by anyone's standards. It powers many of the modern technologies and you have to know how to process XML if you want to work with software of any value in modern industry. In this chapter, we are going to talk about three different methods of parsing XML and turning it into a format that our programs can reason with. We begin with the most rudimentary approach.

## 6.1 Simple API for XML

Simple API for XML (or *SAX*) is the most basic (and oftentimes, most efficient) way of processing XML. It is an *event-based parser* – it starts processing the XML from top-down and generates events as it reads different portions of the file. Examples of events are `START_TAG` and `END_TAG`. Based on the different events (and their sequences), we can decide exactly what data is coming in to us. Let's take a look at an example to understand SAX.

*SAX*

*event-based parser*

### 6.1.1 A SAX Example

Let's consider an XML file that describes a mathematical expression. A basic version of this file (which we will save as `data/expression.xml`) looks like this:

```xml
<?xml version="1.0"?>
<expression>
  <expr func="+">
    <expr func="*">
      <const val="2" />
      <const val="3" />
    </expr>
    <expr func="-">
      <const val="4" />
      <const val="5" />
    </expr>
    <const val="6" />
  </expr>
</expression>
```

This file corresponds to the expression $(2*3)+(4-5)+6$. It's very verbose and if we want to do something useful with this file (such as calculate the result), we need to be able to represent it in some meaningful Java object format.

The first step in doing that would be to define a class `Expresion` in our project.

```java
package org.csrdu.java.ex.xml.expression.model;

import java.util.Vector;

public class Expression {
  public String function = "";
  public Vector operands;

```

```
 9    public Expression() {
10      operands = new Vector();
11    }
12
13    public void printExpression(int level) {
14        System.out.println(nSpaces(level * 2) + "Function:  [" + function + "]");
15
16        for (Object i : operands) {
17          if (i instanceof Integer)
18            System.out.println(nSpaces(level * 2 + 2) + i);
19          else
20            ((Expression)i).printExpression(level+1);
21        }
22    }
23
24    public String nSpaces(int length) {
25      StringBuffer outputBuffer = new StringBuffer(length);
26      for (int i = 0; i < length; i++) {
27        outputBuffer.append(" ");
28      }
29      return outputBuffer.toString();
30    }
31 }
```

The class itself is very simple. It has a field called `function` to store the function name of the expression and a `Vector` of `Expression` objects as the operands. This is needed as the operands can be both constants and other expressions. In case one of the operands is a constant, the `function` field is ignored and the only `operand` will be an `Integer`.

The `printExpression()` function simply outputs the expression in a properly indented format (the indentation coming from the `nSpaces()` function.

*model*
*controller*

Now that we have a *model*, we can define a *controller* for it that reads the data from the XML file and creates an instance of the expression in Java object code.

```
 1 package org.csrdu.java.ex.xml.expression.controller;
 2
 3 import org.csrdu.java.ex.xml.expression.model.Expression;
 4 import org.xml.sax.*;
 5 import org.xml.sax.helpers.*;
 6 import java.io.*;
 7 import java.util.*;
 8
 9 public class ExpressionParser extends DefaultHandler {
10
11    Stack<Expression> st;
12
13    Expression curExpr;
14
15    public ExpressionParser() {
16      st = new Stack<Expression>();
17    }
18
19    // Override methods of the DefaultHandler class
20    // to gain notification of SAX Events.
21    //
22    // See org.xml.sax.ContentHandler for all available events.
23    //
24    public void startElement(String namespaceURI, String localName,
25    String qName, Attributes attr) throws SAXException {
26      if (localName.equals("expr")) {
```

```
27        // we have a new expression.
28        // push the old one on the stack (if any)
29        if (curExpr != null)
30        st.push(curExpr);
31
32        // We need to create a new expression.
33        curExpr = new Expression();
34
35        // and set the function for this expression
36        curExpr.function = attr.getValue("func");
37
38    } else if (localName.equals("const")) {
39        // we have a constant. Let's add it to the expression
40        // first get the value of the constant
41        Integer val = Integer.parseInt(attr.getValue("val"));
42        curExpr.operands.add(val);
43    }
44  }
45
46  public void endElement(String namespaceURI, String localName, String qName)
47    throws SAXException {
48    if (localName.equals("expr")) {
49      // End of current expression
50      // save the expression just ended
51      Expression endedExpr = curExpr;
52      // pop the current expression from the stack
53      if (!st.empty()) {
54        curExpr = st.pop();
55        // add ended expression to parent
56        curExpr.operands.add(endedExpr);
57      }
58    }
59    // we can ignore the end of const tag
60  }
61
62  public static void main(String[] argv) {
63    try {
64        // Create SAX parser...
65        XMLReader xr = XMLReaderFactory.createXMLReader();
66
67        // Set the ContentHandler...
68        ExpressionParser exprs = new ExpressionParser();
69        xr.setContentHandler(exprs);
70
71        // Parse the file...
72        xr.parse(new InputSource(new FileReader("data/expression.xml")));
73        exprs.curExpr.printExpression(0);
74
75    } catch (Exception e) {
76        e.printStackTrace();
77    }
78  }
79 }
```

Let's begin the explanation by describing the `main()` function. It sets up the `XMLReader` which is simply a SAX *parser*. We then create an instance of our *parser* `ExpressionParser` class and parse it to the XMLReader as the *handler*. This means *handler* that whenever the parser finds an "important piece" of the XMl, it will generate an event which will be handled by our handler. This allows our code to be notified whenever an XML event (such as start of tag) occurs.

