

**www.ibm.com/ XML**

Education: [Papers](#)

Building an XML application, step 1: Writing a DTD

Doug Tidwell

IBM XML Technical Strategy Group, TaskGuide Development

Updated January 1999

[Step 2: Generating XML from a Data Store](#)

[Step 3: Converting XML into HTML with the Document Object Model \(DOM\)](#)

Abstract:

One of the main tasks in creating an XML application is writing a Document Type Definition (DTD). The DTD lets us define the different pieces of data we plan to model, along with the relationships between them. The ability to include this semantic information is the source of XML's power, and its main advantage over HTML. In this paper, we'll build a DTD as the first step in building an XML application; future papers will expand on this work.

Mary, Mary, quite contrary, how does your data grow?

When writing a DTD, the first place to start is with your data in its current form. How are your data items arranged currently? Are they in a relational database? In a flat file? On yellow sticky notes plastered to your wall? If you don't have much data, or your data is in an easily manipulatable format, you may be able to restructure your data source to make the job of writing your XML application easier. On the other hand, if you have lots of data, or your current format is inconvenient, you'll probably have to use your data as is.

Sample database structure

For our examples in this paper, we'll create an XML application for the Xtreme Travel Agency, a company specializing in outdoor tours for the active/foolhardy set. This paper involves creating a DTD for the data specified in the `flights` database. The structure of the database is shown in the following table:

Structure of Database Table <code>flights</code>			
Column Name	Sample Data	Column Name	Sample Data
<code>id</code>	0001	<code>DepartFrom_1</code>	Chicago
<code>DepartFrom_2</code>	Palm Springs	<code>DepartTime_1</code>	January 10 1999 6:30 AM
<code>DepartTime_2</code>	January 15 1999 11:50 AM	<code>ArriveIn_1</code>	Palm Springs
<code>ArriveIn_2</code>	Chicago	<code>ArriveTime_1</code>	January 10 1999 11:03 AM
<code>ArriveTime_2</code>	January 15 1999 9:24 PM	<code>Airline_1</code>	American #303
<code>Airline_2</code>	American #1250		
Note: The <code>id</code> field is ignored in this and all further papers about this sample database. The field simply provides a unique key for each record in the database.			

Taking a look at our database, there are several troubling things. One is that several fields contain more than one piece of information. For example, the `DepartTime_1` field contains both the date and time of the departing flight. If we want to look at the starting time and ending time of the flight to calculate the duration of the flight, there's no reliable way to do that. There also seems to be some redundant information (`DepartFrom_2` and `ArriveIn_1` always have the same value, for example).

One of the first questions we would ask is whether any of these concerns affect the data we want to model in our application. If so, it's worth finding out if we can change the structure of the database. Most likely, however, changes to the database won't be allowed, and we'll have to live with what we have.

Modeling Our Data

Assuming that we can't change the database, our first step will be to create a simple DTD that mirrors the structure of the database. We've already listed the fields that make up each record in the `flights` table. There are some other rules we can state about the database:

- The table we're concerned with is called `flights`.
- Each record in the `flights` database represents a complete itinerary: an outbound flight and a returning flight.
- For each itinerary, there are ten fields, the names of which are listed in the previous table.

As a first pass at a DTD, we'll create the tags `<flights>`, `<itinerary>`, etc., specifying the relationships among items we just outlined. Before we do that, we'll discuss the basics of DTD syntax.

DTD Basics

Each statement in a DTD uses the `<XML DTD>` syntax. This syntax begins each instruction with a left angle bracket and an exclamation point, and ends it with a right angle bracket. (The fact that an XML DTD isn't specified in XML is being addressed by several proposals, including the [Document Content Description](#) specification from IBM, Textuality, and Microsoft.) As we mentioned earlier, our first pass at a tag set will look like this:

```
<flights>
<itinerary>
  <departFrom_1 />
  <departTime_1 />
```

```

<arriveIn_1 />
<arriveTime_1 />
<airline_1 />
<departFrom_2 />
<departTime_2 />
<arriveIn_2 />
<arriveTime_2 />
<airline_2 />
</itinerary>
</flights>

```

Defining a Document Element

Our document element, the outermost tag, will be the `<flights>` tag:

```
<!ELEMENT flights (itinerary)+>
```

The element declaration defines the name of the tag (`flights`, in this case), and the *content model* for the tag. The `+` notation above means the `<flights>` tag must contain one or more `<itinerary>` tags.