---
**Algorithm 1** The `startElement` event
---
1: **if** event source = expr **then**
2:     push old expression on stack (if any)
3:     create a new expression object
4:     set function of new expression to attribute 'func'
5: **else if** event source = const **then**
6:     add attribute 'val' as operand to current expression
7: **end if**
---

---
**Algorithm 2** The `endElement` event
---
1: **if** event source = expr **then**
2:     save the expression just ended as endedExpr
3:     **if** stack is not empty **then**
4:         pop the expression at the top of stack as curExpr
5:         add endedExpr as an operand of curExpr
6:     **end if**
7: **end if**
---

*Plain Old Java Object (POJO)*

We then parse the XML file (`data/expression.xml`). This will convert our XML file into a *Plain Old Java Object (POJO)*. We can then print it out or do whatever we need to do with it.

All we have to understand now is how our `ExpressionHandler` handles the events. This is covered in two functions:

**startElement()**

This function defines the majority of our logic. The algorithm for handling the *startElement* event is given in Algorithm 1. We first check if the event was raised for `expr` tag. If so, we know that we are starting a new expression. So, we push the old one on the stack (if one exists) and create a new expression tag. Any constant we see from now on will be an operand of this new expression. We also get the `func` attribute of the expression and set the function name equal to it.

If, on the other hand, we have this event raised for `const`, we know that we got a constant and we can take it's `val` attribute and set it as an operand of the current expression.

The other half of this logic comes into place when we look at what happens when a tag ends.

**endElement()**

When the `const` tag ends, we don't care because it's can't be a parent tag. What we are interested in is when the `expr` tag ends (cf. Algorithm 2). When that happens, we know that the current expression has ended. We need to pop the (parent) expression that we pushed on the stack and then add the expression that just ended to the parent. And that takes care of the whole logic!

If you have trouble understanding, Figure 6.1 might be of help. Refer to this in conjunction with the algorithm and the code.
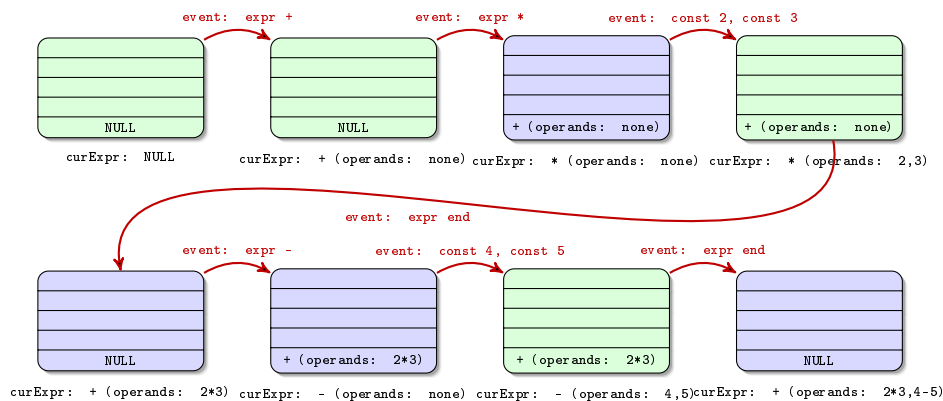
Figure 6.1: SAX Events Dry Run

### 6.1.2 Limitations of SAX

SAX is often very efficient but it has a major limitation. It does not allow random access. You have to parse the whole file to get to a particular piece of information. This becomes very problematic if you have a complicated file and you want to retrieve a particular piece of information from it. One way of solving this problem is through the use of a tree model of XML parsing. We discuss one of these below.

## 6.2 Document Object Model

Document Object Model (more commonly known as *DOM*) is a *tree-based parser* for XML. It allows you to traverse the XML document as if it were a tree (which it really is). The best way to understand DOM is by looking at a simple example. We are going to parse the same expression file that we parsed using SAX. Let's first take a look at the code.

*DOM*

*tree-based parser*

```java
1  package org.csrdu.java.ex.xml.expression.controller;
2
3  import javax.xml.parsers.DocumentBuilderFactory;
4  /* some imports hidden */
5
6  public class ExpressionParserDom {
7
8    public static void main(String argv[]) {
9      try {
10       File file = new File("data/expression.xml");
11       DocumentBuilderFactory dbFactory = DocumentBuilderFactory
12       .newInstance();
13       DocumentBuilder dBuilder = dbFactory.newDocumentBuilder();
14       Document doc = dBuilder.parse(file);
15       doc.getDocumentElement().normalize();
16
17       Element rootElement = doc.getDocumentElement();
18       NodeList nodes = rootElement.getChildNodes().item(1).getChildNodes();
19
20       for (int i = 0; i < nodes.getLength(); i++) {
21         Node node = nodes.item(i);
```

```
22
23            if (node instanceof Element) {
24              // a child element to process
25              Element child = (Element) node;
26              if (child.getTagName() == "expr") {
27                String attribute = child.getAttribute("func");
28                System.out.println("Expression with function ID "
29                    + attribute + " found.");
30              } else if (child.getTagName() == "const"){
31                String attribute = child.getAttribute("val");
32                System.out.println("Constant with value "
33                    + attribute + " found.");
34              }
35            }
36          }
37        } catch (Exception e) {
38          e.printStackTrace();
39        }
40    }
41 }
```

*document builder factory*     First, we create what is called a *document builder factory*. This object can build a document by parsing a file (or string). We then normalize the document to get it into a standardized format. After that, it's fairly simple. We can get the root element of the XML file (using the `getDocumentElement()` function). This returns *nodes* an object of type `Node`. All tags and text in DOM are *nodes*. When we have a node object, we can get its child nodes using the `getChildNodes()` function. This *iterable* returns an *iterable* collection. We can loop over this collection and use the different functions.

## 6.3   XML Pull Parser

One of the more recent parsers has been proposed in the form of *XML Pull Parser* *XML Pull Parser API* *API* [16]. It's a fast implementation that works (at least on the face), similar to SAX. Instead of relying on the push parsing of SAX, the XML Pull Parsers actively poll the document for events. This leads to a faster parsing performance. The XML Pull API is used extensively in the Android framework. However, an implementation of this API is not available in vanilla JREs. In this example, we will use the reference implementation[1] of this API provided by the kXML project [8].

Below, we provide code for parsing our expression XML through this pull parser.

```
1 package org.csrdu.java.ex.xml.expression.controller;
2
3 import java.io.FileReader;
4 import java.io.IOException;
5 import org.xmlpull.v1.XmlPullParser;
6 import org.xmlpull.v1.XmlPullParserException;
7 import org.xmlpull.v1.XmlPullParserFactory;
8
9 /**
10 * An example of an application that uses XMLPULL V1 API.
11 *
12 * @author <a href="http://www.extreme.indiana.edu/~aslom/">Aleksander
13 *          Slominski</a>
```

---

[1]You can download the required JAR file from the lecture server or from the kXML site given in the references.

```
14 */
15 public class ExpressionParserXMLPull {
16   public static void main(String args[])
17     throws XmlPullParserException, IOException {
18     XmlPullParserFactory factory = XmlPullParserFactory.newInstance();
19     factory.setNamespaceAware(true);
20     XmlPullParser xpp = factory.newPullParser();
21     System.out.println("Parser implementation class is "
22         + xpp.getClass());
23
24     ExpressionParserXMLPull app = new ExpressionParserXMLPull();
25
26     System.out.println("Parsing file. ");
27     xpp.setInput(new FileReader("data/expression.xml"));
28     app.processDocument(xpp);
29   }
30
31   public void processDocument(XmlPullParser xpp)
32     throws XmlPullParserException, IOException {
33     int eventType = xpp.getEventType();
34     do {
35       if (eventType == xpp.START_DOCUMENT) {
36         System.out.println("Start document");
37       } else if (eventType == xpp.END_DOCUMENT) {
38         System.out.println("End document");
39       } else if (eventType == xpp.START_TAG) {
40         processStartElement(xpp);
41       } else if (eventType == xpp.END_TAG) {
42         processEndElement(xpp);
43       } else if (eventType == xpp.TEXT) {
44         processText(xpp);
45       }
46       eventType = xpp.next();
47     } while (eventType != xpp.END_DOCUMENT);
48   }
49
50   public void processStartElement(XmlPullParser xpp) {
51     String name = xpp.getName();
52     System.out.println("Start element: " + name);
53   }
54
55   public void processEndElement(XmlPullParser xpp) {
56     String name = xpp.getName();
57     System.out.println("End element: " + name);
58   }
59
60   int holderForStartAndLength[] = new int[2];
61
62   public void processText(XmlPullParser xpp) throws XmlPullParserException {
63     char ch[] = xpp.getTextCharacters(holderForStartAndLength);
64     int start = holderForStartAndLength[0];
65     int length = holderForStartAndLength[1];
66     System.out.print("Characters:    \"");
67     for (int i = start; i < start + length; i++) {
68       if (ch[i] != '\n' && ch[i] != '\t')
69       System.out.println(ch[i]);
70     }
71     System.out.print("\"\n");
72   }
73 }
```

The developer-facing semantics of XML Pull Parser are approximately the same

as SAX. We first create an instance of the parser. This is done using the *factory design pattern*. We can then use our handler to parse this file. The most important line in our handler is the line `xpp.next()` (Line 46). It pulls the next event (hence the name *pull* parser). We can then decide which even it is and call the relevant function based on that. The only difference here from the SAX parser is that we're testing for basic text as well. This is the portion of characters that fall within the element starting and ending tags. This code can be modified to use the stack-based semantics we used in our SAX example.

*factory design pattern*