XML Occurrence Indicators

In addition to the plus sign from our previous example, there are other occurrence indicators:

XML Occurrence Indicator	
Indicator	Meaning
?	The content must appear either once, or not at all.
*	The content can appear one or more times, or not at all.
+	The content must appear at least once, and may appear more than once.
[none]	The content must appear once, exactly as described.

These indicators can be combined with parentheses in any order to create complex expressions. For example, if element `x` is defined with the content model

```
<!ELEMENT x (a, (b | c | d), e)*>
```

all of the following are valid:

```

<x>
  <a />
  <b />
  <e />
</x>

<x>
  <a />
  <d />
  <e />
</x>

<x>
  <a />
  <c />
  <e />
  <a />
  <c />
  <e />
</x>

<x>
  <a />
  <b />
  <e />
  <a />
  <c />
  <e />
  <a />
  <d />
  <e />
</x>

<x>
  </x>

```

Defining More Tags

We've already decided that for our first DTD, we'll simply describe the format of the database. That means we'll define a tag for each field in the row. We'll also rename some of the tags to better reflect

their content. Because rows are represented by the `<itinerary>` tag, we'll include all of the field tags:

```
<!ELEMENT itinerary (outbound-depart-from, outbound-depart-time,
outbound-arrive-in, outbound-arrive-time, outbound-airline,
returning-depart-from, returning-depart-time, returning-arrive-in,
returning-arrive-time, returning-airline)>
```

Notice that the definition of the `itinerary` element doesn't use any occurrence indicators; this means that the elements must occur in exactly this order, and will occur only once.

Now that we've defined the tag that contains the data for each row in the database, we can start defining the individual tags. Each of these ten tags will look like this:

```
<!ELEMENT outbound-depart-from    (#PCDATA)>
<!ELEMENT outbound-depart-time    (#PCDATA)>
...
```

The `#PCDATA` keyword above means that the tag contains parsed character data; this means that the XML parser will find only character data, no tags or entity references (more about these in a minute). There are other keywords, such as `EMPTY`, which means the tag can't contain anything, and `ANY`, which means the tag can contain text, other tags, entity references, etc.

Our DTD - Version 1

Our completed DTD is shown below.

```
<!-- flights.dtd -->
<!ELEMENT flights (itinerary)+>

<!ELEMENT itinerary (outbound-depart-from, outbound-depart-time,
outbound-arrive-in, outbound-arrive-time, outbound-airline,
returning-depart-from, returning-depart-time, returning-arrive-in,
returning-arrive-time, returning-airline)>

<!ELEMENT outbound-depart-from    (#PCDATA)>
<!ELEMENT outbound-depart-time    (#PCDATA)>
<!ELEMENT outbound-arrive-in      (#PCDATA)>
<!ELEMENT outbound-arrive-time    (#PCDATA)>
<!ELEMENT outbound-airline        (#PCDATA)>
<!ELEMENT returning-depart-from    (#PCDATA)>
<!ELEMENT returning-depart-time    (#PCDATA)>
<!ELEMENT returning-arrive-in      (#PCDATA)>
<!ELEMENT returning-arrive-time    (#PCDATA)>
<!ELEMENT returning-airline        (#PCDATA)>
```

Other Things You Can Put in a DTD

There are a number of other things you can put in a DTD; the most common are attribute declarations and entity references.

Attribute Declarations

Attribute declarations allow you to define the attributes that can appear inside a tag, as well as the kinds of data the attributes can contain.

As an example, let's say that we don't like the structure of the `<outbound-airline>` and `<returning-airline>` tags. These tags typically contain data about the airline and the flight number. We've decided that we can reliably parse out the airline and the flight number; the airline number will be the text of the tag, and the flight number will be an attribute inside the tag itself. Here are the definitions for the new tags and their attributes:

```

<!ELEMENT outbound-airline (#PCDATA)>
<!--ATTLIST outbound-airline flightNum CDATA #REQUIRED-->
<!ELEMENT returning-airline (#PCDATA)>
<!--ATTLIST returning-airline flightNum CDATA #REQUIRED-->

```

The #REQUIRED keyword in the attribute definition means that this attribute must be coded for each and every <outbound-airline> and <returning-airline> tag in your document. If an attribute isn't required, you can use the #IMPLIED keyword.

Another type of attribute definition allows you to specify a set of valid values, along with a default. As an example, let's say Xtreme Travel only deals with several airlines, and we've decided that the airline name should be an attribute of the <outbound-airline> and <returning-airline> tags. Here are the definitions for the new tags:

```

<!ELEMENT outbound-airline (EMPTY)>
<!--ATTLIST outbound-airline flightNum CDATA #REQUIRED
carrierName (Alitalia | American | Delta |
Northwest | Pacific |
TWA | United) "American"-->
<!ELEMENT returning-airline (EMPTY)>
<!--ATTLIST returning-airline flightNum CDATA #REQUIRED
carrierName (Alitalia | American | Delta |
Northwest | Pacific |
TWA | United) "American"-->

```

In this somewhat ill-conceived example, the attribute carrierName can have only certain values, all of which are listed in the attribute definition. If no value is specified, the default is American. Also notice that because all of the data contained in these tags is now in the attributes, we used the EMPTY keyword to specify that these tags don't have any content.

Tags or Attributes?

One common question in DTD writing is whether something should be a tag or an attribute. A third approach to our current example would be to create <flightNumber> and <carrierName> tags inside the <outbound-airline> and <returning-airline> tags:

```

<outbound-airline>
  <flightNumber>330</flightNumber>
  <carrierName>American</carrierName>
</outbound-airline>

```

In most cases, the tags versus attributes decision doesn't make any difference. However, if the data we're modelling needs to be reused, data in tags is easier to access. As an example, say the flight number returned by a database query needs to be used as input into another query. Finding and reusing a <flightNumber> tag is much easier than finding and reusing the flightNumber attribute of the <outbound-airline> tag.

Delivering Data in a More Useful Format

As we mentioned earlier, one problem with the underlying database is that it stores the dates and times of the flights as simple text strings. To get around this problem, we will make sure that text is a valid date, then send the parts of that date as XML attributes. This makes it easy for anyone parsing our XML-tagged data to determine exactly what date and time are represented by this markup. Consider these two examples:

```

<outbound-arrive-time>January 10 1999 1:03 PM</outbound-arrive-time>
<outbound-arrive-time year="1999" month="1" day="10" hour="13" minute="3" />

```

The second example has two major advantages:

1. It is very easy to get any component of the date and time.

2. Converting from one format to another (from "January 10 1999 1:03 PM" to "13:03 10 January 1999," for example) is very easy as well.

For now, we'll simply define the attributes required by this markup design; our next topic will discuss parsing the existing data to create the tags. To further simplify our design, we'll use values of the `hour` attribute from 0 to 23 to indicate all the hours of the day, rather than including an `am-pm` attribute. If the date format we're using needs an AM or PM, we'll be able to determine that from the value of the `hour` attribute.

Our new declaration for the time-related tags looks like this:

```
<!ELEMENT outbound-depart-time (EMPTY)>
<!ATTLIST outbound-depart-time  year  CDATA #REQUIRED
                                month CDATA #REQUIRED
                                day   CDATA #REQUIRED
                                hour  CDATA #REQUIRED
                                minute CDATA #REQUIRED>
<!ELEMENT outbound-arrive-time (EMPTY)>
<!ATTLIST outbound-arrive-time  year  CDATA #REQUIRED
                                month CDATA #REQUIRED
                                day   CDATA #REQUIRED
                                hour  CDATA #REQUIRED
                                minute CDATA #REQUIRED>
<!ELEMENT returning-depart-time (EMPTY)>
<!ATTLIST returning-depart-time year  CDATA #REQUIRED
                                month CDATA #REQUIRED
                                day   CDATA #REQUIRED
                                hour  CDATA #REQUIRED
                                minute CDATA #REQUIRED>
<!ELEMENT returning-arrive-time (EMPTY)>
<!ATTLIST returning-arrive-time year  CDATA #REQUIRED
                                month CDATA #REQUIRED
                                day   CDATA #REQUIRED
                                hour  CDATA #REQUIRED
                                minute CDATA #REQUIRED>
```

Entity Declarations

The last thing we'll add to our DTD is an entity declaration. Entities allow you to define symbols that are replaced by other text before they're displayed to the user. Here's an entity that defines the symbol `&xt;` as equivalent to the name "Xtreme Travel."

```
<!ENTITY xt "Xtreme Travel">
```

Markup such as `Welcome to &xt;!` will be rendered as `Welcome to Xtreme Travel!` If you use an entity for a common word or phrase, such as a product name, you can change all occurrences of that word or phrase simply by changing the entity declaration.

Our Final DTD

Here's our final DTD:

```
<!-- flights.dtd -->
<!ELEMENT flights (itinerary)+>
<!ELEMENT itinerary (outbound-depart-from, outbound-depart-time,
                    outbound-arrive-in, outbound-arrive-time,
                    outbound-airline, returning-depart-from,
                    returning-depart-time, returning-arrive-in,
                    returning-arrive-time, returning-airline)>
<!ELEMENT outbound-depart-from (#PCDATA)>
<!ELEMENT outbound-depart-time (EMPTY)>
<!ATTLIST outbound-depart-time  year  CDATA #REQUIRED
                                month CDATA #REQUIRED
```

```

        day      CDATA #REQUIRED
        hour      CDATA #REQUIRED
        minute     CDATA #REQUIRED>
<!ELEMENT outbound-arrive-in    (#PCDATA)>
<!ELEMENT outbound-arrive-time  (EMPTY)>
<!--ATTLIST outbound-arrive-time  year      CDATA #REQUIRED
        month     CDATA #REQUIRED
        day        CDATA #REQUIRED
        hour        CDATA #REQUIRED
        minute     CDATA #REQUIRED>
<!ELEMENT outbound-airline      (EMPTY)>
<!--ATTLIST outbound-airline      flightNum CDATA #REQUIRED
        carrierName (Alitalia | American | Delta |
        Northwest | Pacific |
        TWA | United) "American">
<!ELEMENT returning-depart-from  (#PCDATA)>
<!ELEMENT returning-depart-time  (EMPTY)>
<!--ATTLIST returning-depart-time  year      CDATA #REQUIRED
        month     CDATA #REQUIRED
        day        CDATA #REQUIRED
        hour        CDATA #REQUIRED
        minute     CDATA #REQUIRED>
<!ELEMENT returning-arrive-in    (#PCDATA)>
<!ELEMENT returning-arrive-time  (EMPTY)>
<!--ATTLIST returning-arrive-time  year      CDATA #REQUIRED
        month     CDATA #REQUIRED
        day        CDATA #REQUIRED
        hour        CDATA #REQUIRED
        minute     CDATA #REQUIRED>
<!ELEMENT returning-airline      (EMPTY)>
<!--ATTLIST returning-airline      flightNum CDATA #REQUIRED
        carrierName (Alitalia | American | Delta |
        Northwest | Pacific |
        TWA | United) "American">

<!--ENTITY xt "Xtreme Travel">

```

Sample Files

To study this code, see the html version of this file on the XML web site.

[flights.dtd](#)

The DTD we built in this example

[createdb.bat](#)

An MS-DOS batch file that creates the DB2 database

[flights.txt](#)

File of comma-separated values that represent the sample database.

Summary

This paper has covered the basics of creating a DTD. The most important part of this task is understanding the structure of our source data and the data relationships we want our XML tags to convey. As mentioned earlier, the XML tags we've created add semantic meaning and let us process our data in much more flexible ways. These benefits will be more apparent as we continue to develop our XML application.

[Step 2: Generating XML from a Data Store](#)

[Step 3: Converting XML into HTML with the Document Object Model \(DOM\)](#)

Please send any comments or questions to:

Doug Tidwell

dtidwell@us.ibm.com